

**Б.Ф. Кирьянов, М.С. Токмачев**

# **МАТЕМАТИЧЕСКИЕ МОДЕЛИ В ЗДРАВООХРАНЕНИИ**

---

**МОНОГРАФИЯ**

ВЕЛИКИЙ НОВГОРОД

2009

УДК 51.7:312.6

ББК 22.19 и 5  
К34

## РЕЦЕНЗЕНТЫ

**В.Г. Дегтярёв**

доктор технических наук, профессор,  
академик Международной академии наук высшей школы,  
Заслуженный деятель науки РФ, Председатель Научно-методического  
совета по математике вузов Северо-Запада РФ

**Кафедра прикладной математики и информатики**

Санкт-Петербургского государственного архитектурно-строительного  
университета (заведующий кафедрой – доктор физико-математических  
наук **Б.Г. Вагер**)

**Научный консультант по проблемам здравоохранения:**

**В.А. Медик**

доктор медицинских наук, профессор, член-корреспондент РАМН, ди-  
ректор Новгородского научного центра Северо-Западного Отделения  
РАМН, заведующий кафедрой общественного здоровья, здравоохране-  
ния и общей гигиены Института медицинского образования Новгород-  
ского государственного университета им. Ярослава Мудрого

**Кириянов Б.Ф., Токмачёв М.С.**

К34 Математические модели в здравоохранении: учеб. пособие / Б.Ф.  
Кириянов, М.С. Токмачёв; НовГУ им. Ярослава Мудрого. – Великий  
Новгород, 2009. – 279 с.  
ISBN 978-589896-357-6

Предлагаются и исследуются модели показателей здоровья. Значительное внимание уделяется моделям интегрального показателя здоровья и прогнозированию показателей здоровья. Теоретический материал проиллюстрирован многочисленными примерами из области здравоохранения. Предлагаемые методы и алгоритмы могут быть использованы и для моделей другого назначения.

Излагаемый материал соответствует программе учебной дисциплины «Математические модели в здравоохранении», читаемой магистрантам направления «Прикладная математика и информатика» Новгородского государственного университета им. Ярослава Мудрого. Он полезен аспирантам и специалистам, занимающимся проблемами математического моделирования сложных систем.

УДК 51.7:312.6  
ББК 22.19 и 5

ISBN 978-589896-357-6

© Новгородский государственный  
университет, 2009  
© Б.Ф. Кириянов, М.С. Токмачев,

## ОГЛАВЛЕНИЕ

<b>Введение.</b> Системы и их моделирование. Общественное здоровье в системе качества жизни .....	6
<b>Глава 1.</b> Основные сведения о математических моделях и моделировании .....	12
1.1. Модель и моделирование. Математическое моделирование .....	12
1.2. Разновидности математических моделей .....	17
1.3. Основные этапы разработки математических моделей .....	20
1.4. Проведение исследований на модели и интерпретация полученных результатов .....	23
1.5. Факторный анализ на модели и его применение в здравоохранении .....	26
<b>Глава 2.</b> Основы теории вероятностей и математической статистики ..	31
2.1. Случайные события, вероятности и испытания .....	31
2.2. Случайные величины .....	42
2.2.1. Дискретные случайные величины .....	43
2.2.2. Непрерывные случайные величины .....	51
2.3. Системы случайных величин .....	56
2.3.1. Основные понятия и характеристики .....	56
2.3.2. Совместные распределения случайных величин с примерами из здравоохранения .....	60
2.3.3. Условное распределение. Понятие регрессии .....	63
2.4. Элементы математической статистики .....	65
2.4.1. Выборочный метод .....	65
2.4.2. Оценки параметров распределения. Доверительные интервалы .....	70
2.4.3. Проверка статистических гипотез .....	76
<b>Глава 3.</b> Моделирование показателей здоровья населения .....	80
3.1. Математические модели в здравоохранении и их характеристики	80
3.2. Здоровье населения и статистика его показателей. Базы данных ..	87
3.3. Анализ распределений показателей здоровья населения и показателей работы учреждений здравоохранения .....	94
3.4. Временные ряды в статистике здравоохранения и их характеристики .....	98
3.5. Моделирование показателей здоровья и их временных рядов .....	106
3.5.1. Модели трендов временных рядов .....	107
3.5.2. Моделирование случайной составляющей временных	

рядов показателей здоровья населения и показателей работы учреждений здравоохранения . . . . .	111
3.5.3. Моделирование случайных величин с колоколообразными распределениями . . . . .	114
<b>Глава 4.</b> Модели интегрального показателя здоровья населения на основе “стандартных” показателей здоровья . . . . .	125
4.1. Интегральная оценка здоровья населения. Интегральные показатели . . . . .	125
4.2. Однопараметрические модели . . . . .	129
4.3. Структура многопараметрической модели интегрального показателя общественного здоровья населения . . . . .	133
4.4. Линейные многопараметрические модели . . . . .	136
4.5. Нелинейные многопараметрические модели . . . . .	149
4.6. Чувствительность интегральных показателей здоровья населения . . . . .	154
4.7. Сравнение интегральных оценок здоровья населения регионов Российской Федерации . . . . .	157
<b>Глава 5.</b> Корреляция и регрессия. Модели зависимостей . . . . .	160
5.1. Типы зависимостей . . . . .	160
5.2. Выборочный коэффициент корреляции . . . . .	162
5.3. Проверка независимых признаков . . . . .	167
5.4. Проверка гипотезы о силе линейной связи двух признаков . . . . .	169
5.5. Выборочная регрессия . . . . .	171
5.6. Параметры выборочного уравнения регрессии при линейной зависимости . . . . .	176
5.7. Использование линейной регрессии в случае нелинейной зависимости . . . . .	179
5.8. Мера корреляционной связи. Выборочное корреляционное отношение . . . . .	182
5.9. Простейшие случаи нелинейной регрессии . . . . .	184
5.10. Методика построения модели множественной регрессии . . . . .	186
5.11. Примеры построения регрессионных моделей . . . . .	191
5.11.1. Соотношения параметров физического развития детей . . . . .	192
5.11.2. Вычисление площади поверхности тела человека . . . . .	199
5.11.3. Модели зависимости заболеваемости от степени загрязнения атмосферного воздуха . . . . .	205
5.11.4. Регрессионные модели заболеваемости и смертности населения . . . . .	212
<b>Глава 6.</b> Прогнозирование показателей здоровья населения на основе	

”стандартных“ параметров . . . . .	219
6.1. Классические методы и алгоритмы прогнозирования временных рядов . . . . .	219
6.2. Анализ точности прогнозирования показателей здоровья на основе полиномиальных моделей . . . . .	225
6.3. Прогнозирование показателей здоровья на основе “неполиномиальных” моделей . . . . .	237
6.4. Прогнозирование при “нетипичных” выбросах значений показателей здоровья . . . . .	239
6.5. Прогнозирование состояния здоровья населения Новгородской области . . . . .	246
<b>Глава 7. Моделирование на основе цепей Маркова . . . . .</b>	<b>249</b>
7.1. Основные понятия цепей Маркова . . . . .	249
7.2. Некоторые модели марковских процессов . . . . .	259
7.3. Исследование здоровья населения по статистическим данным . . . . .	267
7.3.1. Состояния здоровья. Классификация . . . . .	267
7.3.2. Формирование стохастических матриц и вычисление безусловных вероятностей . . . . .	273
7.3.3. Вычисление средней продолжительности жизни по группам. Сравнение групп . . . . .	279
7.3.4. Вычисление показателей продолжительности жизни фактического населения . . . . .	283
7.3.5. Показатели, характеризующие состояние здоровья . . . . .	286
Приложения . . . . .	291
Литература . . . . .	298

## **ВВЕДЕНИЕ. СИСТЕМЫ И ИХ МОДЕЛИРОВАНИЕ. ОБЩЕСТВЕННОЕ ЗДОРОВЬЕ В СИСТЕМЕ КАЧЕСТВА ЖИЗНИ**

Человек представляет собой сложную динамическую систему, компоненты которой взаимосвязаны и в ряде случаев реагируют на изменения внешней среды. Трудно назвать какую-либо другую систему из области естествознания, техники, экономики и т. д., превосходящую названную по своей сложности.

Любая система является совокупностью компонентов (объектов), объединённых определённым взаимодействием или некоторой взаимной зависимостью. В крупном плане все системы с точки зрения их природы можно разделить на физические, технические и административные. Следовательно, человек представляет собой физическую систему. Поскольку в этой системе постоянно происходят соответствующие процессы, то она относится к динамическим системам. Кроме динамических, существуют ещё статические системы, в которых на определённом отрезке времени нет изменений. Примером последних может служить система расположения учреждений здравоохранения в некотором городе в течение определённого интервала времени.

По мере развития человечества создавались всё более сложные системы и совершенствовались математические методы, позволяющие исследовать как системы создаваемые человеком, так и системы природного характера. Так было и во времена Архимеда, и во времена Ньютона, и во времена Королёва. Этот процесс наблюдается и в наши дни. Математические методы широко используются не только в управлении, технике, экономике, которые без их развития и применения были бы просто немыслимы, но и в биологии, медицине, метеорологии и др., т.е. в таких областях знаний, в которых полтора-два века тому назад применение математических методов представлялось невозможным и нецелесообразным.

20-й век следует считать родоначальником мощного метода исследования систем – математического моделирования, предполагающего создание математического описания, которое в определённой степени задаёт поведение исследуемых систем. В связи с бурным развитием вычислительной техники такой метод получил широкое применение. В ряде случаев он является единственно возможным методом исследования систем. Указанное математическое описание стали называть моделью рассматриваемой системы.

Не следует думать, что описания ещё неизученных систем предлагались только в 20-м веке. Так, ещё во 2-м веке до нашей эры Клавдий Птоломей предложил геоцентрическую систему (модель) мира, согласно которой центром вселенной является неподвижная Земля, а все остальные небесные тела вращаются вокруг неё. В 16-м веке нашей эры Николаем Коперником была создана гелиоцентрическая система (модель) мира, в которой Земля вращается вокруг Солнца и вокруг своей оси. При проведении в 19-20 веках экспериментальных проверок разрабатываемых лекарственных средств на подопытных животных по существу организм животного рассматривался как модель организма человека.

Начало активного применения математического моделирования в медицине и, в частности, в здравоохранении относится ко второй половине 20-го века. Безусловно, это, как указывалось, связано со значительным расширением возможностей ЭВМ. Однако, кроме того, выяснилось, что математические модели, достаточно адекватные соответствующим моделируемым системам, позволяют просто сравнивать показатели здоровья населения различных регионов, оптимизировать усилия и средства, направленные на повышение эффективности работы медицинских учреждений, проверять различные гипотезы относительно диагноза и развития болезни пациента, прогнозировать протекание болезни в зависимости от результатов обследования больного и методики лечения и т.д. Таким образом, во многих случаях математическое моделирование является мощным средством, с помощью которого можно выполнить значительный объём разнообразных и достаточно

сложных медицинских исследований, а врач может получить хорошего помощника.

Очевидно, в большинстве случаев врач вполне может справиться с решаемыми проблемами, не используя математические модели и моделирование на ЭВМ. Однако несомненно и то, что эти модели и реализующие их машинные программы являются помощниками врача, позволяющими ему более обоснованно принимать то или иное решение или давать соответствующую рекомендацию.

Большинство математических моделей в области здравоохранения используют аппарат теории вероятностей и математической статистики. Отметим, что статистические данные, связанные с медициной, начали собираться и обрабатываться ещё в 17-18-м веках. Так, в 1749-м году в соответствии с постановлением шведского парламента (риксдага) начался сбор текущих сведений о состоянии и движении населения, которые затем были использованы для составления таблиц смертности для обоих полов. Однако медицинская статистика как особая область отделилась от статистики населения позднее. В 1865 году в Германии появилось первое руководство по медицинской статистике. Его автор Fr. Oesterlen собрал рассеянный по всем европейским странам, часто скудный материал и дал ему единую разработку. В России в 1763 году Указом Сената была организована Медицинская коллегия, в которую присылались и обобщались описания редких болезней, научные работы русских врачей. По-видимому, зарождение медицинской статистики в России следует отнести к началу работы указанной коллегии.

В настоящее время отчётные статистические данные всеми региональными отделами здравоохранения через министерство здравоохранения и социального развития РФ направляются в ГОСКОМСТАТ РФ, который ежегодно публикует таблицы региональных статистических данных по различным видам заболеваний, рождаемости, смертности, обеспеченности учреждений здравоохранения врачами и медицинским оборудованием и т.п., а также формирует электронные базы данных. Кроме того, в региональные базы

данных могут входить и более подробные данные, касающиеся соответствующего региона.

Указанные базы статистических показателей, а также основы теории вероятностей, математической статистики, программирования и математического моделирования на ЭВМ являются основой для построения соответствующих математических моделей и проведения исследований на их основе.

Значительная активизация исследований по применению математического моделирования в различных медицинских системах произошла в последнем десятилетии 20-го и в начале 21-го века. Так, в этот период российскими учёными были разработаны и применены в медицинской практике математические модели адаптации организма и патологических процессов [9, 112 ÷ 114], построены оптимизационная модель, позволяющая осуществлять выбор начального плана лечения желтухи, и классификационно-диагностические модели, на основе которых возможно проведение своевременной дооперационной диагностики ряда хирургических заболеваний, осложнённых желтухой [33, 34], а также ряд других моделей, связанных с медициной [31, 32, 37, 39, 61, 98]. Из зарубежных разработок следует отметить модели интегрального показателя оценки здоровья населения [142, 145, 158]. Поиску такого показателя посвящены и работы некоторых российских учёных [66, 75, 80, 116].

В последние годы в медицине начинает применяться и математические модели с анимацией, т.е. с изменяющимися изображениями на дисплее (аналогично мультфильмам в телевидении). Это позволяет наблюдать в ускоренном масштабе времени прогнозируемое развитие ряда болезней, динамику работы сердца и т.д. Причём на дисплеях соответствующие изображения могут представляться в трёхмерном пространстве.

Из разработок математических моделей в здравоохранении центральное место занимают модели интегрального оценивания здоровья населения. Всемирная организация здравоохранения (ВОЗ) ещё в 1971 году сформулировала требования к таким показателям. На необходимость разработки инте-

гральных показателей здоровья населения указывали и известные отечественные учёные Н.П. Амосов, Ю.П. Лисицын и другие. Учитывая, что на здоровье населения, на качество его жизни существенное влияние оказывают социально-экономические условия жизни и образ жизни населения, они ввели термин «общественное здоровье населения». Хотя детально понятие общественного здоровья населения ещё не определено, интуитивно ясно, что оно в значительной степени определяет качество жизни населения. Поэтому при интегральном оценивании здоровья населения необходимо учитывать влияние на него таких факторов как социально-экономические условия и образ жизни населения.

В течение нескольких лет совместно с научным консультантом по проблемам здравоохранения, членом-корреспондентом РАМН, доктором медицинских наук В.А.Медиком авторы работали над проблемами разработки научно-обоснованных, простых и удобных для практического применения математических моделей интегрального показателя оценки общественного здоровья населения, включая проблемы моделирования и прогнозирования показателей здоровья. Разработка указанных показателей, представляющих собой функцию в общем случае от нескольких статистических показателей здоровья, позволяет, в частности, осуществлять обоснованное, интегрированное сравнение здоровья населения различных регионов, облегчая управление системой здравоохранения как на региональном, так и на федеральном уровне. Разработанные авторами модели интегрального показателя здоровья населения используются в системе здравоохранения Новгородского региона. Полученные результаты отражены более, чем в 40 публикациях.

Настоящее учебное пособие предназначено для магистрантов направления «Прикладная математика и информатика» и аспирантов, специальности которых близки к этому направлению. Оно содержит основные сведения о математическом моделировании, описывает ряд разработанных моделей показателей здоровья и характеризует их возможности. В частности, рассматриваются вопросы построения таких моделей на основе цепей Маркова,

а также проблемы прогнозирования значений интегральных и статистических показателей здоровья. Значительная часть материала является оригинальной, часть полученных результатов публикуется впервые.

В связи с указанным настоящее издание может рассматриваться и как монография, рассчитанная для широкого круга читателей, интересующихся вопросами практического приложения математического моделирования, в частности для интегральной оценки качества систем различного профиля, на примере системы здравоохранения. Она полезна студентам математических и вычислительных направлений, а также работникам системы здравоохранения с вузовской математической подготовкой для медицинских специальностей. С отдельными разделами учебного пособия можно знакомиться и вообще не владея указанной математической базой. Вместе с тем авторы посчитали целесообразным включить в него краткое изложение материала по основам теории вероятностей и математической статистики, полезное для тех читателей, которые недостаточно знакомы, но желают ознакомиться с указанным материалом (главы 2 и 6 с примерами из области здравоохранения).

Авторы выражают благодарность научному консультанту, директору Новгородского научного центра Северо-Западного Отделения РАМН, члену-корреспонденту РАМН, доктору медицинских наук, профессору В.А. Медичу, а также рецензентам Председателю научно-методического совета вузов Северо-Запада, Заслуженному деятелю науки РФ, академику международной академии наук высшей школы, доктору технических наук, профессору В.Г. Дегтярёву, кафедре прикладной математики и информатики Санкт-Петербургского государственного архитектурно-строительного университета и лично заведующему этой кафедры доктору физико-математических наук, профессору Б.Г. Вагеру за рекомендации, способствующие улучшению содержания учебного пособия и изложения его материала.

# ГЛАВА 1. ОСНОВНЫЕ СВЕДЕНИЯ О МАТЕМАТИЧЕСКИХ МОДЕЛЯХ И МОДЕЛИРОВАНИИ

## 1.1 Модель и моделирование. Математическое моделирование

Слово «модель» происходит от латинского слова «modulus» – образец, мера, мерило. Модель некоторого объекта или явления, называемого оригиналом, представляет собой такой объект, явление или математическое описание, функционирование которого оказывается в достаточной степени аналогичным функционированию оригинала. Таким образом, свойства модели должны быть аналогичны соответствующим свойствам оригинала. При этом в качестве оригинала не обязательно должна использоваться какая-либо сложная система. Оригиналами могут быть и отдельные элементы такой системы.

Исследование функционирования модели и приписывание оригиналу закономерностей, выявленных при исследовании модели, называется моделированием. Согласно [17] моделирование – это исследование объектов познания на их моделях. Однако возникает вопрос: а зачем вообще нужно моделирование? Почему нельзя исследовать непосредственно оригинал? Дело в том, что это не всегда возможно или допустимо, а во многих случаях и очень дорого. Так в приведённом во введении примере исследования воздействия создаваемых лекарственных средств на организм человека характер этого воздействия вначале изучается на модели, в качестве которой обычно берётся подопытное животное. Проведение же указанных исследований непосредственно на человеке может быть недопустимым. В качестве примера недопустимо дорогого проектирования оригиналов может служить создание нескольких опытных экземпляров кораблей проектируемого класса с различными профилями их корпуса, влияющего на скоростные характеристики корабля. Профиль корабля гораздо проще и дешевле подобрать с помощью модели.

Особым случаем моделирования является подбор на модели параметров для проектируемого или неизвестного оригинала. В этом случае руководствуются общими требованиями к оригиналу. И только в результате исследования различных вариантов модели выбирается наиболее подходящий проект оригинала или наиболее подходящее его математическое описание. Примером такого моделирования является разработка моделей интегрального показателя оценки здоровья населения [66, 83, 87, 90, 116, 142, 152, 158 и др.], т.е. создание и обоснование математического описания для расчёта указанного показателя, который достаточно полно характеризовал бы состояние здоровья населения.

Существуют два основных метода моделирования, а соответственно и два основных класса моделей:

1. Физическое моделирование.
2. Математическое моделирование.

При физическом моделировании используется модель одной с оригиналом физической природы. В этом случае свойства оригинала обычно воспроизводятся полнее, чем при математическом моделировании. Кроме того, имеется возможность использовать для исследования ту же регистрирующую аппаратуру, которая используется для оригинала. Вместе с тем создание и перенастройка физических моделей часто бывают трудоёмкими и дорогостоящими. Типичными примерами физического моделирования являются приведённые выше примеры моделирования профилей кораблей и испытания разрабатываемых лекарств на подопытных животных.

При математическом моделировании используются модели отличной от оригинала физической природы. Однако процессы или явления, проходящие в модели, должны описываться теми же зависимостями, что и основные процессы или явления в оригинале. Так, упомянутая выше задача подбора профилей кораблей может решаться на модели из электропроводной бумаги, в которой вырезан исследуемый вариант профиля корабля. В данном случае используется то, что электрический ток обтекает вырезанный профиль со-

гласно тем же математическим уравнениям, которым подчиняется обтекание водой движущегося корабля.

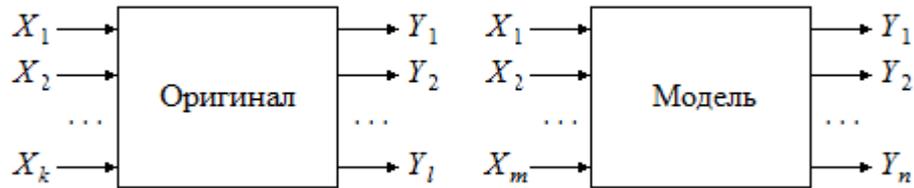
В настоящее время математическое моделирование практически всегда проводится на ЭВМ. В этом случае моделью является моделирующая программа, воспроизводящая исследуемый процесс или явление в соответствии с его математическим описанием. Такое моделирование имеет очень широкие возможности. В отличие от других видов моделирования трудно назвать какой-либо оригинал, который нельзя было бы промоделировать на ЭВМ. В ряде случаев математическое моделирование на ЭВМ является единственно возможным средством исследования систем. Параметры моделирующих программ легко могут быть изменены, а сами программы могут храниться неограниченно долго и практически не занимая места. Обычно математическое моделирование на ЭВМ обходится существенно дешевле других видов моделирования, отличается простотой изменения параметров модели и позволяет получать результаты с более высокой точностью.

Ввиду указанных достоинств математического моделирования на ЭВМ именно оно почти всегда и применяется в последние годы для исследования соответствующих проблем в медицине и, в частности, в здравоохранении [9, 45, 61, 66, 87, 90, 93, 112 – 114, 158 и т.д.]. При этом ежегодный прирост публикаций по этим проблемам имеет устойчивую тенденцию возрастания.

Рассмотрим упрощенное математическое представление модели, необходимое для пояснения важных понятий «*функциональная полнота*» и «*адекватность*» модели оригиналу. Однако при этом не будем пользоваться используемым в математической литературе определением понятия «*математическая модель*», для понимания которого необходимо иметь соответствующую математическую подготовку. Интересующихся можно адресовать, например, к [108].

Пусть некоторый оригинал (рис. 1.1) имеет  $k$  входов, на которые поступают соответствующие входные величины (независимые переменные, ар-

гументы)  $X_i$ . Множество  $\mathbf{X}$  этих величин является упорядоченным в том смысле, что каждая  $i$ -я величина поступает именно на  $i$ -й вход оригинала и



**Рис. 1.1.** Абстрактная система оригинал – модель

поэтому в обозначении множества  $\mathbf{X} = (X_1, X_2, \dots, X_k)$  указывается на  $i$ -м месте. Такие множества элементов, представляемых в виде одной строки, называются векторами. Поэтому можно считать, что на оригинал поступает вектор  $\mathbf{X}$  или входной векторный процесс  $\mathbf{X}(t)$ , где  $t$  – непрерывное или дискретное время. При моделировании на ЭВМ время и координаты входного и выходного векторов изменяются дискретно, и поэтому реализуемые модели называются дискретными.

Аналогично выходы оригинала задаются множеством  $\mathbf{Y} = (Y_1, Y_2, \dots, Y_l)$ , т.е. вектором  $\mathbf{Y}$ . Каждый его элемент  $Y_i$  является функцией вектора  $\mathbf{X}$  и внутреннего состояния оригинала. Последнее в общем случае может изменяться под воздействием как входных, так и выходных величин. Если на вход модели поступает последовательность векторов  $\mathbf{X}$ , а с её выхода снимается последовательность векторов  $\mathbf{Y}$ , то эти последовательности обычно называют векторными процессами – соответственно  $\mathbf{X}(t)$  и  $\mathbf{Y}(t)$ . Процессы  $\mathbf{X}(t)$  и  $\mathbf{Y}(t)$  могут быть детерминированными или случайными.

Приведённая на рис. 1.1 модель может как и оригинал иметь  $k$  входов и  $l$  выходов, формируемых аналогично оригиналу. Такая модель называется функционально полной. Если при этом имеет место ещё и полное совпадение выходных процессов модели и её оригинала, соответствующих любому входному процессу, то говорят, что модель адекватна оригиналу. Однако во многих случаях при моделировании нет необходимости воспроизводить все возможности оригинала. Так при моделировании механизмов обучения в

нейронных сетях мозжечка нет необходимости строить более полную модель системы организации и функционирования памяти человека [98]. Поэтому в данном случае  $m < k$  и  $n < l$ . При этом можно рассматривать степень адекватности такой модели по отношению к оригиналу «система обучения в нейронных сетях мозжечка», которая явно выше, чем по отношению к оригиналу «система организации и функционирования памяти человека». Однако заметим, что понятие «степень адекватности», а также используемые в литературе такие характеристики этой степени как «высокая степень адекватности», «достаточно адекватная» или «относительно адекватная», количественно не определены. Поэтому с математической точки зрения они несостоятельны. Использование же таких понятий и характеристик объясняется тем, что нередко адекватность модели оригиналу просто невозможно проверить. Поэтому при разработке и исследовании математических моделей в области здравоохранения авторы вынуждены пользоваться указанными характеристиками с учётом рекомендаций экспертов-медиков.

В зависимости от числа независимых входных величин, являющимися переменными параметрами, модели нередко называют однопараметрическими, двухпараметрическими, многопараметрическими. Кроме того, выходные величины  $Y_i$  могут зависеть и от постоянных параметров, например, от коэффициентов, на которые умножаются  $X_i$ .

В моделях интегральных показателей значение  $m$  равно числу учитываемых статистических показателей, а значение  $n$  обычно равно единице. В частности, в предложенных авторами моделях интегрального показателя здоровья значение  $m$  в несколько раз меньше числа стандартных статистических показателей здоровья ( $k$ ), определяемых регионами и обобщаемых ГОСКОМСТАТОм России для страны в целом. Общее число таких показателей более 100.

## 1.2. Разновидности математических моделей

Математические модели (или математическое моделирование) можно классифицировать по различным признакам [71, 78 и др.]. По виду математического описания оригинала можно выделить два класса математических моделей: *аналитические* и *имитационные*. В первом случае свойства или поведение оригинала описываются математическими зависимостями, которые и исследуются. Каждая такая зависимость определяет значение соответствующей выходной величины (а не процесса её изменения). Если все указанные зависимости – линейные, то модель называется *линейной*. В простейшем случае, когда в математическом описании используются только операции сложения и вычитания, линейная модель называется *аддитивной*. При наличии в этом описании нелинейных операций модель называется *нелинейной*. Примером нелинейной аналитической модели является аппроксимация статистического ряда распределения некоторого показателя здоровья выбранной функцией, удовлетворяющей выбранному критерию согласия (соответствия) статистическому распределению.

Однако аналитическое описание для всей модели не всегда удаётся получить ввиду сложности системы-оригинала, состоящей, например, из нескольких подсистем, поведение которых зависит от случайных внешних факторов и является неопределённым. В этом случае основой математического описания модели являются алгоритмы, имитирующие соответствующие процессы в оригинале. Такая модель называется имитационной, а исследование её – имитационным моделированием [78, 135], например – моделирование процессов адаптации организма к различным условиям [113], патологических процессов [114], структуры прошлого опыта в памяти человека [98], моделирование процесса роста бляшек в кровеносном тракте [61] и др.

Если  $Y(t)$  зависит только от  $X(t)$  и не зависит от входных векторов, поступивших на модель ранее момента  $t$ , то модель называется статической. В противном случае её называют динамической. В статических моделях все

выходные величины при постоянстве входных величин не изменяются во времени. Типичным примером таких моделей являются модели интегрального показателя здоровья, который изменяется только при изменении учитываемых им статистических показателей. Основой математического описания статических моделей являются алгебраические уравнения. В динамических же моделях выходные величины могут изменяться даже при постоянстве входных величин. Основой их математического описания являются дифференциальные уравнения, описывающие соответствующие процессы [113, 98]. Модели процессов, регистрируемых с дискретным шагом, описываются разностными уравнениями [78, 112].

В зависимости от отсутствия или наличия случайностей в алгоритмах моделирования, т.е. в математическом предписании, определяющем вычислительный процесс получения конечных результатов, модели делятся на *детерминированные* и *стохастические* (реже – *вероятностные, статистические*). В детерминированных моделях используются только детерминированные алгоритмы, однозначно выполняющие необходимые преобразования. При этом входные величины могут принимать и случайные значения, которые однозначно обрабатываются моделью. Примером таких моделей являются модели интегральных показателей здоровья. Стохастические модели отличаются наличием в их алгоритмах преобразований со случайными результатами (стохастических алгоритмов). Обычно такие преобразования осуществляются наряду с детерминированными преобразованиями. Стохастические модели используются для исследования процессов, характеризующихся случайностью каких-либо параметров в соответствующих оригиналах (ежедневный прирост числа заболевших в моделях развития эпидемии, возраст заболевшего ишемической болезнью индивидуума, при котором с начинает расти первая бляшка в левой или в правой сосудистой ветви кровеносного тракта, параметры процесса адаптации организма человека к патологическому процессу и др.).

Для реализации случайностей в программном обеспечении ЭВМ имеется процедура *random*, имитирующая случайные величины, равномерно распределённые в интервале от 0 до 1 или принимающие равновероятные целые значения  $0, 1, 2, \dots, K$ . На основе таких случайных величин можно получать программным путём случайные величины с произвольными законами распределения, а также последовательности таких величин, имитирующих соответствующие случайные процессы.

Особый класс математических моделей представляют *нечёткие* модели, базирующиеся на аппарате нечётких множеств и нечёткой логики. Такие модели используются для исследования систем с расплывчатыми условиями или алгоритмами [14]. В частности такие модели применяются для исследования нечётких нейронных сетей [6].

В зависимости от разновидности решаемой задачи в литературе встречаются и соответствующие названия моделей или вида анализа, отражающие эти разновидности. Укажем некоторые из них, наиболее часто встречающиеся, с соответствующими примерами:

- Приближение функций (представление статистического ряда распределения “подходящей” функцией).
- Оптимизационная модель (выбор расположения учреждений здравоохранения на территории города или области по максимуму или минимуму соответствующего критерия).
- Дисперсионный анализ (анализ оценок дисперсии или среднего квадратического отклонения статистического распределения);
- Разностные модели (модели численного дифференцирования, численного интегрирования, численного решения дифференциальных уравнений, используемые в машинных программах вместо соответствующих оригиналов – методов дифференцирования, интегрирования и решения дифференциальных уравнений).
- Модели прогнозирования (прогнозирование значений интегральных показателей здоровья, протекания некоторой болезни при отсутствии

или наличия лечения, влияния климатических и социально-экономических условий жизни на здоровье населения и др.).

### **1.3. Основные этапы разработки математических моделей**

Перед разработкой модели, как и при проведении любого научного исследования, следует провести поиск литературных источников по интересующей проблеме и оценить целесообразность разработки модели. Если окажется, что все интересующие результаты уже известны, то построение и использование модели может иметь смысл, например, в учебно-демонстрационных целях.

В крупном плане можно выделить следующие 3 этапа разработки математической модели на ЭВМ (несколько другой вариант названий и сущности этих этапов предложен в [115] ):

- Разработка математического описания модели (уточнение задачи моделирования, определение требований к входным и выходным переменным, построение структурной схемы модели, выбор её параметров, выбор методов получения конечных и промежуточных результатов). Результатом этого этапа является формальное математическое описание всех блоков модели. При этом как бы исчезнет физическая сторона исследуемой задачи, т.е. для лиц, не знакомых с решаемой задачей, по математическому описанию модели крайне затруднительно уяснить физическую сущность задачи.
- Разработка машинной программы моделирования (включая алгоритмизацию математического описания, разработку программы в соответствующей среде программирования и проверку правильности её работы – например, по моделированию характерных частных вариантов входных величин, для которых известны или могут быть просто рассчитаны значения выходных величин, по выполнению критериев соответствия

законов распределения моделируемых случайных величин принятым в модели законам и т.д.).

- Испытания и корректировка модели (моделирование исследуемой системы при характерных частных случаях входных величин или параметров модели, оценка достоверности полученных для этих случаев результатов и в случае необходимости – корректировка соответствующих параметров модели).

Разработка концептуальной модели неизбежно связана с выбором степени её функциональной полноты и с попыткой оценить её адекватность оригиналу. При этом стремление обеспечить функциональную полноту модели обычно оказывается нецелесообразным во-первых потому, что модель может стать неоправданно сложной, а во-вторых в связи со слабым во многих случаях влиянием отдельных входных величин на выходные величины модели [78, 134, 139]. Более того, в дальнейшем (гл. 4) будет показано, что при неограниченном увеличении числа входных величин модели может оказаться необходимым неограниченно увеличивать число разрядов после запятой в числовом представлении выходных величин, изменяющихся в конечном интервале. Мероприятия по улучшению степени адекватности модели оригиналу могут проводиться и при корректировке модели. Проведение таких мероприятий иногда называют валидацией модели [78, 105]. К таким мероприятиям относится, например, выяснение необходимой точности задания входных величин и параметров модели для того, чтобы получить выходные величины с желаемой точностью.

Возникает вопрос: специалистами какого профиля должны разрабатываться математические модели? Что касается разработки машинной программы моделирования, то здесь всё ясно – это дело математиков-программистов. А вот 1-й и 3-й основные этапы разработки модели обычно могут выполняться только совместно специалистами в области, к которой относится оригинал, и математиками-прикладниками (специалистами по моделированию на ЭВМ и по теории вероятностей и математической статистике).

При разработке математического описания модели обычно приходится сталкиваться с необходимостью задания в модели значений параметров, которые точно не известны, или даже принятия математических выражений, которые описывали бы протекание недостаточно изученных процессов, задавали бы закон распределения некоторой случайной величины, по которой недостаточно или практически нет статистических данных и т.д. При принятии решений по этим вопросам имеется определённый произвол, который желательно свести к минимуму, заменяя его обоснованными решениями.

Примером указанных выше параметров в моделях интегрального показателя здоровья населения являются коэффициенты, учитывающие исчерпанную (истинную) заболеваемость различных групп населения, а не только заболеваемость этих групп населения, зафиксированную по обращениям в амбулаторно-поликлинические учреждения (включая случаи госпитализации). Случаи заболеваемости, по которым не было обращений в амбулаторно-поликлинические учреждения, выявляются при медицинских осмотрах населения. Такие осмотры являются уникальным источником информации для углублённого изучения заболеваемости населения особенно потому, что они позволяют выявить те хронические заболевания и состояния, которые не явились причиной обращения в указанные учреждения. Проведённые в 2000-м и в 2005-м годах в Новгородской области медицинские осмотры населения показали [68, 95], что учёт исчерпанной заболеваемости, обеспечивающий объективность в учёте числа болевших индивидуумов, может существенно изменить общий показатель заболеваемости. Так, общее число учтённых случаев всех заболеваний населения Великого Новгорода в 2000 году возросло в 2,5 раза, а хронических заболеваний – в 3,25 раза. В 2005-м году общая исчерпанная заболеваемость населения Новгородской области оказалась в 1,898 раза выше общей заболеваемости по обращениям. Однако медицинские осмотры проводятся не во всех регионах, особенно для сельского населения. Если такие осмотры будут проводиться выборочно, то их следует проводить с учётом рекомендаций экспертов.

Аналогично с участием эксперты устанавливаются и значения постоянных коэффициентов в дифференциальных уравнениях, моделирующих процессы адаптации организма к изменяющимся условиям окружающей среды и к патологическим процессам [113], параметры законов распределения статистических показателей здоровья и т.д.

Обычно при разработке модели предусматривается возможность варьирования параметрами её вычислительных алгоритмов. Это обеспечивает определённую универсальность модели, так как с помощью такой модели можно исследовать оригинал при различных его характеристиках. Следовательно, в данном случае задача экспертов заключается в установлении соответствий характеристик оригинала и указанных параметров. В частности, при моделировании процесса протекания какой-либо болезни конкретного индивидуума перечнем значений соответствующих характеристик оригинала является карта обследования больного.

#### **1.4. Проведение исследований на модели и интерпретация полученных результатов**

Приступая к исследованиям с помощью моделирования, целесообразно предварительно хотя бы мысленно наметить план их проведения, т.е. спланировать эксперимент [146]. Одной из целей эксперимента может быть, например, уточнение требований к точности задания входных величин и параметров модели, так как эту точность далеко не всегда можно определить на этапе разработки математического описания модели. То же самое может быть и при выборе числа испытаний при анализе характеристик случайных процессов или явлений.

В процессе исследований может также оказаться, что запрограммированная форма представления получаемых результатов не вполне удобна. В таких случаях приходится корректировать параметры таблиц, масштабы графического представления величин и т.д. Следовательно, приходится вносить изменения и в машинную программу.

Важное значение при исследовании моделей имеет интерпретация получаемых результатов [20, 125]. Слово «интерпретация» происходит от латинского «interpretatio» - разъяснение. Применительно к моделированию оно означает осмысление, разъяснение, толкование результатов моделирования. Основная задача здесь заключается в решении вопроса: могут ли полученные на модели результаты быть получены и с помощью оригинала и насколько они типичны для него? Следовательно, решается вопрос о приемлемости полученных результатов. При этом, как уже указывалось, не всегда на этот вопрос удаётся получить исчерпывающий ответ.

Чаще всего необходимость в интерпретации результатов моделирования возникает при применении стохастических моделей. В таких моделях приходится сталкиваться с определением оценок характеристик получаемых в результате моделирования случайных величин, с проверкой гипотез о соответствии законов распределения указанных величин каким-либо классическим законам распределения, с определением оценок параметров этих законов распределения и т.д. Во многих случаях перечисленные проблемы могут быть решены аналитически с помощью аппарата математической статистики. С методами их решений и с соответствующими примерами можно ознакомиться, например, по публикациям [29, 89, 94, 95, 99, 107] ]. Однако не всегда при оценке приемлемости или точности получаемых путём моделирования результатов можно обойтись без решения с помощью моделирующей программы дополнительных задач. Поясним это на характерном примере.

Пусть необходимо промоделировать и исследовать последовательность  $n$  значений некоторой случайной величины  $X$ , получаемой путём нелинейных преобразований случайных величин  $X_1, X_2, \dots, X_{10}$ . В таком случае, если даже распределения случайных величин  $X_1, X_2, \dots, X_{10}$  известны, определить аналитически не только распределение случайной величины  $X$ , но и его характеристики обычно не удаётся. Однако число  $n$  моделируемых реализаций  $X$  необходимо выбрать так, чтобы получаемая в результате моделирования оценка  $\bar{X}$  значения математического ожидания величины  $X$ , называемая в

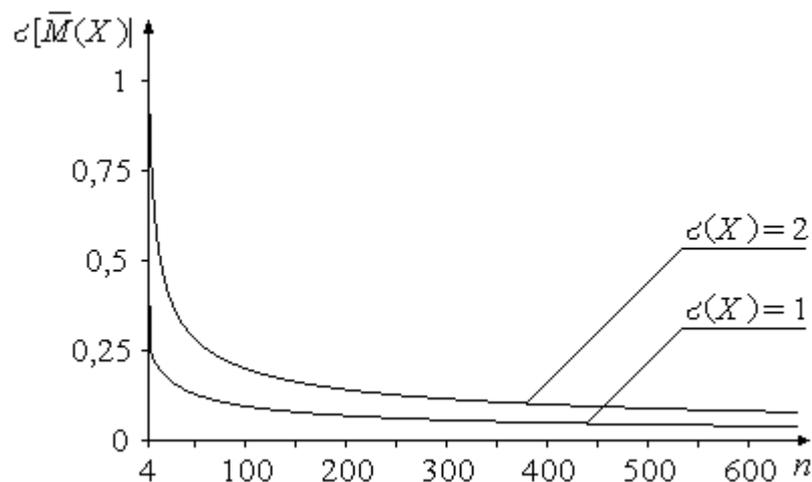
дальнейшем выборочным значением математического ожидания этой величины [29, 91, 95, 107], была бы достаточно близкой к неизвестному значению  $M(X)$ .

Если бы было известно значение среднего квадратического отклонения  $\sigma(X)$ , то, задав в качестве меры близости  $\bar{X}$  к  $M(X)$  допустимое значение  $\sigma(\bar{X})$  среднего квадратического отклонения  $\bar{X}$  от  $M(X)$ , т.е. к точности получения при моделировании значения  $\bar{X}$ , выбор необходимого значения  $n$  можно было бы сделать исходя из выражения [26, 91]

$$\sigma^2(\bar{X}) = \frac{\sigma^2(X)}{n}. \quad (1.1)$$

Действительно, потребовав выполнения условия  $\sigma(\bar{X}) \leq \varepsilon$ , где  $\varepsilon$  - требуемая точность (максимально допустимое значение  $\sigma(\bar{X})$ ), получим:

$$n \geq \frac{\sigma^2(X)}{\varepsilon^2}. \quad (1.2)$$



**Рис. 1.2.** Графики зависимости разброса  $\bar{X}$  от числа реализаций случайной величины  $X$  с  $\sigma(X) = 1$  и с  $\sigma(X) = 2$

Рис. 1.2 иллюстрирует графики зависимости (1.1) для случаев  $\sigma(X) = 1$  и  $\sigma(X) = 2$ . Не сложно заключить, что за увеличение точности определения  $\bar{M}(X)$  приходится расплачиваться увеличением времени моделирования.

Однако поскольку в рассматриваемом примере значение  $\sigma(X)$  не известно, то в выражениях (1.1) и (1.2) вместо  $\sigma^2(X)$  используют значение её оценки  $\overline{\sigma^2(X)}$ , которое определяется на каждом цикле моделирующей программы с получением очередного значения  $X$ , т.е. при каждом увеличении номера шага  $n$  на единицу. Так как  $\sigma^2(X)_n = M(X^2)_n - [M(X)_n]^2$ , то для указанной цели на каждом шаге работы модели вначале вычисляются значения оценок для  $M(X)_n$  и  $M(X^2)_n$ :

$$\overline{M(X)}_n = \frac{n-1}{n} \overline{M(X)}_{n-1} + \frac{X(n)}{n}, \quad \overline{M(X^2)}_n = \frac{n-1}{n} \overline{M(X^2)}_{n-1} + \frac{X^2(n)}{n}.$$

Затем находится и оценка  $\overline{\sigma^2(X)}_n = \overline{M(X^2)}_n - [\overline{M(X)}_n]^2$

Теперь при выполнении условия (1.2) можно было бы и прекратить моделирование очередных  $X$ , ограничившись полученным значением  $n$ . Но для «надёжного выполнения» этого условия лучше предусмотреть, например, его десятикратное выполнение подряд.

Приведённая методика обеспечения необходимой достоверности результатов моделирования использована, например, в математических моделях гидроэнергетической системы [16] и развития ишемической болезни сердца [61].

### 1.5. Факторный анализ на модели и его применение в здравоохранении

Каждая выходная величина  $Y_i$  модели (рис. 1.1) может быть функцией не только входных величин  $X_j$ , но и ряда параметров  $V_1, V_2, \dots, V_s$  модели. Во второй половине 20-го века входные величины модели и её параметры, влияющие на значения выходных величин, нередко стали называть *факторами* или *входными факторами*, а выходные величины – *выходными факторами* или *откликами* [78, 129]. Соответствующие задачи, сущность которых поясняется ниже, получили название задач факторного анализа. В зависимо-

сти от числа факторов, учитываемых моделью, указанные задачи называют задачами *однофакторного, двухфакторного и многофакторного* анализа.

Факторы могут быть либо *количественными*, либо *качественными*. Количественные факторы, как правило, предполагают численные значения, тогда как качественные факторы обычно представляют собой структурные допущения, принятые при построении модели, и не могут измеряться количественно

Если на модели исследуется влияние на выходные величины какой-либо одной входной величин или значений одного из параметров модели, то исследуемая задача является задачей однофакторного анализа. В таких задачах указанную входную величину или указанный параметр модели называют *факторами*. Цель такого анализа состоит в исследовании влияния фактора на отклик. Типичными примерами задач однофакторного анализа являются исследование эффективности разных лекарств одинакового назначения, анализ влияния алгоритмов модели определения интегрального показателя оценки здоровья населения на характеристики этого показателя и др. В этих задачах факторами являются «лекарство» и «алгоритм», а откликом – значение количественной характеристики или значения таких характеристик выходной величины (например, среднее время выздоровления больных, линейность и интервал изменения интегрального показателя).

Если указанных факторов два и отклик является функцией двух этих факторов, то соответствующая задача называется двухфакторным анализом. Примером двухфакторного анализа может служить приведённая выше задача исследования эффективности разных лекарств одинакового назначения, дополненная вторым фактором – временем хранения лекарства с момента его изготовления. Отметим, что в случае, когда имеются два фактора ( $A$  и  $B$ ) и два отклика ( $Y_1$  и  $Y_2$ ), но  $Y_1$  зависит только от  $A$ , а  $Y_2$  – только от  $B$ , то задача распадается на две задачи однофакторного анализа. Так, при анализе распределения смертности населения по переменной «возраст» для мужчин и женщин фактически имеют место две задачи однофакторного анализа, фактора-

ми в которых являются «возраст мужчин» или «возраст женщин» (в данном случае на одной модели могут последовательно исследоваться две указанные задачи).

В случае, когда отклик модели является функцией нескольких факторов (более двух), исследуемая на модели задача считается задачей многофакторного анализа. Такие задачи могут встречаться при использовании моделей интегрального показателя оценки здоровья населения.

Если модель имеет  $m$  входов и  $n$  выходов для скалярных величин (рис. 1.1), то решаемые на ней задачи могут относиться к задачам многофакторного, двухфакторного и однофакторного анализа скалярных факторов. Так, если отклик  $Y_3$  является функцией двух факторов – входных величин  $X_2$  и  $X_5$  или какого-либо параметра модели и одной из указанных входных величин, то задача анализа зависимости  $Y_3$  от этих факторов представляет собой задачу двухфакторного анализа.

При описании работы модели с помощью матричного аппарата можно считать, что модель имеет один входной вектор размерности  $m$ , один выходной вектор размерности  $n$  и один вектор параметров модели, размерность которого равна числу параметров. На такой модели будут решаться задачи только однофакторного или двухфакторного анализа, но для векторных факторов. Очевидно, при этом функционирование модели результаты её исследования не изменятся.

Число значений, которое может принимать тот или иной скалярный фактор, обычно бывает очень большим или даже бесконечным. Поэтому для удобства исследований и представления факторных величин модели интервалы изменения факторов часто делят на участки (группы, классы, разряды, уровни, ранги), т.е. множество возможных значений каждого фактора разбивают на соответствующие подмножества. При этом результирующие данные также разбиваются на соответствующие группы, число которых при однофакторном анализе равно числу групп значений фактора, а при двух и более факторах – произведению числа групп значений каждого фактора.

Табл. 1.1 и 1.2 иллюстрируют возможное представление результатов исследования распределения некоторых случайных величин  $Y$  и  $Z$ , зависящих соответственно от двух ( $A, B$ ) и трёх ( $A, B, C$ ) факторов, на основе репрезентативных выборок их значений. Значения каждого фактора разбиты на три группы. Объём выборки значений величины  $Y$  равен 1000, они распределены на 9 групп. Объём выборки значений величины  $Z$  равен 3000, они распределены на 27 групп.

Таблица 1.1.

Число значений  $y_{ij}$ 

$A \setminus B$	$B_1$	$B_2$	$B_3$
$A_1$	94	105	101
$A_2$	128	129	98
$A_3$	120	106	119

Таблица 1.2. Число значений  $z_{ijl}$ 

$C$	$C_1$			$C_2$			$C_3$		
$A \setminus B$	$B_1$	$B_2$	$B_3$	$B_1$	$B_2$	$B_3$	$B_1$	$B_2$	$B_3$
$A_1$	118	131	100	132	107	96	121	92	109
$A_2$	103	91	93	117	122	104	138	130	105
$A_3$	111	93	108	115	106	97	104	132	115

С увеличением числа факторов и их групп количество групп выходной величины резко возрастает. Так, если в трёхфакторной модели интервалы изменения каждого фактора разбить на 10 групп, то число групп выходной величины окажется равным 1000. Пользоваться такой её группировкой крайне сложно. При дальнейшем увеличении числа групп выходной величины модель становится неработоспособной [78]. Анализ причин её неработоспособности дан в [64] и будет приведён в гл. 4.

В заключение отметим, что в качестве количественных факторов в математических моделях в области здравоохранения используются такие как возраст индивидуума, значения коэффициентов исчерпанной заболеваемости, продолжительность лечения и др., а в качестве качественных – степень возможности повседневной деятельности индивидуума или степень его нетрудоспособности, обеспеченность лечебно-профилактических учреждений врачебным персоналом, распространённость исчерпанной заболеваемости и т.д. При математическом моделировании степень влияния качественных фак-

торов на основе экспертных оценок обычно задаётся с помощью соответствующих коэффициентов, значения которых и делят на группы. В результате качественные факторы представляются в виде количественных.

## ГЛАВА 2. ОСНОВЫ ТЕОРИИ ВЕРОЯТНОСТЕЙ И МАТЕМАТИЧЕСКОЙ СТАТИСТИКИ

### 2.1 Случайные события, вероятности и испытания

**Случайность и здравоохранение.** Реальный мир представляет огромное многообразие различных явлений, процессов, событий, взаимосвязанных между собой бесчисленными нитями причинно-следственных зависимостей. В большинстве явлений хитросплетение множества различных связей настолько сложно, что предсказать заранее осуществление того или иного события не представляется возможным. Имеет место неоднозначность возможных исходов, действий, событий, как для каждого отдельного индивидуума, так и для групп населения, для живой и неживой природы.

Случайность проявляет себя во всех явлениях окружающего нас мира, будь то физические, химические, биологические, социальные и другие процессы. Однако доля неопределенности, доля случая, в разных ситуациях различна. Можно утверждать, что практически все явления происходят с той или иной степенью неопределенности, спектр которой простирается от полной непредсказуемости до несомненной однозначности исходов.

При такой точке зрения естественно ввести некоторую характеристику явлений, выражающую меру неопределенности (меру нашей уверенности) в исходе испытания. Наиболее удобным оказывается представление в виде числа. Меру уверенности в исходе испытания и называют *вероятностью*. Вероятность – объективная характеристика, которая должна зависеть не от нашего личного отношения к испытанию и его исходу, а от определенного набора условий, при котором проводится испытание. Наука, в основу которой положено понятие вероятности и которая занимается выявлением и изучением закономерностей в случайных явлениях, называется теорией вероятностей.

Теория вероятностей – наука математическая, так как широко использует математический аппарат и соответствующие методы исследования, но

имеющая обширные приложения, в том числе в медицине и здравоохранении. Данным фактом и объясняется её повсеместное использование в практике исследования. Численная интерпретация случайности позволяет находить специфические закономерности, специфические вероятностные законы, что в конечном итоге опровергает тезис о хаотичности, бессистемности случайных явлений. Чем далее развивается наука, тем более аргументированными становятся положения о вероятностной основе окружающего нас мира: мир построен на вероятности. Следует отметить, что теория вероятностей изучает только массовые явления, а именно: явления, которые могут быть повторены достаточно большое количество раз, а теоретически и бесконечное число раз. Теория вероятностей изучает события, обладающие статистической устойчивостью, т.е. события, относительная частота появления которых с ростом количества испытаний стабилизируется, колеблется около некоторого значения (именно это значение и будем считать вероятностью рассматриваемого события). Этот факт как раз и выражает закономерность в среде случайности. Теория вероятностей не изучает уникальные события: события, которые заведомо нельзя считать многократно повторяющимися или массовыми.

Также отметим объективную природу случайности, которой в массовых явлениях присущи свои специфические черты, не проявляющиеся в отдельных испытаниях. Классическая наука исследует, как правило, детерминистические закономерности, являющееся частным случаем вероятностных, т.е. события происходят с вероятностью 1 (в 100% испытаний).

Для изучения больших совокупностей относительно равноправных объектов (множество молекул, группы населения), индивидуальные черты не имеют определяющего значения, да и учесть их в полном объеме не представляется возможным. А возникающие закономерности поведения проявляются именно, как следствие массовости, как результат взаимодействия многих разнонаправленных хаотических движений.

Следует подчеркнуть особую роль в вероятностных исследованиях математической статистики, которая поставляет данные, необходимые для применения вероятностных методов. Наличие таких данных и дает основания к использованию этих методов в решении реальных проблем.

Учитывая эту специфику, можно утверждать, что во многих практических задачах в самых различных областях естествознания и деятельности человека вероятностный подход является наиболее целесообразным и соответствующим объективной истине, а вероятностная оценка результата является вполне удовлетворительной.

В частности, современное здравоохранение, имеющее объектом исследования множество индивидуумов, отличающихся друг от друга по всевозможным различным показателям, имеет вероятностную основу и широко использует вероятностные и статистические методы [1, 18, 24-28, 91, 103, 128-130]. Без знания этих методов невозможно осмысление целого ряда медико-социальных и биологических научных дисциплин, и эффективная деятельность медицинского работника существенно зависит от его компетентности в области применения указанных методов.

Методы теории вероятностей и математической статистики широко используются в математическом моделировании. Многие методы исследований путём математического моделирования основаны на соответствующих положениях теории вероятностей и математической статистики. Появляются новые направления исследований, в частности, изучение процессов на компьютерных моделях [130].

**Основные понятия и теоремы.** В теории случайных событий исходным понятием является *испытание*, т.е. всякий опыт, происходящий при некой совокупности условий, неоднозначно предопределяющей его исход. *Случайное событие* – один из возможных результатов рассматриваемого испытания. Например:

- испытание – подбрасывание стандартной монеты, случайные события – выпадение «орла» или «решки»;

- при изучении влияния типа погоды на состояние здоровья больных, страдающих ишемической болезнью сердца, испытание – это смена типа погоды, а случайные события – следующие исходы, связанные с состоянием здоровья: ухудшение, без перемен, улучшение;
- испытание – выбор медицинской карты больного при случайном отборе, случайные события – появление карты больного с диагнозом сахарный диабет, появление карты больного в возрасте 50-59 лет, появление карты больного мужского пола трудоспособного возраста и т.д.

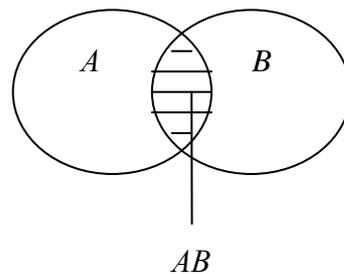
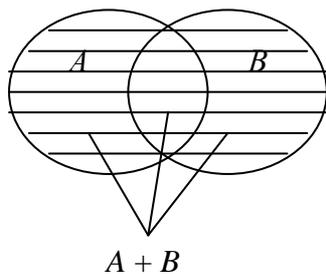
С каждым испытанием связана целая система случайных событий, которые принято обозначать большими буквами латинского алфавита:  $A, B, C, \dots$  Со случайными событиями можно производить различные математические операции: сложение, умножение и др.

Суммой случайных событий  $A$  и  $B$  называют случайное событие  $C$ , состоящее в осуществлении хотя бы одного из событий-слагаемых (т.е. осуществлении либо только  $A$ , либо только  $B$ , либо  $A$  и  $B$  вместе). Обозначение:  $A + B = C$ . Фактически речь идет об объединении множеств (рис. 2.1)

$$A + B = A \cup B.$$

Произведением случайных событий  $A$  и  $B$  называют случайное событие  $C$ , состоящее в осуществлении и события  $A$ , и события  $B$ . Обозначение:  $AB = C$  (или  $A \cdot B = C$ ). Фактически речь идет о пересечении множеств (рис. 2.2).

$$AB = A \cap B.$$



**Рис. 2.1.** Сумма случайных событий      **Рис. 2.2.** Произведение случайных событий

Случайное событие называют *достоверным* по отношению к данному испытанию, если оно осуществляется при любом исходе этого испытания. Условимся достоверное событие обозначать буквой  $U$ .

Случайное событие называют *невозможным* по отношению к данному испытанию, если оно неосуществимо при любом исходе этого испытания. Условимся обозначать невозможное событие буквой  $V$ .

Случайные события  $A$  и  $B$  называют *несовместными*, если их произведение – событие невозможное, т.е. события  $A$  и  $B$  совместно осуществиться не могут:  $AB=V$ , (в обозначениях теории множеств  $A \cap B = \emptyset$ ).

**Пример 2.1.** В больничном учреждении находятся на лечении больные с диагнозами  $a, b, c$ . При этом у одного и того же пациента возможно совместное наличие диагнозов  $a$  и  $b$ , диагноз  $c$  в сочетании с другими не встречается. Испытание – случайным образом выбранный больной. В качестве элементарных (простейших) случайных событий рассмотрим следующие:  $A$  – выбран больной с диагнозом  $a$ ;  $B$  – выбран больной с диагнозом  $b$ ;  $C$  – выбран больной с диагнозом  $c$ .

Тогда,  $AB$  – случайно выбранный больной имеет и диагноз  $a$ , и диагноз  $b$ ;  $A+B$  – случайно выбранный больной имеет или диагноз  $a$ , или диагноз  $b$ , или  $a$  и  $b$  вместе. События  $A$  и  $C$  несовместны, поскольку вместе  $A$  и  $C$  встретиться не могут. Также несовместны  $B$  и  $C$ . Таким образом,  $AC$  и  $BC$  – события невозможные. В то же время  $A+B+C$  – событие достоверное.

Систему случайных событий  $A_1, A_2, \dots, A_n$  называют *полной*, если в результате испытания обязательно осуществляется хотя бы одно из этих событий:

$$A_1 + A_2 + \dots + A_n = U.$$

Систему случайных событий  $A_1, A_2, \dots, A_n$  называют *несовместной*, если все события системы попарно несовместны, т.е.  $A_i A_k = V$  при  $i \neq k$ .

Систему случайных событий  $A_1, A_2, \dots, A_n$  называют *равновозможной*, если никакое из событий системы не имеет преимуществ в осуществлении по отношению к другим событиям системы.

Система случайных событий  $A_1, A_2, \dots, A_n$  называется *полной системой элементарных событий*, если она полная, несовместная и равновозможная.

**Пример 2.2.** При подбрасывании стандартной игральной кости (испытание) полную систему элементарных событий образуют шесть событий:

$A_1$  – выпадение одного очка;

$A_2$  – выпадение двух очков;

.....

$A_6$  – выпадение шести очков.

Остальные случайные события получаются из указанных элементарных с помощью соответствующих математических операций.

Важной характеристикой случайного события является вероятность его появления (мера неопределенности наших знаний об этом событии, шансы осуществления случайного события).

Пусть по отношению к проводимому испытанию задана полная система элементарных событий.

*Вероятностью* случайного события  $A$  называют число, обозначаемое  $P(A)$  и равное отношению  $m/n$ , где  $n$  – число всех событий в полной системе элементарных событий,  $m$  – число элементарных событий в системе, благоприятствующих осуществлению  $A$ .

Следовательно, имеет место простая формула

$$P(A) = \frac{m}{n},$$

согласно которой вероятность является долей элементарных событий, благоприятствующих  $A$ , среди всех возможных.

**Пример 2.3.** Пусть испытание – это извлечение случайным образом карты из стандартной колоды в 36 карт. Введем случайные события:  $A$  – извлечение «туза»;  $B$  – извлечение карты достоинством не ниже «валета». Тогда,  $P(A) = \frac{4}{36} = \frac{1}{9}$ ;  $P(B) = \frac{16}{36} = \frac{4}{9}$ . Вероятность оценивает шанс осуществления случайного события в отдельном испытании и выражается в долях от 1. Как видим, вероятность события  $B$  в 4 раза больше вероятности события  $A$ , т.е. при большом количестве испытаний событие  $B$  по сравнению с событием в  $A$  среднем имеет шанс осуществиться в 4 раза чаще.

Из определения вероятности вытекают следующие свойства:

- вероятность любого случайного события – число неотрицательное;

- вероятность любого случайного события не превосходит 1;
- вероятность достоверного события равна 1;
- вероятность невозможного события равна 0;
- вероятность является безразмерной величиной.

Приведенное выше определение вероятности называют *классическим*. Однако применимость его на практике ограничена, во-первых, условием равновозможности, что достигается в достаточно простых ситуациях или в искусственно организованных опытах (как в примере 1.3); во-вторых, – условием конечного числа элементарных событий в системе.

В реальных задачах указанные условия часто являются невыполнимыми. В этих случаях имеют место обобщения классического определения вероятности: статистическое и аксиоматическое.

Согласно *статистическому* определению в качестве реального значения вероятности может выступать относительная частота  $\mu/n$  осуществления события  $A$  в  $n$  испытаниях, проводящихся при одних и тех же условиях и при достаточно большом количестве испытаний. По определению, с увеличением  $n$  относительная частота  $\mu/n$ , где  $\mu$  – количество осуществлений события  $A$  в  $n$  испытаниях (частота), устойчиво стремится к значению теоретической вероятности. Поэтому в реальных задачах полагают, что  $P(A) \approx \mu/n$ . Таким образом, для нахождения вероятности случайного события  $A$  по статистическому определению предварительно необходимо провести  $n$  испытаний, в каждом из которых фиксируется два альтернативных исхода: произошло  $A$  или не произошло.

**Пример 2.4.** Статистические данные сообщают, что из 825 проведенных кардиохирургических операций, 42 закончились летальным исходом. На основе этого можно утверждать, что вероятность гибели больного составляет  $42/825 \approx 0,0509$ . Полученное значение является приближенным к истинному (теоретическому) значению вероятности летального исхода (случайного события) при операциях (испытаниях), проводимых по определенной методике при данном типе заболевания. Число 0,0509 – шанс умереть для среднестатистического прооперированного больного. Для конкретного больного эти шансы могут быть несколько

иными в зависимости от степени запущенности заболевания, наличия сопутствующих болезней, возраста пациента и т.д. Вероятность, получаемая статистически – это некая усредненная величина, характеризующая больных рассматриваемой совокупности. Она не утверждает, что произойдет в результате каждого конкретного испытания, а показывает шансы возможных исходов для среднестатистического больного.

*Аксиоматическое определение вероятности* базируется на системе аксиом Колмогорова. Согласно этой аксиоматике, вероятность рассматривается как функция множества случайных событий. Классическое определение вероятности является частным случаем аксиоматического определения. Заметим, что, умножив значение вероятности на 100%, шансы осуществления случайного события можно задать в процентах.

При действиях со случайными событиями вероятности изменяются по определенным законам. В частности, имеют место теоремы сложения и умножения.

*Теорема сложения.* По отношению к любому испытанию, для любых случайных событий  $A$  и  $B$  справедливо соотношение

$$P(A+B) = P(A) + P(B) - P(AB).$$

*Следствие.* Если события  $A$  и  $B$  несовместны, то  $P(A+B) = P(A) + P(B)$  (поскольку для несовместных событий  $P(AB) = 0$ ).

Два случайных события называют *противоположными*, если они образуют полную несовместную систему событий. Событие, противоположное событию  $A$ , принято обозначать  $\bar{A}$ . Справедливо равенство  $P(A) = 1 - P(\bar{A})$ .

Например, если событие  $A$  - это излечение больного за время  $t$ , то  $\bar{A}$  - неизлечение больного за то же время  $t$ .

*Условной вероятностью* события  $B$  при гипотезе  $A$  называется вероятность события  $B$  при условии, что событие  $A$  обязательно происходит. Эту условную вероятность будем обозначать  $P(B|A)$ .

Обязательное осуществление события  $A$  связано с изменением комплекса условий, при котором проводится испытание. При изменении ком-

плекса условий, продиктованным обязательным осуществлением события  $A$ , относительно события  $B$  возможны два варианта: вероятность события  $B$  сохраняет прежнее значение или изменяется.

Событие  $B$  называют *независимым* от события  $A$ , если условная вероятность  $B$  при гипотезе  $A$  равна безусловной вероятности события  $B$ , т.е.  $P(B|A) = P(B)$ .

Событие  $B$  называют *зависимым* от события  $A$ , если  $P(B|A) \neq P(B)$ .

*Теорема умножения.* По отношению к любому испытанию и для любых двух случайных событий  $A$  и  $B$  справедливы соотношения

$$P(AB) = P(A)P(B|A) \quad P(AB) = P(B)P(A|B).$$

Данная теорема имеет важные следствия.

*Следствие 1.* Условные вероятности для событий  $A$  и  $B$  при условии, что  $P(A) \neq 0$ ,  $P(B) \neq 0$ , можно вычислять по формулам

$$P(B|A) = \frac{P(AB)}{P(A)} \quad \text{и} \quad P(A|B) = \frac{P(AB)}{P(B)}.$$

*Следствие 2.* Свойство независимости событий  $A$  и  $B$  взаимно: если  $A$  не зависит от  $B$ , то и  $B$  не зависит от  $A$ .

Поэтому, говоря о независимости двух событий, полагаем, что имеет место взаимная независимость.

*Следствие 3.* Для независимых событий  $A$  и  $B$  соотношения в теореме умножения имеют более простой вид:  $P(AB) = P(A)P(B)$ .

**Пример 2.5.** Вернемся к условиям примера 2.1. и укажем конкретные численные данные. Всего в больничном учреждении находятся на лечении 100 пациентов с диагнозами  $a$ ,  $b$ ,  $c$ . Причем диагноз  $c$  в сочетании с другими диагнозами не встречается.

Диагноз	$a$	$b$	$c$	$a$ и $b$ вместе	хотя бы один из $a, b$	$a$ и $c$ вместе	$b$ и $c$ вместе
Число пациентов	52	63	15	30	85	0	0

В данном примере  $A, B, C$  – соответствующие случайные события. Требуется, исходя из вышеизложенной теории, вычислить вероятности различных событий.

Решение.

$$P(A) = \frac{52}{100}; \quad P(B) = \frac{63}{100}; \quad P(C) = \frac{15}{100};$$

$$P(\bar{A}) = \frac{48}{100}; \quad P(\bar{B}) = \frac{37}{100}; \quad P(\bar{C}) = \frac{85}{100};$$

$$P(AB) = \frac{30}{100}; \quad P(AC) = P(BC) = P(ABC) = 0;$$

$$P(A+B) = P(A) + P(B) - P(AB) = \frac{52}{100} + \frac{63}{100} - \frac{30}{100} = \frac{85}{100};$$

$$P(B+C) = P(B) + P(C) = \frac{63}{100} + \frac{15}{100} = \frac{78}{100};$$

$$P(A+C) = P(A) + P(C) = \frac{52}{100} + \frac{15}{100} = \frac{67}{100};$$

$$P(A|B) = \frac{P(AB)}{P(B)} = \frac{30/100}{63/100} = \frac{30}{63} = \frac{10}{21}; \quad P(B|A) = \frac{P(AB)}{P(A)} = \frac{30/100}{52/100} = \frac{30}{52} = \frac{15}{26}.$$

Поскольку  $P(A|B) \neq P(A)$ , то можно утверждать, что  $A$  и  $B$  – зависимые случайные события.

**Последовательность испытаний. Схема Бернулли.** Последовательность испытаний широко используется в исследовательской деятельности для проверки каких-либо гипотез опытным путем. Например, в здравоохранении с помощью последовательности испытаний проверяется эффективность лекарственных средств, новых методов профилактики, диагностики, лечения и т.п.

Пусть испытания рассматриваются по отношению к некоторому исследуемому событию  $A$ , когда в каждом испытании фиксируется лишь один из двух возможных исходов:  $A$  или  $\bar{A}$ . Предполагается, что все испытания проводятся при неизменных условиях, и вероятность осуществления события  $A$  в каждом из них одна и та же:  $P(A) = p$ . Тогда  $P(\bar{A}) = q$ , где  $q = 1 - p$ . Такие испытания с неизменной вероятностью  $p$  называют *однотипными*. Если осуществление или неосуществление события  $A$  в любом испытании не зависит от исходов других испытаний, предыдущих или последующих, то испытания

называют независимыми. Последовательность однотипных и независимых испытаний называют *схемой Бернулли*.

Пусть  $n$  испытаний относительно события  $A$  проводятся по схеме Бернулли. При этом во всей последовательности событие  $A$  осуществляется ровно  $m$  раз и не осуществляется ровно  $(n - m)$  раз, где  $m$  – некоторое целое число ( $0 \leq m \leq n$ ). Тогда справедливо вероятность осуществления события  $A$  ровно  $m$  раз в  $n$  испытаниях, проводимых по схеме Бернулли, можно найти по *формуле Бернулли*

$$P_n(m) = C_n^m p^m q^{n-m},$$

где  $P_n(m)$  – общепринятое обозначение указанной вероятности,  $C_n^m$  – общепринятое обозначение числа сочетаний по  $m$  элементов из множества в  $n$  элементов:  $C_n^m = \frac{n!}{m!(n-m)!}$ . Заметим, что по определению  $n! = 1 \cdot 2 \cdot 3 \cdot \dots \cdot n$  (читается « $n$ -факториал»), а  $0! = 1$ .

**Пример 2.5.** Согласно статистическим наблюдениям после проведенных в центре сердечнососудистой хирургии кардиохирургических операций послеоперационная летальность составляет 7%. Найти вероятность, что после 10 намеченных на следующую неделю операций число выживших будет: а) ровно 8 чел; б) 8 и более чел.

Решение. Так как вероятность смертельного исхода для каждого конкретного пациента равна  $\frac{7\%}{100\%} = 0,07$ , то вероятность выжить оказывается равной  $p = 1 - 0,07 = 0,93$ . Полагаем эти вероятности неизменными для любого из прооперированных больных, т.е. испытания (операции) являются однотипными. Вполне естественно допустить, что результаты каждой операции не зависят от результатов других. Таким образом, есть основания полагать, что испытания проводятся по схеме Бернулли. Тогда искомые вероятности легко вычисляются:

$$\text{а) } P_{10}(8) = C_{10}^8 p^8 (1-p)^2 = \frac{10!}{8! 2!} (0,93)^8 (0,07)^2 = 0,1234$$

$$\text{б) } P_{10}(m \geq 8) = P_{10}(8) + P_{10}(9) + P_{10}(10) = 0,1234 + C_{10}^9 p^9 (1-p) + C_{10}^{10} p^{10} (1-p)^0 = 0,9717.$$

Вывод. Вероятность того, что число выживших после намеченных 10 операций будет ровно 8 чел. составляет 0,1234; 8 и более чел. – 0,9717.

Схема Бернулли является универсальным инструментом многих теоретических и прикладных исследований. Однако использование формулы Бернулли, лежащей в основе данной схемы и, безусловно, справедливой при любых натуральных  $n$ , часто сопряжено с громоздкими расчетами. Исходя из этого, при больших значениях  $n$  вместо вычислений по точной формуле Бернулли принято использовать приближенные соотношения, известные как приближение Муавра-Лапласа и приближение Пуассона.\*

## 2.2 Случайные величины

**Понятие случайной величины.** *Случайная величина* – это переменная, которая в результате испытания однозначно принимает некоторое численное значение, но какое именно априори, т.е. до проведения испытания, неизвестно.

Условимся обозначать случайные величины большими буквами латинского алфавита:  $X, Y, Z, T$  и т.д. и использовать сокращение «с.в.». В результате испытания с.в. может принимать одно из значений, для обозначения которых используем соответствующие малые буквы:  $x, y, z, t$  и т.д. Также будем пользоваться буквами с индексами, например, значения с.в.  $X$  – это  $x_1, x_2, \dots, x_n$ .

С.в. является числовой функцией, заданной на множестве элементарных событий, каждому из которых соответствует численное значение с.в. Над случайными величинами можно производить различные математические операции: сложение, умножение, деление, извлечение корня и др.. Вместе с тем, каждое из возможных значений с.в., связанное со случайным событием, характеризуется своей персональной вероятностью.

Случайные величины классифицируют в зависимости от множества их возможных значений. Наиболее важными для статистических исследований и приложений являются с.в. двух типов: дискретного и непрерывного.

---

\* Подробно см.: Математическая статистика в медицине: учеб. пособие / В.А. Медик, М.С. Токмачев. – М.: Финансы и статистика, 2007, с.64-69.

С.в.  $X$  называют *дискретной*, если множество ее значений конечно или счетно, например: число новорожденных за месяц, число ударов пульса больного в минуту, число заболевших во время эпидемии гриппа.

С.в.  $X$  называется *непрерывной*, если все ее значения заполняют сплошь некоторый интервал, например: температура больного (точное значение  $t = 36,57369745... \text{C}^0$  округляют до ближайшего значения деления термометра  $t = 36,6 \text{C}^0$ ); величина ошибки, допускаемой при округлении точных значений.

**Распределение случайной величины.** Любая случайная величина в результате испытания принимает одно из значений некоторого множества. Однако, знание всех возможных значений с.в. для ее характеристики недостаточно: необходимо знать как часто (с какой вероятностью) она принимает те или иные значения.

*Распределением (законом распределения) с.в.* называют соответствие между ее возможными значениями и соответствующими этим значениям вероятностями. Распределение можно задать таблично, аналитически (в виде формулы) и графически.

### 2.2.1 Дискретные случайные величины

Для дискретной с.в.  $X$  можно перечислить все значения  $x_1, x_2, \dots, x_n, \dots$  и их вероятности  $p_1, p_2, \dots, p_n, \dots$ , которые удобно представлять в виде таблицы

$x_i$	$x_1$	$x_2$	...	$x_n$	...
$p_i$	$p_1$	$p_2$	...	$p_n$	...

где  $p_i = P(X = x_i)$ , При этом обязательно должно выполняться *условие нормировки*:  $\sum_i p_i = 1$ . Таким образом, термин «распределение с.в.» означает

распределение единичной вероятности между всеми возможными значениями. Из бесконечного многообразия распределений с.в. можно выделить некоторые типичные, классические. В частности, отметим ряд дискретных рас-

пределений, имеющих важные приложения в медицинской статистике: дискретное равномерное, биномиальное, распределение Пуассона, геометрическое, отрицательное биномиальное, гипергеометрическое, распределение Пойа.

Функцией распределения с.в.  $X$  называют функцию одной переменной  $F(x)$  такую, что

$$F(x) = P(X < x), \text{ при } x \in (-\infty, \infty).$$

Функция распределения  $F(x)$  равна вероятности попадания с.в.  $X$  в интервал  $(-\infty, x)$ , а, следовательно, как и любая вероятность может принимать значения лишь от 0 до 1. По функции распределения однозначно определяется закон распределения. Отметим, что вероятность попадания с.в.  $X$  в какой-либо интервал  $[a, b)$  равна разности  $F(b) - F(a)$ . График функции распределения  $F(x)$  дискретной с.в.  $X$  всегда имеет ступенчатый вид (рис. 2.3). Величина каждого скачка равна вероятности  $p_i$  появления значения  $x_i$  с.в.  $X$ .

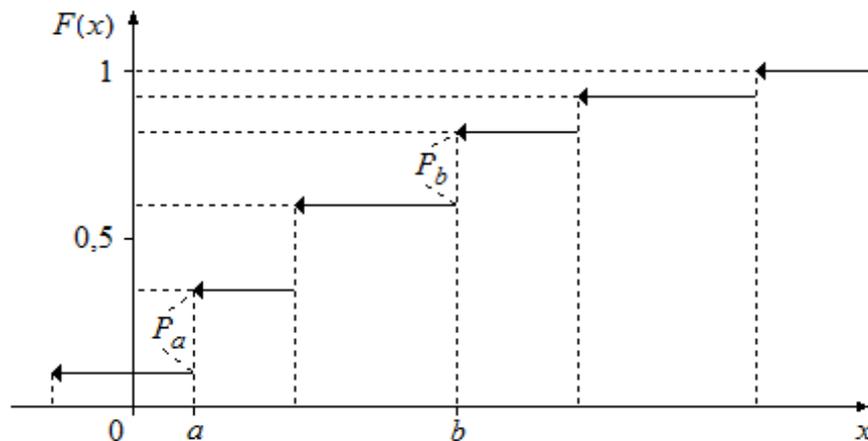


Рис. 2.3. Пример графика функции распределения  $F(x)$  дискретной с.в.  $X$

**Числовые характеристики дискретных случайных величин.** Закон распределения дискретной с.в. исчерпывающе характеризует ее. В то же время, исходя из закона распределения, можно вычислить частные характеристики, которые подчеркивают отдельные стороны поведения с.в.

Начальным моментом  $k$ -го порядка,  $k = 0, 1, 2, \dots$ , дискретной с. в.  $X$  называют число

$$\alpha_k(X) = \sum_i x_i^k p_i,$$

где  $x_i$  представляют все возможные значения с.в.  $X$ , а  $p_i$  – вероятности, соответствующие значениям  $x_i$ .

*Математическим ожиданием* дискретной с. в.  $X$  называют ее первый начальный момент:

$$M(X) = \alpha_1(X) = \sum_i x_i p_i.$$

Смысл математического ожидания:  $M(X)$  – среднее значение с.в.  $X$ , вычисленное с учетом вероятности каждого из значений. Исходя из этого, математическое ожидание часто именуют термином «среднее».

*Центральным моментом*  $k$ -го порядка,  $k = 0, 1, 2, \dots$ , дискретной с.в.  $X$  называют число

$$\mu_k(X) = \sum_i (x_i - M(X))^k p_i,$$

где  $x_i$  – значения с.в.  $X$ ,  $p_i$  – вероятности, соответствующие значениям  $x_i$ ,  $M(X)$  – математическое ожидание с.в.  $X$ . Суммирование проводится по всем возможным значениям  $i$ .

*Дисперсией* дискретной с.в.  $X$  называют ее центральный момент второго порядка:

$$D(X) = \mu_2(X) = \sum_i (x_i - M(X))^2 p_i.$$

Дисперсия всегда неотрицательна и указывает степень отклонения значений с.в. от среднего. Чем больше разброс значений относительно среднего, тем больше дисперсия.

Символ математического ожидания  $M(X)$  удобно использовать и при обозначении моментов:

$$\begin{aligned} \alpha_k(X) &= \sum_i x_i^k p_i = M(X^k), \\ \mu_k(X) &= \sum_i (x_i - M(X))^k p_i = M(X - M(X))^k. \end{aligned}$$

Для дисперсии любой с.в.  $X$  справедливо равенство

$$D(X) = M(X^2) - M^2(X).$$

Средним квадратическим отклонением (стандартным отклонением)

с.в.  $X$  называют число  $\sigma = \sqrt{D(X)}$ .

Стандартное отклонение  $\sigma$  характеризует степень отклонения значений с.в.  $X$  от среднего  $M(X)$  в абсолютных единицах. Как легко заметить, всегда по определению,  $\sigma \geq 0$ .

Медианой с.в.  $X$  называют значение  $x_{0,5}$ , при котором для функции распределения  $F(x)$  выполняется равенство  $F(x_{0,5}) = 0,5$ .

Медиана представляет собой середину распределения, т.е. точку, в которой вся вероятностная масса (единица) делится пополам. Отсюда и индекс 0,5 и сам термин «медиана». Обозначим  $x_{0,5} = Me(X)$ .

Модой дискретной с.в.  $X$  называют значение с.в., вероятность появления которого является наибольшей по сравнению с соседними значениями.

В конкретном распределении мода может быть не единственной, а может и вовсе отсутствовать. Обозначим моду как  $Mo(X)$ .

**Пример 2.6.** Закон распределения с.в.  $X$  имеет вид

$x_i$	-1	0	1	2
$p_i$	0,2	0,4	0,3	0,1

Требуется найти вышеуказанные числовые характеристики с.в.  $X$ .

Решение.

$$M(X) = \sum_i x_i p_i = (-1) \cdot 0,2 + 0 \cdot 0,4 + 1 \cdot 0,3 + 2 \cdot 0,1 = 0,3;$$

$$M(X^2) = \alpha_2(X) = \sum_i x_i^2 p_i = (-1)^2 \cdot 0,2 + 0^2 \cdot 0,4 + 1^2 \cdot 0,3 + 2^2 \cdot 0,1 = 0,9;$$

$$D(X) = M(X^2) - M^2(X) = 0,9 - (0,3)^2 = 0,81; \quad \sigma = \sqrt{D(X)} = \sqrt{0,81} = 0,9;$$

$$Me(X) = 0; \quad Mo(X) = 0.$$

Квантилью порядка  $p$  называют значение  $x_p$  с.в.  $X$ , для которого справедливо равенство  $F(x_p) = p$ , где  $p$  может быть любым числом из интервала  $(0, 1)$ .

Согласно последнему определению медиана является квантилью порядка 0,5.

Квантили порядка 0,25 и 0,75 называют соответственно *нижней и верхней квантилями* ( $x_{0,25}$  и  $x_{0,75}$ ). Разность  $x_{0,75} - x_{0,25}$  называют *интерквантильным расстоянием*.

*Вероятное отклонение* характеризует степень рассеяния значений случайной величины с симметричным распределением и равно  $x_{0,75} - x_{0,25}$ . *Ширина распределения* – это длина интервала между граничными значениями  $x$ , при которых  $F(x) = 0$  и  $F(x) = 1$ .

*Коэффициентом асимметрии* ( $As$ ) распределения с.в.  $X$  называют число  $\gamma_1$ , определяемое выражением:

$$As(X) = \gamma_1 = \frac{\mu_3}{\sigma^3}.$$

Для симметричных распределений  $\gamma_1 = 0$ . При  $\gamma_1 > 0$  автоматически следует  $\mu_3 > 0$ , т.е. положительные отклонения превышают отрицательные, более «длинная часть» распределения находится справа от среднего. Асимметрия положительна. Аналогично при  $\gamma_1 < 0$  асимметрия отрицательна, «длинная часть» распределения находится слева от среднего.

*Коэффициентом эксцесса* ( $Ex$ ) распределения с.в.  $X$  называют следующее число  $\gamma_2$ :

$$Ex(X) = \gamma_2 = \frac{\mu_4}{\sigma^4} - 3.$$

Коэффициент эксцесса характеризует сглаженность распределения с.в. относительно центра, крутость распределения.

Отметим, что для решения ряда статистических задач оказывается достаточным знание лишь конкретных числовых характеристик распределения. Например, для анализа эффективности использования средств, выделяемых на обеспечение льготных категорий граждан бесплатными лекарствами, рассчитывается в основном показатель средней стоимости одного рецепта (математическое ожидание). Или другой пример, для оценки эффективности использования коечного фонда рассматривается длительность пребывания

больного на койке (случайная величина) и для нее рассчитывается показатель средней длительности пребывания больного на койке (математическое ожидание). В данном случае также важно знание и разброса значений этой случайной величины (дисперсии).

### **Некоторые типы распределений дискретных случайных величин.**

Приведем некоторые распределения с.в., которые используются или могут быть полезны для здравоохранения.

*Биномиальное распределение. Распределение Бернулли.* Биномиальным распределением с.в.  $X$  называют распределение вида

$x_m$	0	1	2	...	$n$
$p_m$	$p_0$	$p_1$	$p_2$	...	$p_n$

При этом вероятности  $p_m$  задаются формулой Бернулли

$$p_m = C_n^m p^m q^{n-m}, \quad 0 < p < 1, \quad q = 1 - p.$$

Биномиальное распределение зависит от двух параметров  $n$  и  $p$ , где  $n$  – натуральное число, а  $p$  имеет смысл вероятности в испытаниях, проводимых по схеме Бернулли. Сам факт биномиального распределения с.в.  $X$  с параметрами  $n$  и  $p$  принято обозначать  $X \sim B(n, p)$ \*

Можно доказать, что для биномиального закона приведенного, обычно используемого вида, справедливо:  $M(X) = np$  и  $D(X) = npq$ . В общем случае с.в., вероятности  $p_m$  которой задаются формулой Бернулли, может принимать значения с произвольным шагом. Например:  $x_0 = -3$ ,  $x_1 = 2$ ,  $x_2 = 9$ . В этом случае речь идет о некоторой функции стандартной биномиальной случайной величины, а её числовые характеристики,  $M(X)$ ,  $D(X)$  и др., определяются согласно соотношениям для произвольных распределений.

---

\* Символ  $\sim$  читается «распределена, как» или «распределена по закону».

Биномиальное распределение имеет место тогда, когда  $n$  испытаний осуществляются по схеме Бернулли относительно некоторого случайного события  $A$ . В качестве с.в.  $X$  рассматривается число удачных испытаний.

Важным частным случаем биномиального распределения при  $n = 1$  является распределение Бернулли

$x_m$	0	1
$p_m$	$q$	$p$

Для этого распределения, как легко заметить,  $M(X) = p$ ;  $D(X) = pq$ .

**Пример 2.7.** *Комбинации полов в многодетных семьях.* Среди многодетных семей с пятью детьми вычислить процентные доли семей с различным соотношением детей по половому признаку. Вероятность рождения мальчика считать равной 0,516.

Решение. Среди пятерых детей в семье возможны следующие соотношения:

- все 5 детей – девочки;
- 4 девочки и 1 мальчик;
- 3 девочки и 2 мальчика;
- 2 девочки и 3 мальчика;
- 1 девочка и 4 мальчика;
- все 5 детей – мальчики.

Взяв в качестве события  $A$  рождение в рассматриваемой семье, например, мальчика, оказываемся в условиях схемы Бернулли. Полагаем с. в.  $X$  – количество мальчиков в семье из 5 детей. Тогда с.в.  $X$  имеет биномиальное распределение с параметрами  $n=5$ ,  $p = 0,516$ :

$x_i$	0	1	2	3	4	5
$p_i$	$p_0$	$p_1$	$p_2$	$p_3$	$p_4$	$p_5$

где вероятности  $p_i$  находятся по формуле Бернулли  $p_m = C_n^m (0,516)^m (0,484)^{n-m}$ .

Вычислив указанные вероятности с точностью до  $10^{-5}$ , и переведя их в соответствующие проценты, получим искомые процентные доли семей.

$x_i$	0	1	2	3	4	5	
$p_i$	0,02656	0,14158	0,30188	0,32184	0,17156	0,03658	$\Sigma = 1$
Доли в %	2,656	14,158	30,188	32,184	17,156	3,658	$\Sigma = 100\%$

Вывод. Таким образом, в 2,656% семей с пятью детьми все дети – девочки, в 14,158% семей - только один мальчик и т.д. Найденные процентные соотношения являются теоретическими. Отклонение реальных данных от теоретических при достаточно большом наблюдаемом количестве семей и при вероятности  $p = 0,516$

не должно быть существенным. В противном случае есть повод усомниться в предположении  $p = 0,516$  и необходимо уточнить значение этой вероятности статистическими методами.

*Распределение Пуассона.* Распределением Пуассона с.в.  $X$  называют распределение вида

$x_m$	0	1	2	...	$m$	...
$p_m$	$p_0$	$p_1$	$p_2$	...	$p_m$	...

где вероятности  $p_m$  вычисляются по формуле

$$p_m = \frac{a^m}{m!} e^{-a}, (a > 0).$$

Распределение Пуассона зависит лишь от одного параметра  $a$ . Можно доказать, что для пуассоновской с.в.  $X$  оказывается  $M(X) = a$ ;  $D(X) = a$ ,

$$As(X) = \frac{1}{\sqrt{a}}, \quad Ex(X) = \frac{1}{a}.$$

Следует отметить, что распределение Пуассона является предельным для биномиального распределения при достаточно больших значениях  $n$  и малых значениях вероятности  $p$  (редких событиях).

**Пример 2.8.** *Травматизм на производстве.*

Рассматривается количество случаев травматизма с ВУТ (временной утратой трудоспособности) на крупном промышленном предприятии за месяц. Все работники предприятия разбиты на  $N$  подразделений с примерно равным количеством рабочих. С.в.  $X$  – количество случаев травматизма в подразделении (редкое событие) – распределена по закону Пуассона с параметром  $a=0,8$ , где  $a$  – статистическое среднее по отрасли количество случаев травматизма с ВУТ в таком подразделении.

$x_i$	0	1	2	3	4	5	> 5
$p_i$	0,44933	0,35946	0,14379	0,03834	0,00767	0,00123	0,00018
$\omega_i$	$\omega_0$	$\omega_1$	$\omega_2$	$\omega_3$	$\omega_4$	$\omega_5$	$\omega_{>5}$

Вывод. Реальные значения относительных частот  $\omega_i$  находятся по формуле  $\omega_i = \frac{m_i}{N}$ , где  $m_i$  – количество подразделений со значением  $x_i$ . Сравнивая наблюдаемые значения относительных частот ( $\omega_i$ ) с теоретическими ( $p_i$ ), можно судить

об уровне техники безопасности в подразделении по отношению со средним уровнем по отрасли.

### 2.2.2 Непрерывные случайные величины

Множество возможных значений непрерывной с.в. – это множество значений, заполняющих сплошь некоторый интервал. Например, время ожидания приезда бригады скорой медицинской помощи – это непрерывная величина со значениями из интервала 5-25 минут (в соответствии с установленными нормативами).

Вероятность любого конкретного значения непрерывной с.в. равна 0. Поэтому вероятности значений непрерывных с.в. рассматриваются не для отдельных значений (точек на числовой оси), а для совокупности значений (интервалов на числовой оси). В отличие от дискретных с.в. для изучения непрерывных с.в., в силу их специфики, применяются методы дифференциального и интегрального исчисления. Вероятность распределения непрерывной с.в. описывается функцией – плотностью распределения вероятностей.

*Плотностью распределения вероятностей* (функцией плотности) непрерывной случайной величины  $X$  называют неотрицательную, определенную при всех  $x$ , функцию  $f(x)$ , удовлетворяющую условию:

$$P(X \in [a, b)) = \int_a^b f(x) dx.$$

Как следует из определения, функция распределения  $F(x)$  выражается через функцию плотности распределения  $f(x)$  по формуле

$$F(x) = \int_{-\infty}^x f(t) dt.$$

Плотность распределения  $f(x)$  восстанавливается по функции распределения  $F(x)$  дифференцированием:

$$f(x) = F'(x).$$

Для плотности распределения  $f(x)$  выполнено обязательное *условие нормировки*

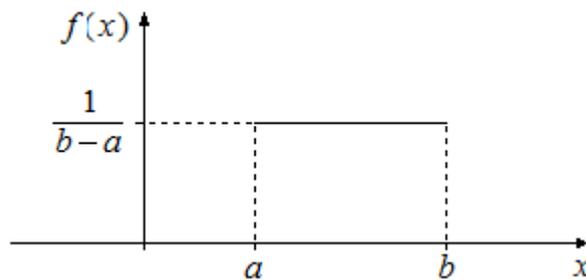
$$\int_{-\infty}^{+\infty} f(x)dx = 1.$$

Для непрерывных случайных величин, как и для дискретных, вводятся понятия математического ожидания ( $M(X)$ ), дисперсии ( $D(X)$ ), начальных и центральных моментов, имеющие ту же самую интерпретацию. Однако при этом вместо операции суммирования для непрерывных случайных величин используется операция интегрирования. В частности,

$$M(X) = \int_{-\infty}^{+\infty} x f(x)dx; \quad D(X) = \int_{-\infty}^{+\infty} (x - M(X))^2 f(x)dx.$$

**Некоторые типы распределений непрерывных случайных величин.** Для непрерывных с.в., также как и для дискретных, существует бесконечное количество возможных распределений, поскольку каждое распределение – это способ рассредоточить единичную вероятность между значениями какого-то интервала. Приведем некоторые распределения, типичные для прикладных исследований.

**Равномерное распределение на отрезке.** Простейшим является случай, когда вероятность распределена «поровну» между всеми возможными значениями интервала, что соответствует постоянной величине плотности во всех точках заданного интервала. Такое распределение называется равномерным на отрезке (рис. 2.4).



**Рис. 2.4.** График плотности равномерного распределения

В соответствии с условием нормировки функция плотности распределения однозначно определяется выражением

$$f(x) = \begin{cases} \frac{1}{b-a}, & \text{при } x \in [a, b] \\ 0, & \text{при } x \notin [a, b]. \end{cases}$$

При этом граничные значения интервала  $a$  и  $b$  являются параметрами распределения, разность  $(b - a)$  – длина интервала,

$$M(X) = \frac{a+b}{2}, \quad D(X) = \frac{(b-a)^2}{12}.$$

**Пример 2.9.** Анализируется работа станции скорой медицинской помощи в городе N. Непрерывная случайная величина  $X$  – время ожидания приезда бригады скорой медицинской помощи имеет равномерное распределение на отрезке  $[5; 25]$  минут. Найти плотность распределения вероятностей этой с.в.,  $M(X)$ ,  $D(X)$ .

Решение. В соответствии с приведенными формулами, получаем:

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{при } x \in [a; b]; \\ 0 & \text{при } x \notin [a; b] \end{cases} = \begin{cases} \frac{1}{25-5} & \text{при } x \in [5; 25]; \\ 0 & \text{при } x \notin [5; 25] \end{cases} = \begin{cases} 0,05 & \text{при } x \in [5; 25]; \\ 0 & \text{при } x \notin [5; 25]. \end{cases}$$

$$M(X) = \frac{a+b}{2} = \frac{5+25}{2} = 15.$$

Таким образом, пациенты ожидают приезда бригады скорой медицинской помощи в среднем 15 минут. При этом

$$D(X) = \frac{(b-a)^2}{12} = \frac{(25-5)^2}{12} = 33,33; \quad \sigma = \sqrt{D(X)} = 5,77.$$

Вывод. Отклонение от среднего на величину 5,77 минуты (т.е. значения на отрезке  $[9,23; 20,77]$ ) считается допустимым. С учетом рекомендуемого времени приезда бригады скорой медицинской помощи на вызов (20 минут) работу станции скорой медицинской помощи города N по этому показателю можно считать удовлетворительной.

**Нормальное распределение.** Данное распределение является важнейшим из вероятностных распределений, наиболее часто встречающимся в прикладных исследованиях, в том числе в медицинских.

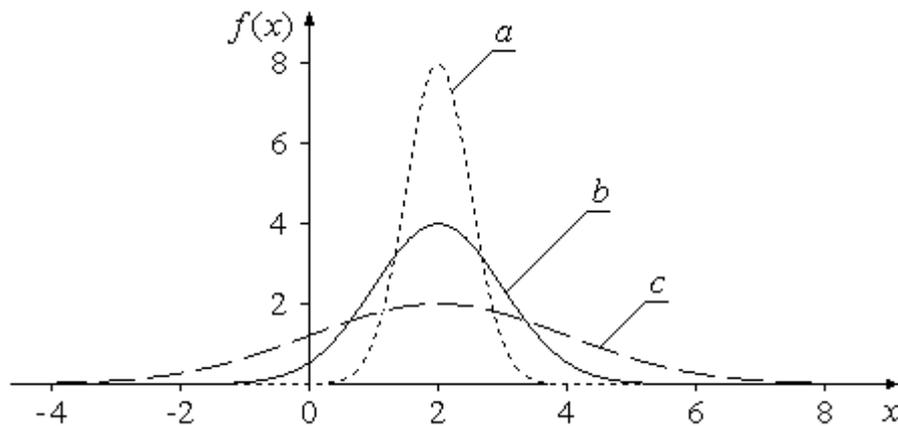
Непрерывная случайная величина  $X$  называется *нормально распределенной* с параметрами  $m$  и  $\sigma$  ( $\sigma > 0$ ), если ее плотность распределения вероятностей имеет вид:

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-m)^2}{2\sigma^2}},$$

что обозначается следующим образом:  $X \sim N(m; \sigma)$ .

Нормальное распределение является двухпараметрическим распределением. В случае  $m = 0$ ,  $\sigma = 1$  нормальное распределение называют *стандартным нормальным распределением*.

График функции  $f(x)$  называют нормальной кривой или *кривой Гаусса*. Вид этого графика при различных значениях параметров приведен на рис. 2.5.



**Рис. 2.5.** Кривая Гаусса при  $m=2$  и  $\sigma=0,5$  (a),  $\sigma=1$  (b),  $\sigma=2$  (c)

Можно доказать, что для нормального распределения  $M(X) = m$ ,  $D(X) = \sigma^2$ ,  $As(X) = 0$  и  $Ex(X) = 0$ .

Для функции распределения  $F(x)$  справедливо соотношение

$$F(x) = \Phi\left(\frac{x-m}{\sigma}\right),$$

где  $m = M(X)$ ,  $\sigma = \sqrt{D(X)}$ ,  $\Phi\left(\frac{x-m}{\sigma}\right)$  — известная функция Лапласа

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-0,5t^2} dt,$$

значения которой затабулированы (табл. П 1).

Для стандартного нормального распределения  $X \sim N(0, 1)$  функция распределения полностью совпадает с функцией Лапласа:  $F(x) = \Phi(x)$ .

Широкая распространенность нормального распределения в теоретических и прикладных исследованиях объясняется центральной предельной теоремой, согласно которой при некоторых условиях утверждается: *если случай-*

ная величина  $X$  является суммой достаточно большого количества взаимно независимых, произвольно распределенных случайных величин, вклад каждой из которых во всю сумму незначителен, то распределение с.в.  $X$  приближенно нормально с соответствующими параметрами  $m$  и  $\sigma$ .

Среди других непрерывных распределений с.в. выделим: показательное, логарифмически нормальное (логнормальное), гамма-распределение. Для статистики особое место, кроме нормального, занимают распределения:  $\chi^2$  – распределение (хи-квадрат распределение Пирсона),  $t$ -распределение (распределение Стьюдента),  $F$ -распределение (распределение Фишера-Снедекора, распределение дисперсионного отношения), рассматриваемые в разделе 3.3.

**Числовые характеристики распределений непрерывных случайных величин** определяются по тем же принципам, что и для дискретных с.в., но вычисляются иначе. Основными являются уже упоминавшиеся моменты (начальные и центральные), среди которых особая роль отводится математическому ожиданию и дисперсии\*.

*Модой* распределения непрерывной с.в.  $X$  называют значение с.в., при котором функция плотности  $f(x)$  достигает своего максимума.

*Квантилью порядка  $p$*  распределения непрерывной с.в.  $X$  называют значение  $x_p$  такое, что для функции распределения  $F(x)$  справедливо равенство  $F(x_p) = p$ , где  $0 < p < 1$ .

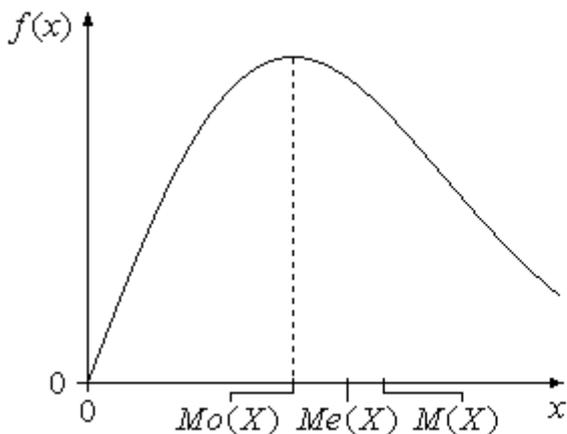
*Медианой* распределения называют квантиль порядка 0,5. Квантили порядка 0,25 и 0,75 называют соответственно *нижней и верхней квантилями*. Очевидно, чем больше значение рассматриваемой  $F(x)$ , тем больше и соответствующий ему  $x_p$ . Квантили  $x_{0,1}, x_{0,2}, \dots, x_{0,9}$  называют *децилями*, а  $x_{0,01}, x_{0,02}, \dots, x_{0,99}$  – *процентильями (персентильями)*.

---

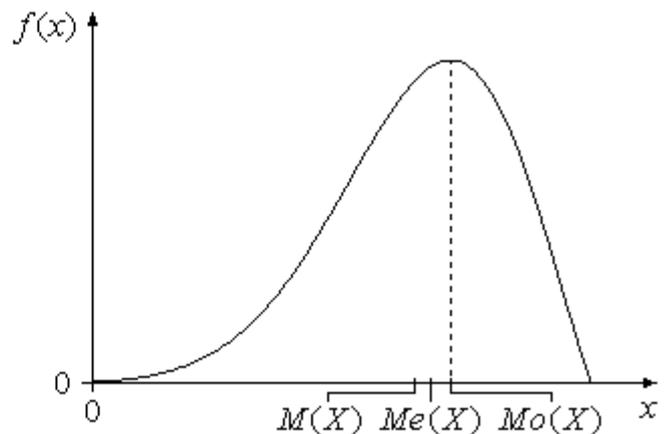
\* Подробнее см.: Математическая статистика в медицине: учеб. пособие / В.А. Медик, М.С. Токмачев. – М.: Финансы и статистика, 2007. – С. 136-137.

Коэффициенты асимметрии и эксцесса определяются теми же выражениями, что и для дискретных с.в.:  $As(X) = \frac{\mu_3}{\sigma^3}$ ,  $Ex(X) = \frac{\mu_4}{\sigma^4} - 3$ .

Для симметричных распределений значения  $M(X)$ ,  $Mo(X)$  и  $Me(X)$  совпадают. Для асимметричных распределений математическое ожидание, мода и медиана различны. Для распределения с правосторонней (положительной) асимметрией  $Mo(X) < Me(X) < M(X)$  (рис. 2.6), а для распределений с левосторонней (отрицательной) асимметрией  $M(X) < Me(X) < Mo(X)$  (рис. 2.7).



**Рис. 2.6.** Распределение с положительной асимметрией (распределения Вейбулла, Релея, Фишера и др.) и положение его характеристик



**Рис. 2.7.** Распределение с отрицательной асимметрией (часть бета-распределений, типа синуса и др.) и положение его характеристик

Отметим, что указанные выше числовые характеристики для конкретных распределений легко найти с помощью ЭВМ.

## 2.3 Системы случайных величин

### 2.3.1 Основные понятия и характеристики

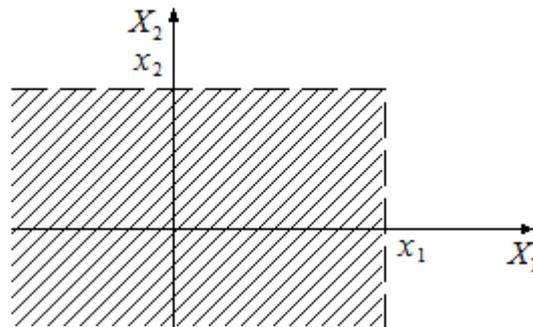
Системой случайных величин называют  $n$  случайных величин  $X_1, X_2, \dots, X_n$ , рассматриваемых в совокупности. Например, многомерный показатель общей заболеваемости – система случайных величин, где  $X_1, X_2, \dots, X_n$  – отдельные нозологические формы. Значения случайных величин обозначают соответственно  $x_1, x_2, \dots, x_n$ , т.е.  $x_k$  – переменная величина, включающая в себя

все возможные значения с.в.  $X_k$ . Каждому значению  $x_1$  одной случайной величины  $X_1$  соответствует точка на прямой, каждой паре значений  $(x_1, x_2)$  системы с.в.  $X_1, X_2$  соответствует точка на плоскости, каждой тройке значений  $(x_1, x_2, x_3)$  – точка в трехмерном пространстве. При  $n > 3$  формально говорят о точке  $n$ -мерного пространства.

**Функцией распределения** системы случайных величин  $X_1, X_2, \dots, X_n$  называют функцию  $n$  переменных  $F(x_1, x_2, \dots, x_n)$ , равную вероятности произведения событий  $X_1 < x_1, X_2 < x_2, \dots, X_n < x_n$ :

$$F(x_1, x_2, \dots, x_n) = P(X_1 < x_1, X_2 < x_2, \dots, X_n < x_n).$$

В частности, при  $n = 2$  функция распределения  $F(x_1, x_2)$  – это вероятность попадания значений случайной точки  $(X_1, X_2)$ , в заштрихованную область на рис. 2.8.



**Рис. 2.8.** Область, вероятность попадания в которую случайной точки  $(X_1, X_2)$  равна  $F(x_1, x_2)$

Для системы  $n$  непрерывных случайных величин  $X_1, X_2, \dots, X_n$  вводится понятие *совместной плотности распределения*  $f(x_1, x_2, \dots, x_n)$ .

Случайные величины  $X_1, X_2, \dots, X_n$  называют *независимыми*, если вероятность произведения событий  $X_1 < x_1, X_2 < x_2, \dots, X_n < x_n$  равна произведению вероятностей событий-сомножителей для всех значений  $x_1, x_2, \dots, x_n$ :

$$P(X_1 < x_1, X_2 < x_2, \dots, X_n < x_n) = P(X_1 < x_1) \cdot P(X_2 < x_2) \cdot \dots \cdot P(X_n < x_n).$$

Если указанное равенство не выполнено хотя бы для какого-то набора значений  $x_1, x_2, \dots, x_n$ , то случайные величины  $X_1, X_2, \dots, X_n$  называют *зависимыми*.

Для независимых случайных величин  $X_1, X_2, \dots, X_n$ , справедливы соотношения

$$F(x_1, x_2, \dots, x_n) = F_1(x_1) \cdot F_2(x_2) \cdot \dots \cdot F_n(x_n),$$

$$f(x_1, x_2, \dots, x_n) = f_1(x_1) \cdot f_2(x_2) \cdot \dots \cdot f_n(x_n),$$

где  $F_1(x_1), F_2(x_2), \dots, F_n(x_n)$  – функции распределения;  $f_1(x_1), f_2(x_2), \dots, f_n(x_n)$  – плотности распределения каждой из случайных величин  $X_1, X_2, \dots, X_n$  соответственно. Справедливы и обратные утверждения: если имеют место приведенные равенства, то рассматриваемые случайные величины независимые.

Для определения степени зависимости случайных величин используется понятие корреляции. *Корреляция* (от лат. correlatio – соотношение) – численная характеристика взаимозависимости двух с.в.

Корреляция может быть выражена в виде корреляционного момента и коэффициента корреляции.

*Корреляционным моментом* двух случайных величин  $X, Y$  называют функцию  $K(X, Y)$ , которая для дискретных случайных величин находится по формуле

$$K(X, Y) = \sum_i \sum_j (x_i - M(X))(y_j - M(Y)) p_{ij},$$

где  $x_i, y_j$  – все возможные значения случайных величин  $X$  и  $Y$ , а  $p_{ij} = P(X = x_i, Y = y_j)$ .

Для непрерывных случайных величин корреляционный момент вычисляется с помощью операции интегрирования

$$K(X, Y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - M(X))(y - M(Y)) f(x, y) dx dy.$$

Другое название корреляционного момента – *ковариация*\*. Соответствующее обозначение –  $\text{cov}(X, Y)$ .

Можно доказать, что для независимых случайных величин их корреляционный момент равен нулю. Обратное, вообще говоря, неверно: из равенства  $K(X, Y) = 0$  не следует независимость случайных величин  $X, Y$ . Несмотря на условие  $K(X, Y) = 0$  с.в.  $X$  и  $Y$  могут оказаться зависимыми. Если  $K(X, Y) \neq 0$ , то случайные величины  $X$  и  $Y$  обязательно зависимы.

Случайные величины  $X$  и  $Y$  называются *некоррелированными*, если выполнено условие  $K(X, Y) = 0$ . В противном случае, при  $K(X, Y) \neq 0$ , случайные величины  $X$  и  $Y$  называют *коррелированными*.

Условие некоррелированности не свидетельствует о независимости случайных величин. Однако для нормально распределенных с.в. условия независимости и некоррелированности совпадают. Поэтому из некоррелированности этих с.в. следует их независимость.

Корреляционный момент может использоваться как характеристика степени зависимости случайных величин: чем больше  $K(X, Y)$  отличается от нуля, тем больше зависимость  $X$  и  $Y$ . Однако значения корреляционных моментов зависят и от соответствующих стандартных отклонений, что может повлиять на корректность сравнений. В связи с вышесказанным вводится другая, более универсальная, характеристика зависимости случайных величин – коэффициент корреляции.

*Коэффициентом корреляции* случайных величин  $X$  и  $Y$  называют число  $r_{xy}$ , определяемое выражением

$$r_{xy} = \frac{K(X, Y)}{\sqrt{D(X)D(Y)}} = \frac{K(X, Y)}{\sigma_x \sigma_y}.$$

---

\* В некоторых изданиях величину, определяемую согласно приведенному выражению, называют только ковариацией, а корреляционные моменты находят не вычитая в этом выражении математических ожиданий.

Модуль числителя этого выражения не может быть больше знаменателя. Поэтому значение коэффициента корреляции изменяется по модулю от 0 до 1. Оно характеризует степень зависимости случайных величин: независимость и некоррелированность ( $r_{xy} = 0$ ), зависимость и некоррелированность ( $r_{xy} = 0$ ), зависимость и коррелированность с.в. ( $r_{xy} \neq 0$ ). Увеличение  $|r_{xy}|$  свидетельствует об увеличении степени зависимости случайных величин, и при  $|r_{xy}| = 1$  имеет место линейная зависимость вида  $Y = aX + b$ .

Таким образом, коэффициент корреляции является показателем степени зависимости случайных величин. При этом, если  $r_{xy} = 1$ , то  $X$  и  $Y$  связаны прямой линейной зависимостью ( $a > 0$ ): с ростом значений  $X$  растут и значения  $Y$ ; если  $r_{xy} = -1$ , то  $X$  и  $Y$  связаны обратной линейной зависимостью ( $a < 0$ ): чем больше значения  $X$ , тем меньше значения  $Y$ .

При  $|r_{xy}| \neq 1$  обычно однозначного соответствия возрастания (убывания) значений одной случайной величины с ростом другой нет. Однако имеется тенденция к возрастанию (убыванию) соответствующих значений. Например, случайные величины рост человека ( $X$ ) и масса его тела ( $Y$ ) характеризуются положительным значением коэффициента корреляции, отличным от единицы. Несмотря на возможные отдельные исключения (человек большего роста обладает меньшей массой), соблюдается общая тенденция: «большому значению роста соответствует большая масса тела». В этом случае наблюдается прямая зависимость, но функциональной, а тем более линейной зависимости, нет. Зависимость такого рода называется корреляционной.

### 2.3.2 Совместные распределения случайных величин

Рассмотрим распределения Пирсона, Стьюдента и Фишера-Снедекора, которые являются распределениями определенного вида функций от многих

случайных переменных. Эти распределения широко используются в статистических исследованиях.

**Распределение Пирсона (распределение  $\chi^2$ ).** Пусть  $X_1, X_2, \dots, X_n$  – система случайных величин, которые независимы и имеют стандартное нормальное распределение.

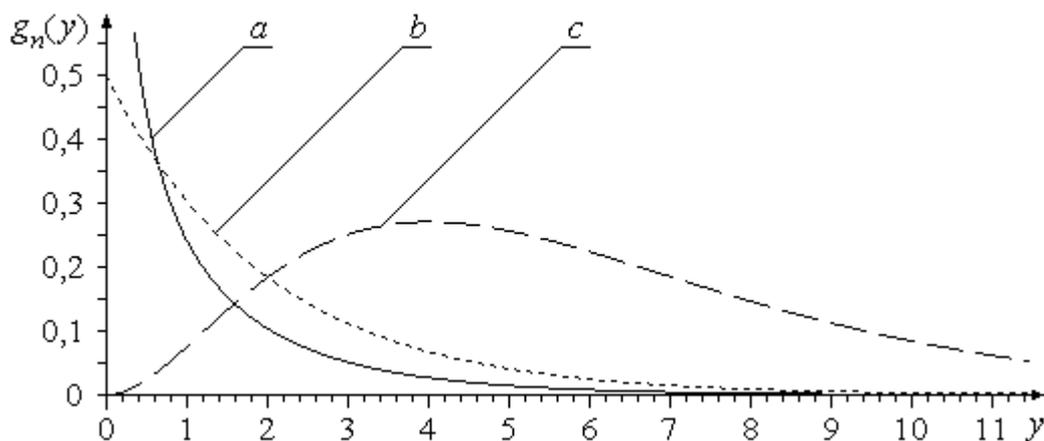
Введем новую случайную величину как функцию вида  $X_1^2 + X_2^2 + \dots + X_n^2$ . Данную сумму квадратов обычно обозначают  $\chi^2$  (читается «хи-квадрат») или с учетом количества слагаемых:

$$\chi_n^2 = X_1^2 + X_2^2 + \dots + X_n^2.$$

Обозначим значения случайной величины  $\chi^2$  через  $y$  (отметим, что  $y \geq 0$ ). Тогда исходя из совместной плотности системы независимых случайных величин (см. соотношения для плотности в разделе 2.3.1), имеющих стандартное нормальное распределение

$$f(x_1, x_2, \dots, x_n) = f_1(x_1) \cdot f_2(x_2) \cdot \dots \cdot f_n(x_n) = \frac{1}{(\sqrt{2\pi})^n} e^{-\frac{1}{2}(x_1^2 + x_2^2 + \dots + x_n^2)},$$

можно найти функцию распределения и плотность  $g_n(y)$  распределения случайной величины  $y = \chi_n^2$ . График функции плотности  $g_n(y)$  при различных  $n$  представлен на рис. 2.9.



**Рис. 2.9.** Графики функции  $g_n(y)$  при  $n=1$  (a),  $n=2$  (b) и  $n=6$  (c)

Распределение случайной величины  $\chi_n^2$  зависит лишь от одного параметра  $n$ , называемого *числом степеней свободы*. При этом  $M(\chi_n^2) = n$ ,  $D(\chi_n^2) = 2n$ .

**Распределение Стьюдента ( $t$ -распределение).** Рассмотрим случайную величину  $X \sim N(0, 1)$  и случайную величину  $\chi_n^2$ . Пусть эти величины независимы. Тогда их совместная плотность равна произведению плотностей  $f(x)$  и  $g_n(y)$ , где  $f(x)$  - плотность стандартного нормального распределения.

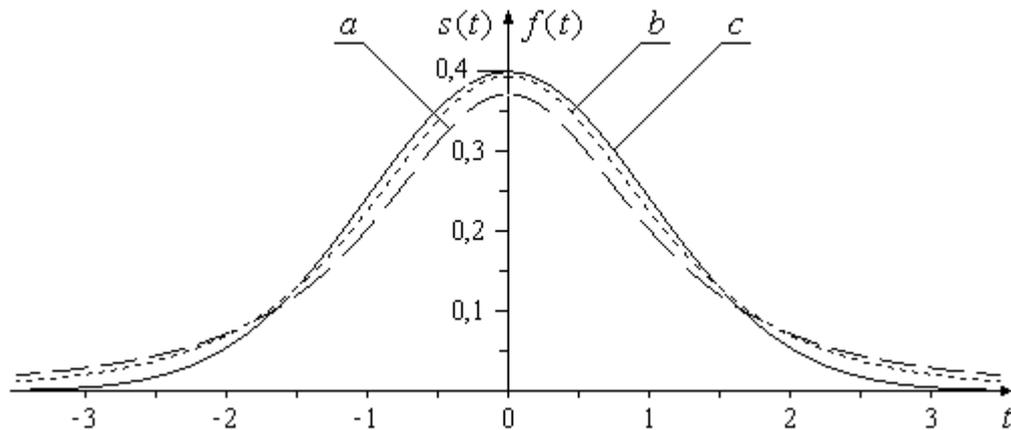
Рассмотрим новую случайную величину

$$T = \sqrt{n} \frac{X}{\sqrt{\chi_n^2}},$$

принимающую значения  $t$ .

Используя совместную плотность распределения случайных величин  $X$  и  $\chi_n^2$ , можно найти функцию распределения  $F(t)$  и плотность  $s_n(t)$  распределения случайной величины  $T$ .

Распределение с.в.  $T$  называется распределением Стьюдента ( $t$ -распределением). Единственный параметр распределения  $n$  называют числом степеней свободы. График функции плотности  $s_n(t)$  симметричен относительно оси ординат (рис 2.10).



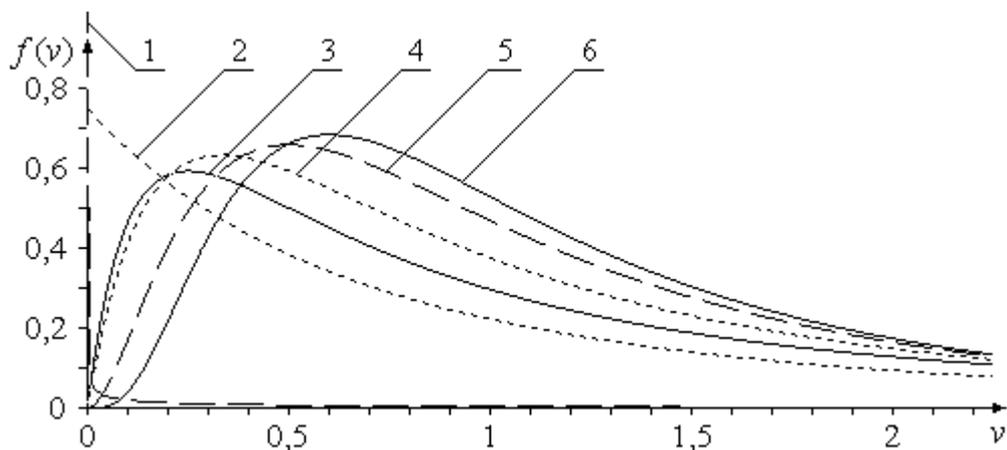
**Рис. 2.10.** Графики плотности распределения Стьюдента в сравнении со стандартной нормальной кривой: а)  $n = 2$ , б)  $n = 4$ , в) нормальная кривая  $N(0, 1)$

Распределение Стьюдента с ростом  $n$  стремится к нормальному распределению с нулевым средним и уже при  $n = 30$  кривая  $s_n(t)$  практически совпадает с соответствующей кривой Гаусса.

**Распределение Фишера-Снедекора.** Данное распределение называют также  $F$ -распределением. Оно является распределением с.в.  $F$ , которая имеет следующую структуру:

$$F = \frac{\chi_m^2 / m}{\chi_n^2 / n} = \frac{n}{m} \cdot \frac{\chi_m^2}{\chi_n^2},$$

где  $\chi_m^2, \chi_n^2$  – независимые с.в., имеющие распределение  $\chi^2$  с соответствующим числом степеней свободы ( $m$  и  $n$  – параметров распределения). Случайная величина  $F$  принимает только положительные значения  $v$ . Графики функции плотности  $f(v)$   $F$ -распределения представлены на рис. 2.11.



**Рис. 2.11.** Графики функции плотности  $F$ -распределения при  $m=1, n=4$  (1),  $m=2, n=4$  (2),  $m=4, n=2$  (3),  $m=n=4$  (4),  $m=n=6$  (5) и  $m=10, n=6$  (6)

### 2.3.3 Условное распределение. Понятие регрессии

Как известно, для условной вероятности случайных событий справедлива формула

$$P(B/A) = \frac{P(AB)}{P(A)}.$$

Для событий, связанных со случайными величинами,  $X = x_i$  или  $Y = y_j$ , приведенная условная вероятность запишется в виде

$$P(Y = y_j | X = x_i) = \frac{P(X = x_i, Y = y_j)}{P(X = x_i)}.$$

Условная вероятность  $P(Y = y_j | X = x_i)$  означает вероятность того, что  $Y = y_j$  при обязательном осуществлении события  $X = x_i$ . Заметим, что фигурирующая в той же формуле вероятность  $P(X = x_i, Y = y_j)$  – это вероятность произведения событий  $X = x_i$  и  $Y = y_j$ . Используя условные вероятности, можно из совместного распределения с.в.  $X, Y$  получить условные законы распределения. Зная условные распределения, по стандартным формулам можно вычислить условное математическое ожидание и условную дисперсию. С условным математическим ожиданием связано одно из важнейших понятий статистики: регрессия.

Регрессия – зависимость математического ожидания какой-либо величины от некоторой другой величины или нескольких величин. Т.е. регрессия в отличие от корреляции характеризует зависимость случайных величин не одним числом, а в виде функции.

*Регрессией* с.в.  $Y$  на с.в.  $X$  называют условное математическое ожидание  $M(Y|X=x)$ . При каждом значении  $x$  величина  $M(Y|X=x)$  принимает свое персональное значение, т.е. регрессия  $M(Y|X=x)$  является функцией от  $x$ . Отметим, что для независимых с.в.  $X$  и  $Y$  регрессия является постоянной:  $M(Y|X=x) = M(Y) = const$ . Для зависимых с.в. регрессия рассматривается в определенном классе функций (линейных, квадратичных, логарифмических и др.).

По аналогии с понятием условной вероятности случайных событий для непрерывных с.в. вводится понятие условной плотности распределения, ис-

пользуя которую, можно по стандартным формулам вычислить условные математическое ожидание и дисперсию\*.

Если регрессия  $Y$  на  $X$  и регрессия  $X$  на  $Y$  суть линейные функции, то зависимость между  $Y$  и  $X$  называют *линейной корреляционной зависимостью*. Можно доказать, что любые коррелированные нормально распределенные с. в.  $X$  и  $Y$  связаны линейной корреляционной зависимостью, причем

$$M(Y|X = x) = m_y + r \frac{\sigma_y}{\sigma_x} (x - m_x),$$

$$M(X|Y = y) = m_x + r \frac{\sigma_x}{\sigma_y} (y - m_y),$$

где  $x, y$  – переменные, а  $m_x, \sigma_x, m_y, \sigma_y, r$  – числовые характеристики распределений.

Если с.в.  $X$  и  $Y$  некоррелированы, т.е. их коэффициент корреляции  $r = 0$ , то регрессия оказывается постоянной:  $M(Y|X = x) = m_y, M(X|Y = y) = m_x$ . В случае произвольного распределения зависимых с.в.  $X$  и  $Y$  они не обязаны иметь именно линейную корреляционную зависимость.

## 2.4. Элементы математической статистики

Математическая статистика занимается методами сбора, обработки, анализа и интерпретации экспериментальных данных. Во многих случаях указанные данные являются параметрами соответствующих математических моделей. Кроме того, алгоритмы математического моделирования нередко основываются на соответствующих положениях математической статистики.

### 2.4.1 Выборочный метод.

Одним из важнейших понятий в статистике является понятие генеральной совокупности (г.с.). *Генеральной совокупностью* называют множество качественно однородных объектов, объединенных по какому-либо признаку

---

\* Подробнее см.: Математическая статистика в медицине: учеб. пособие / В.А. Медик, М.С. Токмачев. – М.: Финансы и статистика, 2007. – С. 191.

или группе признаков. Генеральная совокупность может быть конечной (например, множество граждан России) или бесконечной (например, множество натуральных чисел), может быть реальной или гипотетической (например, множество исходов при неограниченном количестве испытаний).

Любое подмножество объектов г.с. называют *выборочной совокупностью (выборкой)*. Каждый элемент совокупности, обладающий рассматриваемым признаком, является *единицей совокупности* (генеральной или выборочной). Количество элементов совокупности называют ее *объемом*. Целью статистического исследования является получение максимальной информации о г.с., исходя из выборочной совокупности. Формирование выборочной совокупности, “достаточно хорошо” отражающей характеристики генеральной совокупности, является важнейшим этапом статистического исследования.

В математическом смысле генеральная совокупность понимается как случайная величина  $X$  с некоторым распределением вероятностей. Выборка определяется двояко: с одной стороны – это набор случайных величин  $X_1, X_2, \dots, X_n$ , каждая из которых является той же самой с.в.  $X$ , но в соответствующем испытании; с другой стороны, выборка – это набор чисел  $x_1, x_2, \dots, x_n$ , значений (реализаций) случайных величин  $X_1, X_2, \dots, X_n$ .

Тот факт, что выборка  $X_1, X_2, \dots, X_n$  полномочно представляет г.с. в последовательности из  $n$  испытаний, реализуется в следующем определении.

Определение. Выборка  $X_1, X_2, \dots, X_n$  из г.с.  $X$  называется *репрезентативной* (полномочно представительной), если  $X_1, X_2, \dots, X_n$  независимы и распределены так же, как и г.с. Реальное обеспечение репрезентативности выборки гарантируется способом случайного отбора элементов в выборку.

Все суждения о г.с. по выборочным данным справедливы лишь для репрезентативных выборок. Любое нарушение репрезентативности при формировании выборки ставит под сомнение достоверность выводов касательно г.с.

Подлежащие обработке элементы выборки (числа), как правило упорядочивают в виде *вариационного* или *статистического* ряда. Вариационным рядом называют выборку  $x_1, x_2, \dots, x_n$ , все элементы которой упорядочены по возрастанию ( $x_{i+1} \geq x_i$ ). Элементы такого ряда часто называют *вариантами*. Очевидно, в вариационном ряде могут быть совпадающие элементы.

Статистическим рядом называют множество элементов выборки с сопоставленными им числами  $m_i$  их появления в выборке:

$$\begin{array}{cccc} x_1 & x_2 & \cdots & x_k \\ m_1 & m_2 & \dots & m_k \end{array}$$

При этом все числа  $x_1, x_2, \dots, x_k$  различны и упорядочены по возрастанию,  $m_1 + m_2 + \dots + m_n = n$ .

Задачей, наиболее часто решаемой на основе выборок, является *вычисление средних*, которые являются обобщёнными характеристиками рассматриваемых показателей. Например, инкубационный период вирусной инфекции для каждого отдельного больного является величиной индивидуальной, зависящей от защитных свойств организма, от внешних условий. А среднее значение, вычисленное на основании многих наблюдений, оказывается свободным от случайных индивидуальных факторов и в значительной мере характеризует контагиозность инфекционного процесса и вирулентность конкретного типа вируса.

Однако средние, найденные по выборке, отличаются от средних генеральной совокупности по объективным причинам: недостаточному числу наблюдений, ошибкам репрезентативности при формировании выборки, ошибкам регистрации. Т.е. средние значения, вычисляемые по выборке, являются лишь оценками соответствующих характеристик генеральной совокупности. При большом разбросе выбранных значений рассматриваемого параметра (при большой дисперсии этих значений) среднее становится малоинформативной величиной. Такой величиной является, например, средняя зарплата в стране, где различия в численных значениях этого показателя в десятки раз.

В подобных случаях необходима группировка данных и соответственно вычисление групповых средних, как более осмысленных характеристик.

Среднее значение чаще всего вычисляют как среднее арифметическое. Но если в полученной объёмом  $n$  имеются повторяющиеся значения, что характерно для выборок дискретных случайных величин, то при вычислении среднего вначале можно определить числа  $m_1, m_2, \dots, m_k$  появления значений  $x_1, x_2, \dots, x_k$  рассматриваемой с.в.  $X$ . Тогда для среднего  $\bar{x}$  получаем три эквивалентных выражения:

$$\bar{x} = \frac{1}{n} \sum_{i=0}^n x_i = \frac{1}{n} \sum_{i=0}^k m_i x_i = \sum_{i=0}^k \omega_i x_i,$$

где  $m_1 + m_2 + \dots + m_k = n$ , а  $\omega_i = m_i/n$  – *относительная частота* или *статистическая вероятность* появления значения  $x_i$  в последовательности значений с.в.  $X$ .

Можно доказать, что при неограниченном увеличении числа испытаний и при репрезентативности выборки некоторого показателя относительные частоты появления его значений как угодно близко приближаются к вероятностям появления этих значений, а среднее значение этого показателя как угодно близко приближается к его математическому ожиданию.

**Выборочное распределение.** Пусть из г.с.  $X$  в результате  $n$  независимых испытаний (наблюдений) получена некоторая выборка объёма  $n$ . Среди  $n$  значений  $x_i$  могут быть и равные. Сами числа в выборке, вообще говоря, располагаются в порядке их появления.

*Статистическим распределением* (распределением выборки) случайной величины  $X$  называют последовательность её значений  $x_1, x_2, \dots, x_k$ , расположенных в возрастающем порядке с указанием относительных частот  $\omega_1, \omega_2, \dots, \omega_k$ , с которыми они содержатся в выборке:

$$\begin{array}{cccc} x_1 & x_2 & \dots & x_k \\ \omega_1 & \omega_2 & \dots & \omega_k \end{array}$$

Статистическое распределение является своего рода приближением теоретического (истинного) распределения г.с.  $X$ . Сумма всех  $\omega_i$  равна единице. По своей структуре статистическое распределение соответствует распределению дискретной с.в. Для его описания статистического используются те же характеристики, что и для теоретического распределения, но с учетом, что они относятся к выборочному распределению. Выборочные характеристики, определяемые по статистическому распределению, являются аналогами, оценками аналогичных характеристик генеральной совокупности. Для числовых характеристик выборки можно использовать те же выражения, что и для теоретического распределения дискретных с.в., если в этих выражениях вместо вероятностей  $p$  использовать соответствующие относительные частоты. Так, дисперсию выборки с.в.  $X$  находят согласно выражению

$$\bar{\mu}_2(X) = \bar{D}_{\text{выб}}(X) = \frac{1}{n} \sum_{i=1}^k (x_i - \bar{x})^2 = \sum_{i=1}^k \omega_i (x_i - \bar{x})^2.$$

Отметим, что г.с.  $X$  может иметь непрерывное распределение, распределение же выборки всегда дискретно. К характеристикам выборки принято добавлять прилагательные «выборочная», «эмпирическая» или «статистическая», в отличие от соответствующих теоретических (генеральных) характеристик. Выборочные характеристики обозначаются теми же буквами, что и теоретические, но с чертой сверху или с индексом «выб» снизу. Иногда используют оба эти обозначения. Например:  $M(\bar{X})$ ,  $\sigma(\bar{X})$ ,  $\bar{\sigma}(X)$ ,  $\bar{\sigma}(X)_{\text{выб}}$  и т.д.

*Эмпирической функцией распределения* выборки объема  $n$  называют функцию  $\hat{F}_n(X)$ , значения которой определяются выражением

$$\hat{F}_n(x) = \begin{cases} 0 & \text{при } x < x_1, \\ \frac{1}{n} \sum_{i=1}^k n_i = \sum_{i=1}^k \omega_i & \text{при } x_1 \leq x < x_{k+1} \leq x_n, \\ 1 & \text{при } x \geq x_n. \end{cases}$$

График этой функции имеет ступенчатый вид, аналогичный графику теоретической  $F(x)$  на рис. 2.3.

#### 2.4.2. Оценки параметров распределения. Доверительные интервалы

**Понятие оценки. Точечные оценки.** Одной из основных задач статистического исследования является нахождение параметров распределения (а также иных характеристик) генеральной совокупности  $X$ , исходя из выборочной совокупности. Если известен тип распределения г.с.  $X$ , то с нахождением параметров распределения г.с.  $X$  оказывается полностью определенной. Однако по выборке невозможно определить точные значения параметров генеральной совокупности. Выборка может дать лишь приближенные, с различной степенью точности значения, называемые *оценками параметров*. Рассмотрим виды оценок, их качество, методы построения. При этом в характеристике оценки будем исходить не из конкретных чисел, которые меняются от выборки к выборке, а из ее структуры, из вида и свойств соотношения, применяемого для построения оценки по выборке.

Пусть  $X_1, X_2, \dots, X_n$  – репрезентативная выборка из генеральной совокупности  $X$ . Согласно определению репрезентативности последовательность  $X_1, X_2, \dots, X_n$  является системой независимых случайных величин, распределенных так же, как и с. в.  $X$ . Любую функцию или функционал от выборочных значений  $X_1, X_2, \dots, X_n$  принято называть *статистикой*.

Обозначим точное значение параметра генеральной совокупности  $X$  посредством  $\Theta$ . Таким образом,  $\Theta$  – число, которое, вообще говоря, исследователю неизвестно. Любая статистика выборочных значений, приближенно равная оцениваемому параметру  $\Theta$ , т.е. представляющая собой одно число, называется точечной оценкой (или просто оценкой) параметра  $\Theta$  и может быть обозначена как  $\bar{\Theta}$  или как  $\bar{\Theta}(X_1, X_2, \dots, X_n)$ . Значение  $\bar{\Theta}$  зависит от  $n$  случайных величин и, естественно, от вида оцениваемого параметра распределения г.с.

Приведём выражения для определения оценок основных параметров исследуемых распределений с помощью выборок объёма  $n$ :

$$M(\bar{X}) = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i, \quad D(\bar{X}) = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2, \quad \sigma(\bar{X}) = \sqrt{D(\bar{X})},$$

$$\bar{R}(X, Y) = \frac{\frac{1}{n-1} \sum_{i=1}^n X_i Y_i - \bar{X} \bar{Y}}{\sigma(X) \sigma(Y)}.$$

Важными характеристиками оценок являются *состоятельность*, *несмещённость* и *эффективность*. Все приведённые выше статистики обеспечивают состоятельность и несмещённость соответствующих оценок.

Оценка  $\bar{\Theta}_n$ , составленная по репрезентативной выборке  $X_1, X_2, \dots, X_n$  из г.с.  $X$  называется *состоятельной оценкой* параметра  $\Theta$  г.с., если эта оценка с ростом  $n$  сходится по вероятности к оцениваемому параметру, т.е. при любом положительном числе  $\varepsilon$   $\lim_{n \rightarrow \infty} P(|\bar{\Theta}_n - \Theta| < \varepsilon) = 1$ . При такой оценке с ростом  $n$  значения  $\bar{\Theta}_n$  имеют тенденцию концентрироваться около параметра  $\Theta$ . Следовательно, разброс значений  $\bar{\Theta}_n$  уменьшается, а  $\lim_{n \rightarrow \infty} D(\bar{\Theta}_n) = 0$ .

Оценка  $\bar{\Theta}$ , составленная по репрезентативной выборке  $X_1, X_2, \dots, X_n$  из г.с.  $X$  называется *несмещенной оценкой* параметра  $\Theta$  г.с., если  $M(\bar{\Theta}) = \Theta$ . Несмещенность оценки означает, что при использовании вместе  $\Theta$  его оценки  $\bar{\Theta}$  не допускается систематическая ошибка, т.е. ошибка одного знака в сторону завышения или занижения истинного значения. Несмещенность оценки особенно важна при малых выборках.

Оценка  $\bar{\Theta}$ , составленная по репрезентативной выборке  $X_1, X_2, \dots, X_n$  называется *эффективной оценкой* параметра  $\Theta$  генеральной совокупности  $X$ , если она обладает наименьшей дисперсией по сравнению с дисперсиями других оценок параметра  $\Theta$ , составленных по той же выборке. Эффективность оценки  $\bar{\Theta}$  означает, что при одном и том же объеме выборки  $n$  эффективная оценка обладает наименьшим разбросом от выборки к выборке своих возможных значений относительно  $\Theta$ . Так, для математического ожидания нор-

мально распределённой с.в.  $X$  приведённая выше статистика является эффективной. Однако в целом далеко не всегда возможно подобрать оценку, наилучшую во всех отношениях, поэтому допустимо использование оценок с малыми смещениями или оценок, не являющихся эффективными.

**Интервальные оценки.** Наряду с точечными оценками, недостатком которых является подверженность их случайным колебаниям, особенно при малых выборках, в статистике широкое применение нашли *интервальные оценки* параметров распределения г.с.

Интервальные оценки задают двумя выборочными значениями: границами интервала, в котором оказывается исследуемый параметр  $\Theta$  с некоторой заранее оговоренной вероятностью (*надежностью*). Обозначая эти границы интервала  $\bar{\Theta}_1$  и  $\bar{\Theta}_2$ , получаем выражение

$$P(\bar{\Theta}_1 < \Theta < \bar{\Theta}_2) = \gamma,$$

где вероятность  $\gamma$ , называемая *доверительной вероятностью* или *надежностью* выбирается исследователем. Значения  $\gamma$  обычно выбираются близкими к 1. Стандартными считаются числа 0,9, 0,95 и 0,99.

Интервал  $(\bar{\Theta}_1, \bar{\Theta}_2)$ , соответствующий выбранной доверительной вероятности  $\gamma$  называют *доверительным интервалом*, а  $\bar{\Theta}_1$ ,  $\bar{\Theta}_2$  - доверительными границами.

Из смысла введенных понятий следует, что оценка в виде доверительного интервала характеризуется двумя величинами: шириной доверительного интервала (чем меньше интервал, тем лучше оценка) и надежностью (чем больше значение  $\gamma$ , тем достовернее оценка). Заметим, что надежность и размеры доверительного интервала взаимосвязаны. Улучшая одну из этих характеристик при прочих одинаковых условиях, мы неминуемо ухудшим другую (например, полагая большее значение надежности, получим как следствие увеличение доверительного интервала). Одновременно улучшить обе этих качественных составляющих оценки возможно лишь при увеличении объема выборки, что вполне логично.

Доверительный интервал для параметра  $\Theta$ , найденный по выборке, называют интервалом значений  $\Theta$ , не противоречащих опытными данным. Обычно он выбирается симметричным относительно соответствующей точечной оценки. В этом случае обе доверительных границы характеризуются лишь одним числом  $\Delta$  – расстоянием от центра интервала  $\bar{\Theta}$ , т.е.  $\bar{\Theta}_1 = \bar{\Theta} - \Delta$ ,  $\bar{\Theta}_2 = \bar{\Theta} + \Delta$ .

Таким образом, для нахождения доверительного интервала некоторого параметра  $\Theta$  г.с.  $X$  необходимо по выборке  $X_1, X_2, \dots, X_n$  вычислить соответствующую точечную оценку  $\bar{\Theta}$  и по назначаемому исследователем значению  $\gamma$  определить величину  $\Delta$ . Полученное значение  $\Delta$  является точностью интервальной оценки.

Значение  $\Delta$ , соответствующее выбранной вероятности  $\gamma$ , зависит от закона распределения рассматриваемой с.в.  $X$ . Рассмотрим случаи определения доверительных интервалов для  $M(X)$  и  $\sigma(X)$  нормально распределённой г.с.  $X$ .

**Оценка математического ожидания.** Можно доказать, что при известной дисперсии  $\sigma^2$  математическое ожидание  $M(X) \equiv \Theta$ , оценивается с помощью двойного неравенства

$$\bar{X} - \frac{\sigma}{\sqrt{n}} \varepsilon_\gamma < \Theta < \bar{X} + \frac{\sigma}{\sqrt{n}} \varepsilon_\gamma,$$

где  $n$  – объем выборки, а коэффициент  $\varepsilon_\gamma$  определяется по таблице функции Лапласа (табл. П 1) для  $\sigma = 1$  при заданной надежности  $\gamma$ :

$$\gamma = \int_0^{\varepsilon_\gamma} e^{-\frac{z^2}{2}} dz, \quad \Phi(\varepsilon_\gamma) = \frac{\gamma + 1}{2}.$$

Величину  $\varepsilon_\gamma$  называют *коэффициентом доверия* (табл. 2.1). Она зависит только от  $\gamma$ . Умножая её на  $\sigma/\sqrt{n}$ , получают значение  $\Delta$ , равное половине доверительного интервала и соответствующее известному значению и выбранному объёму выборки  $n$ .

Таблица 2.1. Значения коэффициентов доверия

$\gamma$	$\varepsilon_\gamma$	$\gamma$	$\varepsilon_\gamma$	$\gamma$	$\varepsilon_\gamma$	$\gamma$	$\varepsilon_\gamma$
0,6826	1	0,85	1,439	0,91	1,694	0,96	2,053
0,80	1,282	0,86	1,475	0,92	1,750	0,97	2,169
0,81	1,310	0,87	1,513	0,93	1,810	0,98	2,325
0,82	1,340	0,88	1,554	0,94	1,880	0,99	2,576
0,83	1,371	0,89	1,597	0,95	1,960	0,9973	3
0,84	1,404	0,90	1,643	0,9544	2	0,999	3,290

Если значение  $\sigma$  не известно, то вначале оно оценивается по выборке с помощью точечной оценки  $S$ , где

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2, \quad S = \sqrt{S^2}.$$

Затем  $M(X) \equiv \Theta$  оценивается двойным неравенством

$$\bar{X} - \frac{S}{\sqrt{n}} t_\gamma < \Theta < \bar{X} + \frac{S}{\sqrt{n}} t_\gamma,$$

где  $t_\gamma$  – коэффициент доверия. Однако вычисляется  $t_\gamma$  иначе, чем в предыдущем случае: с привлечением распределения Стюдента. Значение  $t_\gamma$  при заданных  $n$  и  $\gamma$  можно вычислить, исходя из соотношения  $2 \int_0^{t_\gamma} s_{n-1}(t) dt = \gamma$ . На практике обычно это значение находят из таблиц (табл. П 2). Указанный доверительный интервал, как правило, используют в случае выборок малого объема. С увеличением  $n$  коэффициент доверия  $t_\gamma$  стремится к соответствующему значению  $\varepsilon_\gamma$ , а выборочное стандартное отклонение  $s$  – к значению генерального  $\sigma$ . Таким образом, доверительные интервалы для выборок большого объема, получаемые как при известной, так и при неизвестной дисперсии, становятся малоразличимыми.

**Оценка дисперсии.** Для нормально распределенной г.с.  $X$  доверительный интервал дисперсии  $D(X)$  при достаточно больших  $n$  (приемлемые значения, как правило, получаются уже для  $n > 30$  можно найти по формуле

$$S^2 - \sqrt{\frac{2}{n-1}} S^2 \varepsilon_\gamma < D(X) < S^2 + \sqrt{\frac{2}{n-1}} S^2 \varepsilon_\gamma,$$

где  $S^2$  - исправленная выборочная дисперсия, а коэффициент доверия  $\varepsilon_\gamma$ , находится по табл. 2.1.

Аналогично для равномерно распределенной г.с.  $X$  интервальная оценка  $D(X)$  выглядит следующим образом:

$$S^2 - \sqrt{\frac{0,8n+1,2}{n(n-1)}} S^2 \varepsilon_\gamma < D(X) < S^2 + \sqrt{\frac{0,8n+1,2}{n(n-1)}} S^2 \varepsilon_\gamma,$$

где значение  $S^2$  вычисляется, а  $\varepsilon_\gamma$ , при заданном  $\gamma$  определяется из табл. 2.1.

**Оценка вероятности по относительной частоте.** Пусть испытания проводятся по схеме Бернулли относительно некоторого события  $A$ , которое в каждом испытании может произойти с вероятностью  $p$  или не произойти с вероятностью  $q = 1 - p$ . Вероятность  $p$  и является оцениваемым параметром.

Рассматриваемый случай является частным случаем оценивания математического ожидания с.в.  $X$  (количество осуществлений события  $A$  в одном испытании), принимающей лишь два возможных значения, 0 и 1, с соответствующими вероятностями  $q$  и  $p$  (распределение Бернулли). Полагаем  $X$  в качестве генеральной совокупности. Тогда точечная оценка параметра  $p$ , определяемая по выборке  $X_1, X_2, \dots, X_n$ , равна:

$$\bar{p} = \omega = \frac{1}{n} \sum_{i=1}^n X_i.$$

Так как в схеме Бернулли дисперсия равна  $pq$ , то за её точечную оценку принимаем величину  $\omega(1 - \omega)$ . Подставляя эту величину в выражение для интервальной оценки математического ожидания, получаем:

$$\omega - \sqrt{\frac{\omega(1-\omega)}{n}} \varepsilon_\gamma < p < \omega + \sqrt{\frac{\omega(1-\omega)}{n}} \varepsilon_\gamma.$$

Значение коэффициента доверия  $\varepsilon_\gamma$  по-прежнему определяется с помощью табл. 2.1 по выбранной доверительной вероятности  $\gamma$ .

Отметим, что в рассмотренной методике вместо биномиального распределения для г.с.  $X$  использовано нормальное распределение, к которому стремится биномиальное распределение при увеличении объёма выборки  $n$ . Поэтому для использования найденных соотношений значения  $np$  и  $nq$  должны быть не менее 10.

### 2.4.3 Проверка статистических гипотез

*Статистической гипотезой* называют предположение о неизвестном законе распределения генеральной совокупности либо о параметрах известных распределений.

Проверка статистической гипотезы осуществляется статистическими методами, исходя из выборочных данных. К статистической проверке гипотез сводятся задачи сравнительной проверки и оценки различных процессов: эффективности лечения, продолжительности болезни и восстановительного периода, степени тяжести заболевания, сравнение различных характеристик процесса, сравнение групп больных, и т.д.

Статистические гипотезы, не использующие допущений о конкретном законе распределения, называют *непараметрическими*. А гипотезы о параметрах при известном распределении – *параметрическими*.

Основную гипотезу, подлежащую проверке, принято называть *нулевой гипотезой*. Её обычно обозначают  $H_0$ .

Для проверки нулевой гипотезы выдвигается гипотеза альтернативная, противоречащая нулевой. Такую гипотезу называют *конкурирующей гипотезой* и обозначают  $H_1$ .

Из двух взаимоисключающих гипотез объективно справедлива лишь одна. В результате статистического исследования принимается решение о

принятии той или иной гипотезы. Это решение может быть либо верным, либо неверным. Во втором случае возможны два вида ошибок: первого и второго рода.

*Ошибка первого рода* заключается в том, что верная нулевая гипотеза  $H_0$  отвергается, а принимается конкурирующая ложная гипотеза  $H_1$ .

*Ошибка второго рода* заключается в том, что ложная гипотеза  $H_0$  принимается, хотя на самом деле верна конкурирующая гипотеза  $H_1$ .

Отметим, что гипотезы  $H_0$  и  $H_1$ , в исследовании не равноправны. Статистическая проверка осуществляется только для нулевой гипотезы  $H_0$ , поэтому её и называют основной. Проверить нулевую гипотезу необходимо так, чтобы возможности ошибок обоих типов свести к минимуму.

Вероятность – допустить ошибку первого рода обозначим  $\alpha$ . Число  $\alpha$  называют *уровнем значимости*. Аналогично, вероятность допустить ошибку второго рода обозначим  $\beta$ .

Вероятность не допустить ошибку второго рода, т.е. при справедливости конкурирующей гипотезы  $H_1$ , вероятность принять эту гипотезу, называют *мощностью критерия* (иногда говорят "*чувствительность критерия*"). Мощность критерия равна  $1-\beta$ . Чем больше это значение, тем лучше, качественнее работает используемый критерий проверки исследуемой гипотезы.

Задача исследователя – минимизировать обе вероятности: и  $\alpha$ , и  $\beta$ . Но обе вероятности оказываются взаимосвязанными, и, уменьшая одну из них при фиксированных условиях, мы неизбежно это уменьшение компенсируем ростом вероятности другой ошибки. Единственный способ одновременного уменьшения вероятностей обеих ошибок – это увеличение объема выборки.

Обычно поступают следующим образом: фиксируют уровень значимости  $\alpha$ , т.е. задают границу вероятности отклонить нулевую гипотезу  $H_0$ , когда она верна, и пытаются провести исследование так, чтобы значение  $\beta$  оказалось наименьшим. Стандартными уровнями значимости  $\alpha$ , для которых по-

строены соответствующие таблицы, считаются числа 0,2; 0,1; 0,05; 0,02; 0,01; 0,005; 0,002; 0,001.

Разумное соотношение между  $\alpha$  и  $\beta$  находят, исходя из тяжести последствий каждой из ошибок в рассматриваемой задаче.

Принятие или отбрасывание нулевой гипотезы происходит с определенной вероятностью и не является логическим доказательством ее справедливости или ложности. Например, принятие нулевой гипотезы  $H_0$  при уровне значимости  $\alpha = 0,01$  означает, что наша гипотеза не противоречит наблюдаемым опытным выборочным данным и для других выборок того же объема из той же генеральной совокупности принятая гипотеза в среднем будет справедлива в 99 случаях из 100 (точнее, в среднем не менее чем в 99 случаях на 100 рассматриваемых). Таким образом, принимая эту нулевую гипотезу, мы рискуем ошибиться в среднем не чаще чем в 1% случаев.

**Статистический критерий.** Статистические гипотезы требуют проверки. Для проверки согласованности теории с опытными данными используется частотная интерпретация, согласно которой случайное событие, имеющее малую вероятность, в длинной последовательности испытаний будет осуществляться достаточно редко. Можно полагать практически несомненным, что в единичном испытании это редкое событие не произойдет. И наоборот, событие с вероятностью осуществления, близкой к 1, в единичном испытании обязательно должно произойти, т.е. такое событие является практически достоверным. На этих принципах и основана теория проверки статистических гипотез. Теория допускает возможность появления ошибок и предлагает методы, минимизирующие вероятности появления ошибок.

Рассмотрим стандартную схему проверки гипотез. Нулевая гипотеза обычно связывается либо с каким-то распределением г.с.  $X$ , либо с параметром при уже известном законе распределения. Для проверки нулевой гипотезы  $H_0$  вводится некоторая случайная величина, называемая *статистическим критерием*. При справедливости  $H_0$  статистический критерий оказывается полностью определенным, с известным распределением и параметрами.

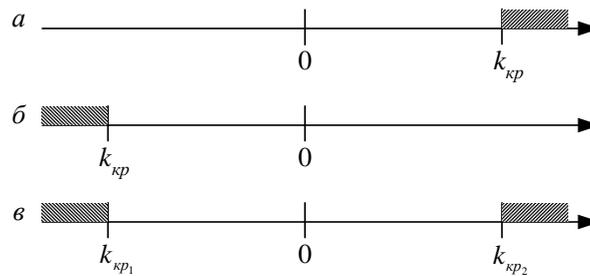
Обычно критерий выбирают таким, чтобы он имел одно из следующих стандартных распределений: нормальное,  $\chi^2$ , распределение Стьюдента, распределение Фишера.

Выбрав критерий  $K$ , все множество значений, принимаемых с.в.  $X$ , разбиваем на два подмножества. Значения из первого подмножества принимаются  $K$  с вероятностью, близкой к 1, т.е. это практически достоверное событие при предполагаемом теоретическом распределении  $X$ , соответствующим гипотезе  $H_0$ . Значения из второго подмножества принимаются с вероятностью, близкой к 0, т.е. это маловероятное событие является практически невозможным опять же при теоретическом распределении  $K$ , соответствующим гипотезе  $H_0$ .

Первое из указанных множеств значений называют *областью принятия гипотезы*. Если вычисляемое по выборке значение критерия  $K$  попадает в эту область, то принимается гипотеза  $H_0$ , как не противоречащая опытным данным.

Второе из обозначенных множеств называют *критической областью*. Если вычисленное по выборке значение  $K$  попадает в критическую область, т.е. происходит событие практически невозможное при справедливости  $H_0$ , то появляется повод усомниться именно в нулевой гипотезе. В этом случае принимается конкурирующая гипотеза  $H_1$ .

Поскольку все возможные значения критерия образуют интервал, то область принятия гипотезы и критическая область полагаются в виде непересекающихся интервалов. Виды стандартных критических областей (правосторонняя, левосторонняя и двусторонняя) представлены на рис. 2.12 (заштрихованные интервалы).



**Рис. 2.12.** Виды критических областей: а) правосторонняя; б) левосторонняя; в) двусторонняя

Граничные точки, отделяющие критическую область от области принятия гипотезы, называют *критическими точками*. Критические точки полностью определяют указанные области. Вероятность попадания критерия  $K$  в критическую область при справедливости  $H_0$  – достаточно малое число. Данная вероятность называется *уровнем значимости* и обозначается  $\alpha$ . Значение  $\alpha$  – вероятность ошибки первого рода. Эта величина выбирается исследователем. В случае двусторонней критической области вероятность попадания в каждый из интервалов критической области полагают равной  $\alpha/2$ .

Итак, в результате исследования наблюдаемое значение критерия  $K$  оказывается либо в области принятия гипотезы, либо в критической области. В первом случае принимается гипотеза  $H_0$ , как не противоречащая опытным данным. Различия между наблюдаемыми значениями и истинными обусловлены случайными причинами и признаются *не значимыми* (не принципиальными). Во втором случае нулевая гипотеза отвергается, принимается гипотеза  $H_1$ . Различия между наблюдаемыми значениями и теоретическими (согласно нулевой гипотезе) значимы, т.е. обусловлены принципиальными причинами: ошибочностью нулевой гипотезы.

Критерии общего характера проверки статистических гипотез называют *критериями значимости*. Критерии проверки соответствия выборочного и теоретического распределений принято называть *критериями согласия*.

## ГЛАВА 3. МОДЕЛИРОВАНИЕ ПОКАЗАТЕЛЕЙ ЗДОРОВЬЯ НАСЛЕНИЯ

### 3.1. Математические модели в здравоохранении и их характеристики

В задачах моделирования динамики развития и лечения различных болезней, анализа точности алгоритмов прогнозирования здоровья населения, моделирования медико-демографических процессов и др. возникает необходимость в моделировании показателей здоровья населения (ПЗ) и работы медицинских учреждений. Для решения данной задачи следует выяснить характер указанных распределений, возможности аппроксимации их классическими распределениями, а в случае необходимости подобрать и подходящие частные распределения, удобные для моделирования соответствующих случайных величин. Перечисленные вопросы и являются предметом рассмотрения в настоящей главе.

Математические модели в здравоохранении в основном характеризуются числом параметров (одно-, двух- и многопараметрические), отсутствием или наличием случайностей в алгоритмах моделирования (детерминированные или стохастические) и постоянством или изменением во времени моделируемых величин (статические и динамические). Однопараметрические модели обычно являются наиболее простыми. Они могут входить в многопараметрические модели, т.е. являться как бы программными блоками многопараметрических моделей. Это имеет место, например, в моделях прогнозирования интегральных показателей (индикаторов) здоровья населения, в состав которых входят однопараметрические модели статистических показателей здоровья, рассматриваемых в данной главе.

Если при использовании многопараметрической модели фиксируются значения всех её параметров, кроме одного, то условно модель можно считать и однопараметрической. Вообще можно зафиксировать от одного до

$n-1$  параметров  $n$ -параметрической модели, то получим модель, в которой при исследовании варьируются только от одного до  $n-1$ -параметров.

Учитывая, что одно и то же явление часто можно представить моделями с разным числом параметров проиллюстрируем методику разработки двухпараметрической и однопараметрической моделей, моделирования и интерпретации результатов прогнозирования среднего времени  $T$  пребывания больного на больничной койке в 2007-м году в Новгородской области, предполагая, что прогнозирование проводилось по значениям временного ряда средних значений  $T$  в 2002 ÷ 2006 годах. Сравним полученные модели и прогнозируемое время с фактически полученным (алгоритмы прогнозирования рассматриваются в 6-й главе). Воспользуемся приведёнными в табл. 3.1 средними статистическими значениями времени пребывания больного на больничной койке для Российской Федерации, Северо-Западного федерального округа и Новгородской области в 2000 – 2007 годах.

Т а б л и ц а 3.1. Среднее время пребывания больного на больничной койке (в днях).

Год	2000	2001	2002	2003	2004	2005	2006	2007
Российская Федерация	15,40	15,20	14,80	14,50	14,40	13,80	13,60	13,10
Северо-Западный ФО	15,60	15,40	14,70	14,30	14,10	13,90	13,70	13,30
Новгородская область	15,32	15,16	15,10	14,64	14,21	13,79	13,72	13,69

Построим вначале двухпараметрическую модель. Поскольку значения  $T$  времени пребывания больного на больничной койке не могут быть отрицательными и меньшими некоторого положительного значения  $T_0$ , раньше которого больного не выписывают (обычно один день), а распределение значений  $T$ , по-видимому, имеет колоколообразный график на конечном отрезке, т.е. график, имеющий один максимум и монотонно спадающий до нуля в обе стороны от местоположения максимума, то попытаемся в качестве классического моделирующего распределения взять распределение Вейбулла [107,

136]. Функция плотности  $f(T)$  и функция распределения  $F(T)$  этого распределения при  $T > T_0$  имеют вид:

$$f(T) = \frac{m(T-T_0)^{m-1}}{Q} e^{-\frac{(T-T_0)^m}{Q}}, \quad F(T) = 1 - e^{-\frac{(T-T_0)^m}{Q}}, \quad (3.1)$$

где  $m$  и  $Q$  – параметры распределения. При этом математическое ожидание  $M(T)$ , среднее квадратическое отклонение  $T$  от  $M(T)$  и  $Q$  связаны выражениями

$$M(T) = \Gamma(1/m+1)Q^{1/m}, \quad \sigma^2(T) = (\Gamma(2/m+1) - \Gamma^2(1/m+1))Q^{2/m} \quad \text{и}$$

$$Q = \left( \frac{M^2(T) + \sigma^2(T)}{\Gamma^m(2/m+1)} \right)^{2/m},$$

в которых используется гамма-функция  $\Gamma(z)$  от параметра  $z$  [26, 29, 138]. При  $T \leq T_0$   $f(T) = F(T) = 0$ .

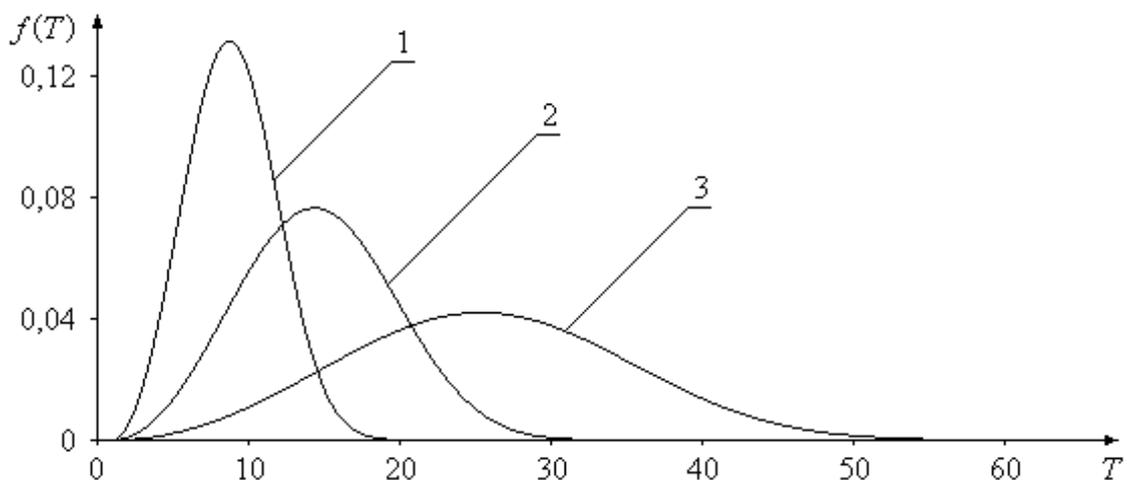
В медицине распределение Вейбулла используется, например, для моделирования процессов, имеющих один экстремум, до которого моделируемая величина монотонно возрастает от нулевого значения, а после него – монотонно убывает, стремясь к нулевому значению [61]. Для рассматриваемого примера в соответствии со статистическими данными наиболее подходящим значением параметра  $m$  является 3. В этом случае  $\Gamma(1/3+1) = 0,8938$ , а  $\Gamma(2/3+1) = 0,9028$ . Выражая  $Q$  через  $M(T)$  и подставляя найденную зависимость в выражения (3.1), для  $T > T_0$  получаем:

$$f(T) = \frac{2,1421(T-T_0)^2}{M^3(T)} e^{-\frac{(T-T_0)^3}{M^3(T)}}, \quad F(T) = 1 - e^{-\frac{0,7140(T-T_0)^3}{M^3(T)}}. \quad (3.2)$$

Выражение (3.2) задаёт двухпараметрическую модель с фиксированным значением одного параметра ( $m = 3$ ).

Согласно статистическим данным времени пребывания больного на больничной койке для Новгородской области в 2000 ÷ 2006 годах было получено:  $M(T) = 14,56$ . С помощью линейного алгоритма прогнозирования (гл. 6) находим:  $\bar{M}(2007) = 13,65$  дня. Фактически же было 13,69 дня. Прогноз значения среднего квадратического отклонения:  $\bar{\sigma}(2007) = 0,36 \cdot 13,65 = 4,91$ .

На рис. 3.1 приведены графики функции плотности (3.2) при  $M(T) = 13,65$  и при двух других значениях  $M(T)$ . Второй график – график прогнозируемой  $f(2007)$ .

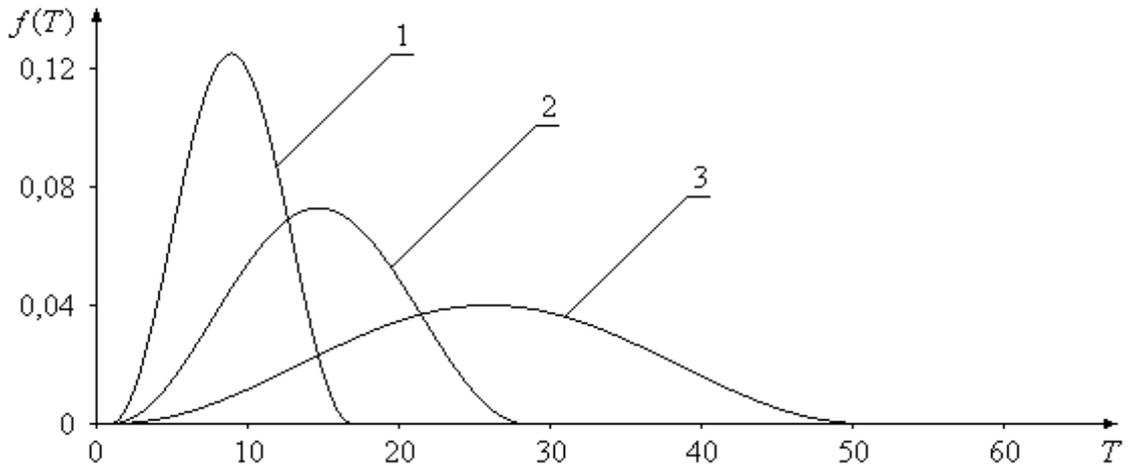


**Рис. 3.1.** Распределения (3.2) времени пребывания больного на больничной койке (в днях) при  $M(T) = 8$  (1),  $M(T) = 13,65$  (3) и  $M(T) = 13,65$  (2)

В качестве однопараметрического распределения можно воспользоваться, например, синусоидальным распределением [54] с параметром  $a = 2M(T) > 0$ .

$$f(T) = \begin{cases} \frac{0,5}{M(T)} \left( 1 + \sin \left( \frac{\pi T}{M(T)} - \frac{\pi}{2} \right) \right) & \text{при } T \in (T_0, T_0 + 2M(T)), \\ 0 & \text{в противном случае.} \end{cases} \quad (3.3)$$

Рис. 3.2 иллюстрирует графики распределения (3.3) при тех же значениях  $M(T)$ , что и на рис. 3.1. Графики на этих рисунках при одинаковых  $M(T)$  мало отличаются друг от друга.



**Рис. 3.2.** Распределения (3.3) времени пребывания больного на больничной койке (в днях) для Новгородской обл. при  $M(T) = 8$  (1),  $M(T) = 13,74$  (2) и  $M(T) = 25$  (3)

Обычно при моделировании одного и того же явления адекватность модели и оригинала улучшатся с увеличением числа характеристик оригинала, которые удаётся учесть в модели. Однако в рассмотренном примере сложно сделать вывод о том, что модель Вейбулла лучше. Дело в том, что в обеих моделях не учитывается важная характеристика моделируемого явления – статистическая оценка среднего квадратического отклонения значения  $T$ . Если бы значение параметра  $m$  в этой модели можно было бы выбрать с учётом указанной оценки, то, по-видимому, данная модель улучшила бы адекватность оригиналу.

Степень адекватности разработанных моделей ПЗ статистическим данным можно оценить по результатам прогнозирования и сравнения этих результатов с фактическими за несколько лет (глава 6-я). Другие двухпараметрические модели, алгоритмы реализации которых позволяют более просто учитывать значения как необходимого математического ожидания, так и необходимого среднего квадратического отклонения моделируемых ПЗ, рассматриваются в § 3.4.2.

Что можно делать с помощью моделей (3.2) или (3.3)? Такие модели помогают, например, осуществлять планирование среднего числа дней работы койки в рассматриваемом году с учётом гарантированного обеспечения

пациентов койками с надёжностью  $P$ . Действительно, если согласно прогнозу с помощью рассмотренных моделей в указанном году оценка математического ожидания значения  $T$  ожидается равной  $\overline{M(T)}$ , а вероятность (надёжность) обеспечения поступающего в некоторую больницу пациента больничной койкой должна быть не менее  $P$  (например, 0,95), то для гарантированного обеспечения больного койкой с вероятностью не менее  $P$  следует планировать, что среднее время пребывания больного на койке равно не  $\overline{M(T)}$ , а  $T_p$ , определяемое из условия

$$P = F(T_p) = F(T) = 1 - e^{-\frac{0,714 T_p^3}{M^3(T)}}.$$

После соответствующих преобразований из приведённого условия следует:

$$T_p = \sqrt{-\frac{\ln(1-P)}{0,714}} M(T).$$

Для случая  $P = 0,95$  согласно последнему выражению получаем:  $T_p = 1,613M(T)$ . По такому времени пребывания больного на больничной койке следует рассчитывать и необходимое число больничных коек для контингента жителей рассматриваемой административной единицы.

В зависимости от характера использования рассмотренная модель может быть статической или динамической, детерминированной или стохастической. Если принятое выражение для функции плотности применяется только для соответствующих аналитических вычислений, то модель является статической и детерминированной. Если же это выражение используется для моделирования последовательности случайных значений времени пребывания больного на больничной койке, то в целом модель становится динамической и стохастической.

### 3.2. Здоровье населения и статистика его показателей. Базы данных

Среди математических моделей в области медицины и, в частности, здравоохранения [9, 31 – 34, 37, 39, 61, 98, 112 – 114] значительное место занимают модели здоровья населения. При этом рядом авторов [66, 75, 80, 116, 142, 145, 158] рассматривались вопросы интегральной оценки здоровья населения.

Для рассмотрения моделей здоровья населения необходимо вначале уточнить: а что такое “здоровье населения”? Из известных публикаций следует, что различные специалисты понимали понятие “здоровье” неоднозначно. Известно более ста определений этого понятия [95], которые в крупном плане можно классифицировать следующим образом:

- здоровье – это отсутствие болезней;
- здоровье – это нормальная жизнь, хорошее самочувствие;
- здоровье – это единство морфологических, психоэмоциональных и социально-экономических констант.

Как указывает академик Ю.П. Лисицын, общим для приведенных определений и подходов является то, что здоровье понимается как нечто противоположное понятию “нездоровье” и зависит от распространенности тех или иных болезней, дефектов развития, уровня смертности и т.д. Р. Капра определил здоровье как благополучие, являющееся следствием динамического равновесия, которое учитывает как физические и психологические аспекты существования организма, так и взаимодействие с природной и социальной окружающей средой. Н. Амосов определяет здоровье как математическую функцию – “сумму резервных мощностей” основных функциональных систем индивидуумов, которые следует выразить через “коэффициент резерва, как максимальное количество функции, отнесённое к её нормальному уровню”. Согласно В.В. Парину, здоровье – это такое состояние организма, при котором обеспечивается максимальная адаптивность индивидуума, т.е. наиболее универсальное свойство всего живого, которое лежит в основе фило- и

онтогенеза человека. М.Ф. Сауткин считает [95], что “здоровье характеризуется способностью организма сохранять основные параметры гомеостаза в условиях изменения его внутренней среды, противостоять воздействиям инфекции, физических, химических и психических факторов, являясь интегральным показателем жизнедеятельности организма как в данный конкретный момент, так и на протяжении всей жизни”.

Отправной точкой для медико-социальной интерпретации здоровья является определение, принятое Всемирной организацией здравоохранения (ВОЗ): “Здоровье является состоянием полного физического, духовного и социального благополучия, а не только отсутствием болезней и физических дефектов”. Таким образом, здоровье рассматривается как состояние, позволяющее вести активную в социальном и экономическом плане жизнь. Современное понятие здоровья базируется на пяти общепринятых критериях: адаптивность, равновесие, гармоничность, благополучие и способность функционировать. При этом ведущим критерием является адаптивность [133], “так как остальные являются её важнейшими производными”.

При оценке здоровья принято выделять 3 его уровня [75, 95, 157]:

- здоровье отдельного человека (индивидуальное здоровье);
- групповое здоровье (здоровье социальных, этнических групп, населения административных территорий);
- общественное здоровье (здоровье общества, субпопуляции в целом).

Согласно Ю.П. Лисицыну, “Общественное здоровье – не только совокупность характеристик и признаков индивидуального здоровья, но и интеграция социально-экономических черт, делающих его жизненно необходимой частью того социального организма, каким является общество. Общественное здоровье – результат социально опосредованных воздействий, проявляющихся через образ жизни человека, группы населения”. В.П. Казначеев определяет общественное здоровье как “процесс социально-исторического развития социально-природной, антропо-экологической жизнеспособности населения в ряду поколений, повышения его социально-трудовой активности

в общественно значимых целях, совершенствования психофизиологических возможностей человека”.

Различные количественные характеристики здоровья (рождаемость, смертность, степень нетрудоспособности и др.) принято называть показателями здоровья населения. При этом показатели группового и общественного здоровья определяются как средние статистические значения по одноимённым показателям для всех жителей рассматриваемой административной единицы, т.е. как оценки их математических ожиданий. На практике для краткости такие статистические показатели здоровья населения принято называть просто *показателями здоровья*. С точки зрения факторного анализа моделей группового и общественного здоровья ПЗ являются соответствующими факторами.

Очевидно, мониторинг здоровья населения должен базироваться на хорошо организованной статистике показателей здоровья, использующей современные компьютерные технологии и единое, по крайней мере для РФ, информационное пространство. Хорошая организация статистики прежде всего предполагает единую систему определения ПЗ и их достоверность. В последние годы достигнуты определённые успехи в создании государственной базы данных показателей здоровья РФ.

Все регионы РФ (области, края, республики, округа) по окончании года составляют отчёт по полученным значениям ПЗ, перечень которых определён ГОСКОМСТАТОМ РФ. Отчётные медицинские статистические данные регионов относятся к категории государственных статистических данных. Из ежегодных отраслевых отчётных статистических данных ГОСКОМСТАТ РФ формирует и публикует статистическую базу данных страны и федеральных округов. Региональные организации здравоохранения могут пользоваться значениями ПЗ из этой базы для любых регионов. Кроме того, каждый регион может вводить для себя и использовать в работе своих организаций дополнительные ПЗ.

Информация по ПЗ населения является частью существенно возросших в последние годы медицинских информационных ресурсов [12, 159, 160 и др.], представляющих собой совокупность сведений, характеризующих деятельность элементов, подсистем и всей системы здравоохранения в целом по отношению к пациентам. Кроме указанных ПЗ к этим ресурсам относятся: качество медицинской помощи, лекарственное обеспечение, кадровое обеспечение организаций здравоохранения и др.

Каждый регион (округ), область, отрасль, вид ПЗ и т.д. фигурируют в государственной базе данных с соответствующими номерами, т.е. хранящаяся в базе данных информация структурирована по указанным номерам. В отличие от обычных информационных файлов организация баз данных предусматривает возможность автоматизированной обработки хранящейся в них структурированной информации. Так, из государственной базы статистики медицинских показателей можно просто получать выборки ПЗ по различным административным единицам и годам, например: рождаемость за определённые годы по Северо-Западному округу РФ (с данными по каждой области и республике этого округа и по городу Санкт-Петербург) или перечень областей РФ, в которых в 2000-м году смертность превышала интересуемое значение. Возможности автоматического получения подобных сведений значительно облегчают использование информации, находящейся в базе данных.

Значения всех ПЗ в государственной базе статистики медицинских показателей приводятся в расчёте на 1000, 10000 или на 100000 жителей соответствующих административных единиц, т.е. в виде относительной величины, умноженной на указанное число жителей. При этом, например, количество случаев определённого вида заболевания в интересуемом году равно  $0,001 \times \text{ПЗ} \times N$ , где  $N$  – число жителей рассматриваемой административной единицы в этом году, а значение ПЗ получено в расчёте на 1000 жителей. Указанный способ представления значений ПЗ, во-первых, упрощает сравнение однотипных ПЗ для регионов с разной численностью населения и, во-вторых, позволяет избежать необходимости хранения величин, значительно

меньших единицы, уменьшая, в частности, объём памяти, требующийся для хранения статистики ПЗ.

При создании государственной базы медицинской статистики было предусмотрено представлять в ней значения 117-ти ПЗ с точностью до 4-го знака после запятой, так как такой точности вполне достаточно для различения отличий в данных от разных административных единиц, а более высокая точность требует увеличения объёма памяти для хранения базы. В ряде случаев значения ПЗ были представлены и отражены в базе данных с точностью лишь до 2-го знака после запятой. Общее количество ПЗ, представлявшихся в государственную статистику до 2005-го года, равно 95; затем оно возросло.

Отметим, что в регионах кроме государственной базы медицинской статистики дополнительно могут использоваться локальные, персонифицированные базы медицинских статистических данных. Примерами могут служить базы отдельных организаций здравоохранения, районов большого города или базы, содержащие виды ПЗ, отсутствующих в государственной базе. В частности, в Новгородском научном медицинском центре РАМН разработана единая персонифицированная база данных «Здоровье населения Новгородской области», позволяющая отслеживать в памяти компьютера несколько показателей здоровья каждого жителя области (§ 7.3.1, [34]).

В дальнейшем все ПЗ, приводимые в государственной базе медицинской статистики, будем называть *стандартными*, а остальные ПЗ – *дополнительными*. Смысл введения в модели здоровья населения дополнительных, “нестандартных” ПЗ заключается в стремлении дифференцировать характер влияния различных заболеваний и их последствий на возможности повседневной жизнедеятельности индивидуумов, на создание соответствующих ограничений в этих возможностях, а также учесть возрастные данные индивидуумов. В конечном счёте это преследует цель – получить более точное соответствие значений ИП фактическому состоянию здоровья населения.

В качестве примера нестандартных ПЗ, можно привести классы нетрудоспособности (табл. 3.2), на которые в работе [158] предлагалось разбить всех индивидуумов, испытывающих какие-либо ограничения в повседневной жизни:

Т а б л и ц а 3.2. Классы нетрудоспособности

Класс	Описание	Вес
1	Ограниченная способность (на 50% и более) исполнять по крайней мере одну деятельность в одной из следующих областей: отдых, образование, воспроизведение или профессиональная деятельность.	0,096
2	Ограниченная способность исполнять большинство действий в одной из следующих областей: отдых, образование, воспроизведение или профессиональная деятельность.	0,220
3	Ограниченная способность исполнять действия в двух или более следующих областей: отдых, образование, воспроизведение или профессиональная деятельность.	0,400
4	Ограниченная способность исполнять большинство действий во всех следующих областях: отдых, образование, воспроизведение или профессиональная деятельность.	0,600
5	Требуется помощь для ежедневной инструментальной деятельности типа подготовки пищи, посещения магазина или работы по дому.	0,810
6	Требуется помощь для ежедневной деятельности типа приема пищи, персональной гигиены или использования туалета.	0,920

Отметим, что похожее на приведённое разделение на классы индивидуумов с ограничениями в повседневной жизни приводится и в [66, 134]. Введение веса класса задаёт степень нетрудоспособности. Значения весов устанавливаются экспертами. При этом можно считать, что имеется не несколько ПЗ, представляющих указанные классы, а только один интегральный показатель здоровья, принимающий число значений, равное числу классов.

Вопросы оценки здоровья населения на основе нестандартных показателей рассматриваются в главе 7.

Важной характеристикой здоровья населения, показатели которой не приводятся в государственной медицинской статистике, является физическое развитие населения [93, 133]. Под физическим развитием понимают комплекс непрерывно происходящих в организме биологических процессов, фенотипическим проявлением которых (в отдельных возрастных периодах) являются индивидуальные особенности размеров частей тела, массы, силы, уровня работоспособности и других физических характеристик человека. Показатели физического развития делятся на антропометрические (показатели внешнего осмотра), антропометрические (рост, масса тела, периметр грудной клетки и др.) и антропофизиометрические (показатели силы кисти руки, становой силы, жизненной ёмкости лёгких, физической работоспособности человека). Все эти показатели и методики их измерения подробно описаны, например, в [93].

При оценке физического развития индивидуумов часто пользуются интегральными показателями физического развития. В частности, Г.А. Апанасенко предложены шкалы соматического здоровья для мужчин и женщин, имеющие 5 градаций соматического здоровья и учитывающие рост, массу тела, жизненную ёмкость лёгких, частоту сердечных сокращений в покое, силу кисти руки, уровень систолического давления и время восстановления частоты сердечных сокращений после функциональной пробы [93].

В.А. Медик, А.Г. Швецов и М.С. Токмачёв предложили оценивать физическое развитие индивидуумов с помощью одного интегрального показателя – индекса физического состояния (ИФС). При этом под физическим состоянием человека понимается [96, 133] “степень готовности его организма переносить внешние воздействия различного характера в данный конкретный отрезок времени в зависимости от уровня его физических (двигательных) качеств, особенностей физического развития, функциональных возможностей

отдельных систем организма, наличием или отсутствием заболеваний и травм”. ИФС вычисляется согласно выражению:

$$\text{ИФС} = 0,2\text{СИ} + 0,3\text{ПСИ} + 0,5\text{КСИ},$$

где СИ – соматический индекс, ПСИ – пульмоно-соматический индекс, КСИ – кардио-соматический индекс. Методики определения значений СИ, ПСИ и КСИ приводятся в [93]. С 2005 года ИФС входит в перечень показателей здоровья системы здравоохранения Новгородской области.

Очевидно, в различных административных единицах могут создаваться базы нестандартных ПЗ и в случае введения одинакового перечня таких ПЗ во всех регионах страны – также использоваться для построения ИП здоровья населения.

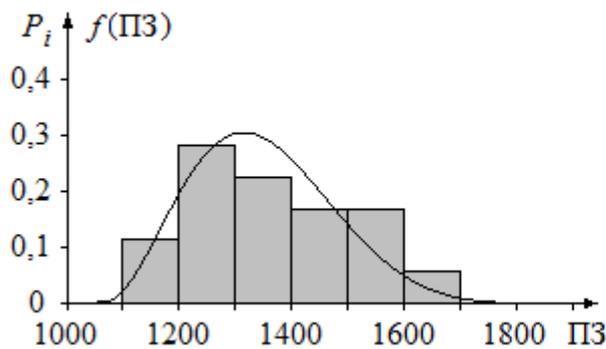
Следует отметить, что статистикой любых ПЗ имеет смысл пользоваться для административных единиц с достаточной численностью населения. В противном случае будет иметь место значительный разброс значений ПЗ от года к году, так как среднее квадратическое отклонение выборочного математического ожидания ПЗ в отчётном году обратно пропорционально квадратному корню от численности населения рассматриваемой административной единицы [29, 89, 94, 107].

### **3.3. Анализ распределений показателей здоровья населения и показателей работы учреждений здравоохранения**

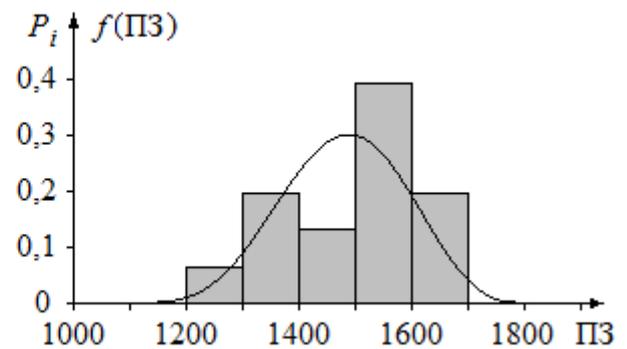
Для моделирования рассматриваемых показателей важно знать тип распределений, которым подчиняются различные показатели. Статистический анализ распределений показателей здоровья [82, 159 и др.] и показателей работы учреждений здравоохранения, являющихся случайными величинами, показывает, что в большинстве случаев распределения этих можно аппроксимировать распределениями колоколообразного вида. В подтверждение данного утверждения на рис. 3.3 и 3.4 приводятся гистограммы  $P_i$  показателей здоровья населения административных единиц различных регионов и

деятельности организаций здравоохранения по данным в 2007 г., а также аппроксимирующие их классические распределения.

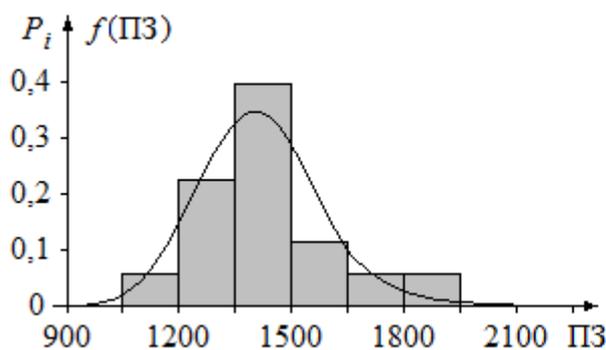
В качестве классических распределений во всех случаях, кроме распределения на рис. 3.3в, использовано бета-распределение, поясняемое в § 3.5.2. На рис. 3.3в для аппроксимации гистограммы использовано распределение Вейбулла (3.3), оказавшееся более удобным. Значения всех ПЗ указаны в расчёте на 1000 человек соответствующего населения, то есть рождаемость, общая заболеваемость и общая смертность – в расчёте на 1000 человек населения данного района или города, заболеваемость детского населения и его первичная инвалидность – в расчёте на 1000 детей и т.д. Значения охвата медосмотром населения, относящегося к указанным на рис. 3.4.3 возрастным группам, приведены в процентах (Пр).



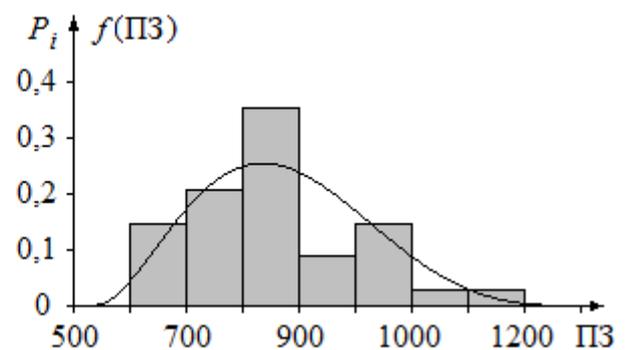
а) Общая заболеваемость взрослого населения по обращениям (Центральный федеральный округ)



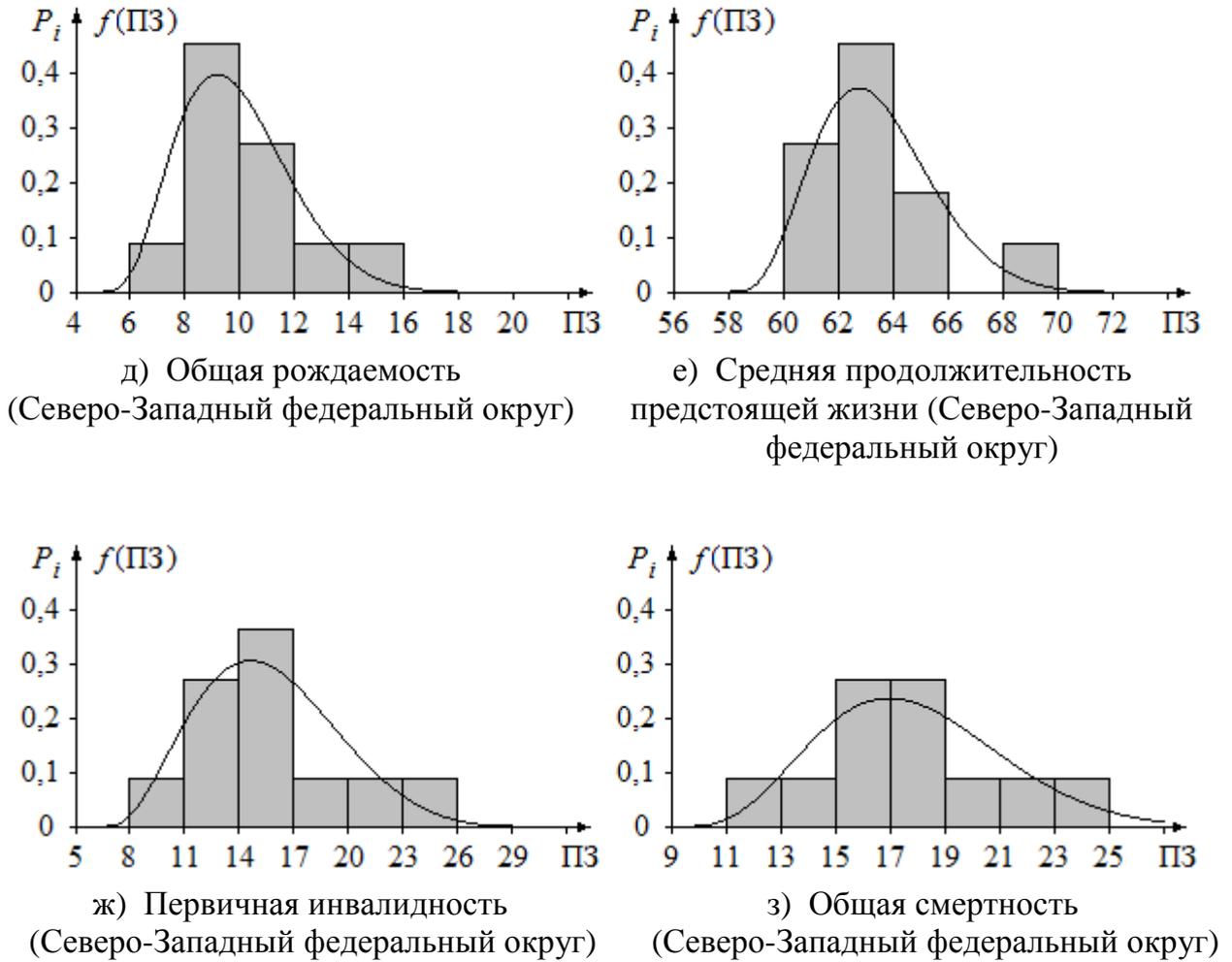
б) Общая заболеваемость взрослого населения по обращениям (Приволжский федеральный округ)



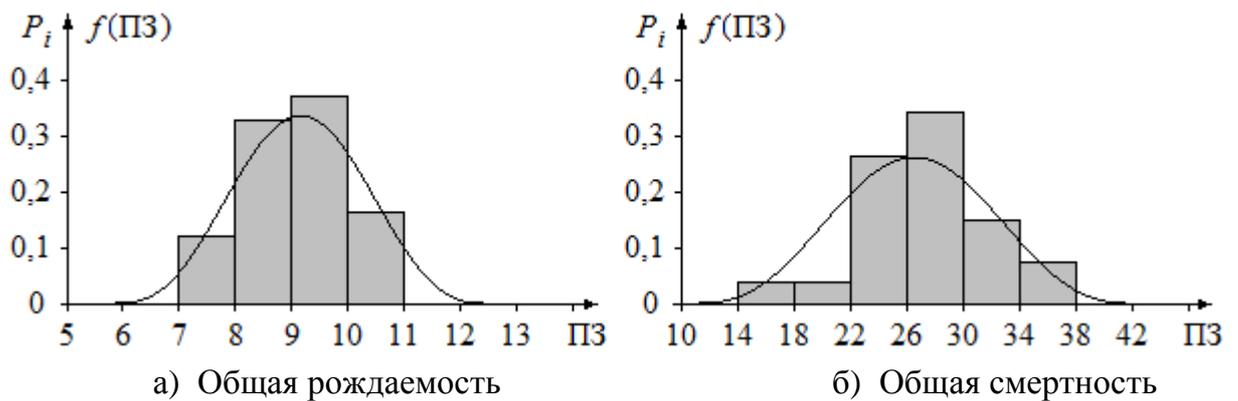
в) Общая заболеваемость взрослого населения по обращениям (Сибирский федеральный округ)

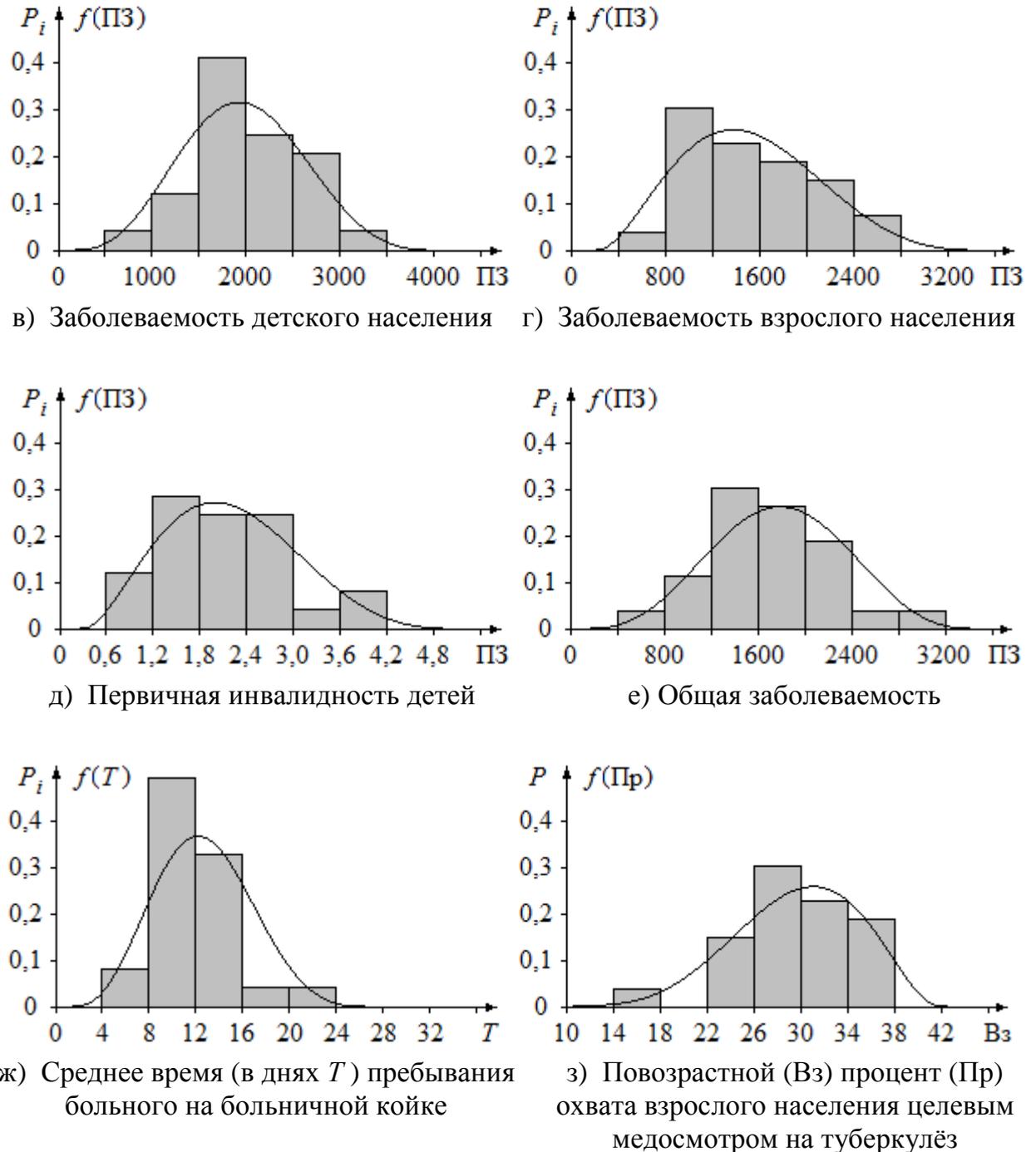


г) Первичная заболеваемость по классу новообразований (регионы Волжского бассейна)



**Рис. 3.3.** Статистические распределения показателей здоровья населения административных единиц федеральных округов РФ в 2007 г.





**Рис. 3.4.** Статистические распределения показателей здоровья населения и деятельности организаций здравоохранения Новгородской области в 2007 г.

С течением времени параметры рассмотренных распределений могут изменяться, однако их колоколообразный характер сохраняется. Указанное изменение связано с изменением соответствующих характеристик временных рядов, образуемых показателями здоровья населения и деятельности организаций здравоохранения.

### 3.4. Временные ряды в статистике здравоохранения и их характеристики

При проведении статистических исследований получаемые результаты часто представляются в виде упорядоченных последовательностей значений этих результатов, называемых элементами последовательности. Упорядочение заключается в том, что каждому элементу последовательности присваивается соответствующий номер. При этом полученные результаты записываются в порядке возрастания их номеров.

Если элементами указанной последовательности являются значения  $X(t)$  некоторого протекающего во времени процесса, то её обычно называют временным рядом [2, 5, 8, 15, 22, 27, 29, 46, 47, 120, 129, 156]. Поскольку при этом время  $t$  является дискретным, то под ним можно понимать номер шага ряда.

Временные ряды получили широкое применение в медицинской статистике. Их примерами являются упорядоченные последовательности значений температуры больного, содержания сахара в его крови, ПЗ и ИП здоровья населения и др. Причём элементами временного ряда могут быть не только скалярные величины (числа), но и векторные величины. Такие временные ряды являются многомерными (векторными). Примером многомерного временного ряда является ряд из 117 значений ПЗ, публикуемых ежегодно ГОСКОМСТАТом РФ, т.е. ряд реализаций 117-координатного вектора, каждая координата которого представляет собой соответствующий ПЗ.

На практике приходится иметь дело не с бесконечными рядами, а с их выборками. Выборки многомерных временных рядов обычно представляются в виде таблиц, а скалярных временных рядов – и в виде графиков (кусочно-линейных или ступенчатых).

Временные ряды могут быть детерминированными (неслучайными) и стохастическими (случайными). В медицинской практике почти всегда приходится иметь дело со стохастическими рядами, характеризующимися опре-

делённой случайностью значений их элементов. Если интервал изменения образующих временной ряд значений рассматриваемого показателя  $X$  разбиваются на несколько уровней, то и элементами ряда являются уровни этого показателя. Такие уровни часто используются в медицине, экономике и т. д.

К основным характеристикам временных рядов относятся математическое ожидание, среднее квадратическое отклонение и корреляционные функции (автокорреляционные и взаимно корреляционные).

Математическим ожиданием стохастического временного ряда  $X(t)$  называется детерминированный временной ряд  $m(t)$  такой, что для каждого момента времени  $t$

$$m(t) = M[X(t)].$$

Функцию  $m(t)$  часто называют средним значением ряда (текущим средним). Однако в большинстве случаев эта функция бывает неизвестна. Но пользуясь реализацией ряда, можно построить сглаженную функцию, принимаемую за оценку  $\bar{m}(t)$ . Такую функцию часто получают согласно алгоритму

$$\bar{m}(t) = \frac{1}{n} \sum_{k=0}^{n-1} x(t-k),$$

реализующему метод скользящего среднего [5, 15, 129].

При рассмотрении соответствующих временных рядов нередко осуществляют декомпозицию их на детерминированную долгосрочную составляющую, представляющую общую тенденцию изменения ряда и называемую *трендом* [120, 129, 156], на детерминированную сезонную составляющую, циклическую составляющую и случайную составляющую. Циклическую составляющую чаще всего считают детерминированной. Обычно временной ряд имеет ещё и случайную составляющую. Для временных рядов, характеризующих состояние здоровья, сезонная составляющая возникает, например, при шаге ряда в один месяц. В этом случае на значения элементов ряда будет влиять изменение интенсивности заболеваний в различные времена года. Например, количество простудных заболеваний возрастает в холод-

ное время года и уменьшается в тёплый период. Но при шаге в один год, установленном ГОСКОМСТАТОм РФ для формирования статистики ПЗ, такие сезонные явления не сказываются на годовых значениях ПЗ. Отсутствует в рядах ПЗ и циклическая составляющая, имеющая место, например, в рядах соответствующих метеорологических параметров.

Поскольку в государственной статистике значения ПЗ приводятся по итогам каждого календарного года, то в большинстве рассматриваемых в монографии задач будем считать, что временные ряды ПЗ имеют только детерминированную составляющую, т.е. тренд, и случайную составляющую. Тренд (реже – *дрейф*) получается путём сглаживания временного ряда ПЗ. Он представляет собой плавно изменяющуюся детерминированную компоненту, учитывающую медленно изменяющиеся факторы. Примерами трендов могут служить возрастные изменения в организме человека, демографические процессы, влияние на эффективность работы системы здравоохранения условий её финансирования и др. Тренд характеризует общую тенденцию развития процесса и обычно зависит от многих факторов. Так, показатель (уровень) заболеваемости населения некоторого региона зависит природных и экологических условий в данном регионе, от профилактических мероприятий, проводимых органами здравоохранения, и т.д.

Для рассматриваемых временных рядов можно предложить следующее математическое определение понятия «тренд»: *трендом временного ряда, не имеющего сезонной и циклической составляющих, называется функция, приближающая функцию математического ожидания этого ряда, определённая на том же множестве значений аргумента, что и временной ряд, и удовлетворяющая выбранному критерию её близости к математическому ожиданию ряда.*

Анализ текущих значений оценок математического ожидания или трендов рассматриваемых временных рядов является важной задачей в исследовании экономических, демографических, медико-социальных, экологических и других процессов. Но значения трендов, определяемых при таком

исследовании, совпадают с текущими значениями математических ожиданий временных рядов только при стационарности этих рядов [15], т.е. для рядов с не изменяющимися во времени характеристиками. Для нестационарных рядов, даже не имеющих сезонной и циклической составляющих, ввиду наличия инерционности при определении текущего среднего значения ряда получаются соответствующие погрешности [16]. Поэтому фактически получают тренд, несколько отличающийся от текущего среднего. Причём разным алгоритмам определения тренда соответствуют и несколько отличающиеся друг от друга тренды.

Автокорреляционной функцией стохастического временного ряда  $X(t)$  для интервала  $\tau$  называется детерминированный временной ряд  $R(t, \tau)$  такой, что для каждой пары значений  $t$  и  $\tau$ .

$$R(t, \tau) = \frac{M [(X(t) - m(t))(X(t + \tau) - m(t + \tau))]}{\sigma(t)\sigma(t + \tau)} = \frac{\text{cov}[X(t), X(t + \tau)]}{\sigma(t)\sigma(t + \tau)}. \quad (3.4)$$

Числитель выражения (3.4) является ковариационной функцией, т.е. ковариацией (п. 2.3.1), изменяющейся с изменением параметров  $t$  и  $\tau$ . Значения  $R(t, \tau)$  для каждой пары  $(t, \tau)$  называют коэффициентами корреляции. Следовательно, автокорреляционная функция временного ряда является последовательностью коэффициентов корреляции этого ряда. Иногда её называют нормированной ковариационной (нормированной автокорреляционной) функцией, так как значения коэффициентов корреляции изменяются в интервале  $[-1, 1]$ . Для независимых случайных величин  $X(t)$  и  $X(t + \tau)$  ковариация и коэффициент их корреляции равны нулю. Чем больше значение модуля этого коэффициента, тем сильнее взаимосвязь этих величин. При  $\tau = 0$  ковариация случайных величин  $X(t)$  является дисперсией этих величин:  $\text{cov}[X(t), X(t)] = D[X(t)]$ . Заметим, что в случае, когда математические ожидания  $m(t)$  и  $m(t + \tau)$  не известны, вместо их в выражении (3.4) используют значения тренда  $Tr[X(t)]$  и  $Tr[X(t + \tau)]$  для шагов  $t$  и  $t + \tau$ .

Для векторных временных рядов, кроме автокорреляционной функции, к основным характеристикам временного ряда относят ещё взаимные корреляционные функции, характеризующие связь соответствующих координат вектора. Для  $i$ -й и  $j$ -й координат вектора коэффициент корреляции  $X_i(t)X_k(t+\tau)$  определяется согласно выражению

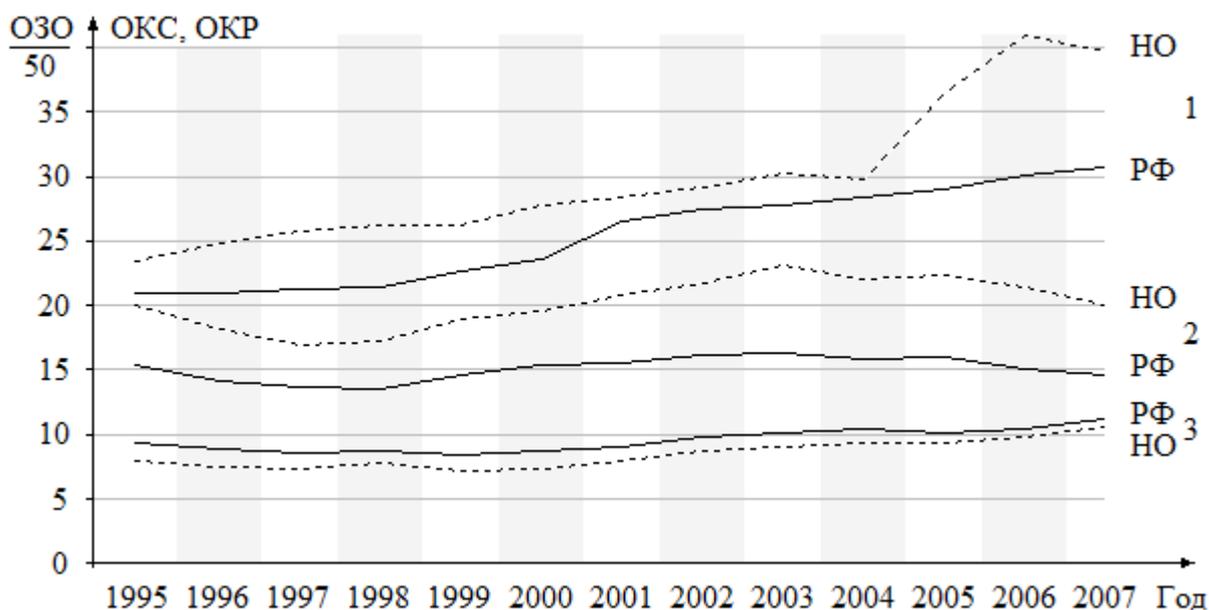
$$R[X_i(t)X_k(t+\tau)] = \frac{M[(X_i(t)-m_i(t))(X_k(t+\tau)-m_k(t+\tau))]}{\sigma_i(t)\sigma_k(t+\tau)} = \frac{\text{cov}[X_i(t), X_k(t+\tau)]}{\sigma_i(t)\sigma_k(t+\tau)},$$

числитель которого является ковариацией случайных величин  $X_i(t)$  и  $X_k(t+\tau)$ . Для стационарных временных рядов корреляционная функция зависит только от интервала  $\tau$ .

Знак коэффициента корреляции соответствующих величин указывает на их положительную или отрицательную корреляцию. При положительной корреляции рассматриваемые величины имеют идентичные или относительно идентичные тенденции изменения, а при отрицательной – противоположные или относительно противоположные. Указанное поясняется выборками из временных рядов ПЗ, графики которых приведены на рис. 3.5. Они свидетельствуют о том, что ПЗ населения больших регионов (РФ) обычно менее динамичны, чем аналогичные ПЗ населения малых регионов (Новгородская область).

На интервале с 1995 по 2007 год тенденции изменения ПЗ общей заболеваемости, общей смертности и рождаемости имели относительно одинаковые тенденции изменения. Поэтому все выборочные значения коэффициентов их корреляции при  $\tau=0$  оказались положительными. Для РФ  $R_{12} = 0,895$ ,  $R_{13} = 0,795$ ,  $R_{23} = 0,782$ , а для Новгородской области  $R_{12} = 0,884$ ,  $R_{13} = 0,807$ ,  $R_{23} = 0,845$ .

Отметим, что при определении оценок текущих характеристик временных рядов по ограниченным выборкам их элементов нередко предполагают стационарность рассматриваемых рядов (такое предположение было сделано



**Рис. 3.5.** Динамика статистических показателей: общей заболеваемости по обращениям (ОЗО - 1), общего коэффициента смертности (ОКС - 2) и общего коэффициента рождаемости (ОКР - 3) на 1000 человек населения для Российской Федерации (сплошные линии) и Новгородской области (НО)

и в рассмотренном примере). В этом случае для определения значений соответствующих оценок по выборке  $n$  последовательных значений  $X(t)$  используются выражения [25, 91, 94, 107 и др.]:

$$\text{cov}(X_i X_k) = \frac{1}{n-1} \sum_{t=1}^n [X_i(t) X_k(t) - \bar{M}(X_i) \bar{M}(X_k)],$$

$$\bar{M}(X_i) = \frac{1}{n} \sum_{t=1}^n X_i(t), \quad \bar{R}_{ik} = \frac{\overline{\text{cov}}(X_i X_k)}{\sqrt{\bar{D}(X_i) \bar{D}(X_k)}}.$$

В дальнейшем увидим, что получаемую таким образом последовательность оценок для  $M(X)$  в ряде случаев можно принять за функцию, являющуюся трендом.

Значения коэффициентов корреляции нескольких случайных величин и оценок этих коэффициентов обычно представляют в виде корреляционных матриц  $\mathbf{R}$ . Так, в рассмотренном примере получено:

$$\bar{\mathbf{R}} = \begin{vmatrix} \bar{R}_{11} & \bar{R}_{12} & \bar{R}_{13} \\ \bar{R}_{21} & \bar{R}_{22} & \bar{R}_{23} \\ \bar{R}_{31} & \bar{R}_{32} & \bar{R}_{33} \end{vmatrix} = \begin{vmatrix} 1 & 0,741 & -0,661 \\ 0,741 & 1 & -0,956 \\ -0,661 & -0,956 & 1 \end{vmatrix}$$

Кроме указанных выше основных характеристик временных рядов, при исследовании временных рядов медицинских показателей нередко используются и другие характеристики, в частности, характеристики, использующие разбиение бесчисленного множества возможных значений рассматриваемого показателя  $X$  на конечное число уровней:  $X_0, X_1, \dots, X_n$ . К таким характеристикам относятся, например, приращения  $X$  от уровня  $X_i$  до уровня  $X_j$ , а также характеристики на основе этих приращений\*.

Основными задачами анализа временных рядов в медицине являются:

1. Сглаживание и фильтрация полученных рядов статистических параметров, выделение тренда и случайной составляющей (декомпозиция ряда).
2. Моделирование временных рядов (композиция моделей тренда и случайной составляющей).
3. Прогнозирование временных рядов.
4. Корреляционный анализ элементов одного, двух и более рядов.

При решении первой из перечисленных задач прежде всего выделяется тренд. Затем путём вычитания тренда из полученного ряда находят случайную компоненту ряда. Однако, если временной ряд имеет циклическую или сезонную компоненту, то вначале из ряда выделяются эти компоненты.

Необходимость моделирования временных рядов возникает в моделях, на входы которых могут поступать соответствующие временные ряды. При этом, изменяя характеристики моделируемых рядов, можно проанализировать влияние этих характеристик на получаемые путём моделирования результаты. Моделируя стохастический временной ряд, используют полученные при его декомпозиции тренд и случайную составляющую, которые моделируемые отдельно. Затем осуществляется композиция временного ряда в

---

\* Более подробно см. главу 7.

соответствии с принятой моделью временного ряда [129, 156]. В качестве такой модели чаще всего используют *аддитивную* модель ряда, согласно которой *ряд является суммой указанных его компонент*. Такая модель является наиболее подходящей для моделирования временных рядов ПЗ. Поэтому в дальнейшем авторы используют только аддитивную модель. Заметим, что в некоторых прикладных задачах используется *мультипликативная* модель ряда, рассматривающая ряд как произведение детерминированной и случайной компонент. Если такое произведение прологарифмировать то получим аддитивную модель относительно логарифмов этих компонент.

При построении моделей тренда и случайной компоненты следует обеспечить достаточную их адекватность моделируемому оригиналу (если временной ряд оригинала имеется). Обычно это пытаются сделать в процессе моделирования и сравнении результатов моделирования с фактическими данными.

Третья из перечисленных задач анализа временных рядов, прогнозирование, имеет важное значение во многих областях и отраслях. Алгоритмы прогнозирования представляют собой модели развития соответствующих процессов. Если, например, модель позволит врачу лучше оценивать перспективы развития болезни и её лечения, то ценность такой модели неоспорима. При решении задач прогнозирования по существу получают наиболее вероятный отрезок тренда ряда ПЗ или какого либо другого рассматриваемого показателя на предстоящих шагах модели.

Корреляционный анализ статистических временных рядов, как уже было показано на примере выборок из трёх рядов, графики которых приведены на рис. 3.3, позволяет выяснить характер и степень взаимной зависимости элементов внутри каждого ряда (автокорреляционные характеристики) и взаимного влияния элементов разных рядов друг на друга (взаимно корреляционные характеристики). Результаты такого анализа важны для обоснованного построения модели исследуемой системы.

### 3.5. Моделирование показателей здоровья и их временных рядов

Как уже указывалось, при построении и исследовании соответствующих математических моделей, например, моделей прогнозирования значений статистических и интегральных показателей здоровья, возникает необходимость в моделировании временных рядов этих показателей. Но перед моделированием неисследованного ряда всегда необходимо оценить характеристики этого ряда, игнорирование которых при моделировании может привести к неадекватности результатов моделирования искомому решению.

В связи с указанным при моделировании временных рядов ПЗ, приходится решать по крайней мере две следующие задачи:

1. Выяснить основные характеристики тренда ряда, на основе которого строится модель ряда, выбрать модель генерирования тренда, построить или выбрать алгоритм и разработать программу его моделирования, проверить соответствие получаемого при моделировании тренда тренду имеющегося временного ряда.
2. Установить одномерный, в первом приближении считающийся независимым от тренда ряда закон распределения случайной составляющей ряда, элементы которой считаются независимыми, и построить или выбрать алгоритм реализации этой составляющей с выбранными одномерным распределением.

Разумеется, при решении каждой из перечисленных задач следует руководствоваться имеющимися статистическими данными. Генерируемый моделью временной ряд, представляющий собой сумму тренда и случайной составляющей, также полезно оценить на предмет адекватности его характеристик характеристикам реального ряда. При такой оценке можно проверить на сколько корреляционная функция суммы полученных тренда и независимой случайной составляющей соответствует корреляционной функции исходного ряда. При исследовании временных рядов и построении их моделей можно выбрать вариант модели, сразу генерирующей последовательность соответствующих показателей, или вариант модели, формирующей вначале отдельно

тренд ряда и его случайную составляющую, объединяемые затем в последовательность элементов ряда. Второй вариант считается предпочтительным в том случае, когда имеют место сезонная или циклическая компоненты ряда. Если временной ряд является стационарным, то распределения значений элементов ряда (показателей здоровья) и случайной составляющей ряда отличаются только тем, что математическое ожидание случайной составляющей равно нулю, т.е. для показателей здоровья оно меньше, чем у второго распределения, на величину математического ожидания ряда. Причём в данном случае значение математического ожидания ряда совпадает с трендом ряда.

### **3.5.1. Модели трендов временных рядов**

Задача построения моделей трендов временных рядов ПЗ может решаться принципиально по-разному в зависимости от длины, т.е. от числа шагов моделируемого ряда. Малая длина ряда (обычно менее 10 шагов) используется в задачах прогнозирования ПЗ. Практически все модели трендов таких рядов основаны на использовании метода наименьших квадратов [10, 25 и др.], с использованием которого осуществляется сглаживание соответствующих участков статистических временных рядов с обеспечением минимума суммы квадратов отклонений сглаживающей функции (тренда ряда) от статистического ряда. В связи с указанным построение трендов рядов малой длины рассмотрим в главе 6, посвящённой прогнозированию показателей здоровья населения.

При построении и исследовании сложных моделей может возникнуть необходимость моделирования рядов весьма большой длины. Так, в [61] при моделировании процессов патогенеза и лечения ишемической болезни сердца было необходимо моделировать временные ряды соответствующих параметров длиной в несколько десятков тысяч шагов. При этом часто конечные результаты решения, получаемые с помощью подобных моделей, достигаются на последнем шаге работы модели. Поэтому важно, чтобы именно эти ре-

зультаты были адекватными результатам, полученным на основе реальных статистических временных рядов. Что же касается текущих характеристик генерируемого моделью ряда, то в указанном случае требования к ним могут быть ослаблены.

Подобные модели временных рядов ПЗ незаменимы, например, при оценке точности различных методов и алгоритмов прогнозирования, так как в этом случае наиболее состоятельными оценками являются статистические оценки отклонения прогнозов от фактических результатов, получаемых через интервал прогнозирования. Ряды же ПЗ длиной хотя бы в несколько десятков тысяч шагов можно получить только с помощью модели.

В литературе, например, в [120, 129], рассматривается несколько моделей трендов, точнее, аппроксимирующих из функций: линейная, параболическая, кубическая, экспоненциальная, логарифмическая и др. Но для длинных временных рядов обычно не удаётся аппроксимировать весь ряд какой-то одной функцией. Поэтому в этом случае тренд разделяется на участки, для каждого из которых подбирается соответствующая модель, т.е. в целом аппроксимирующая функция строится как сплайн (последовательность аппроксимирующих функций). В частности, в ряде приложений для построения трендов используются экспоненциальные сплайны, когда сглаживание временного ряда производится соответствующими участками экспоненты [120, 129].

Согласно проведённому анализу, для моделирования отдельных участков трендов  $Y(t)$  временных рядов ПЗ и ИП для различных регионов России применима и четырёх параметрическая синусоидальная модель этих участков, определяемая выражение

$$Y(t) = A \sin\left(\frac{2\pi t}{T} + \theta\right) + B \quad (3.5)$$

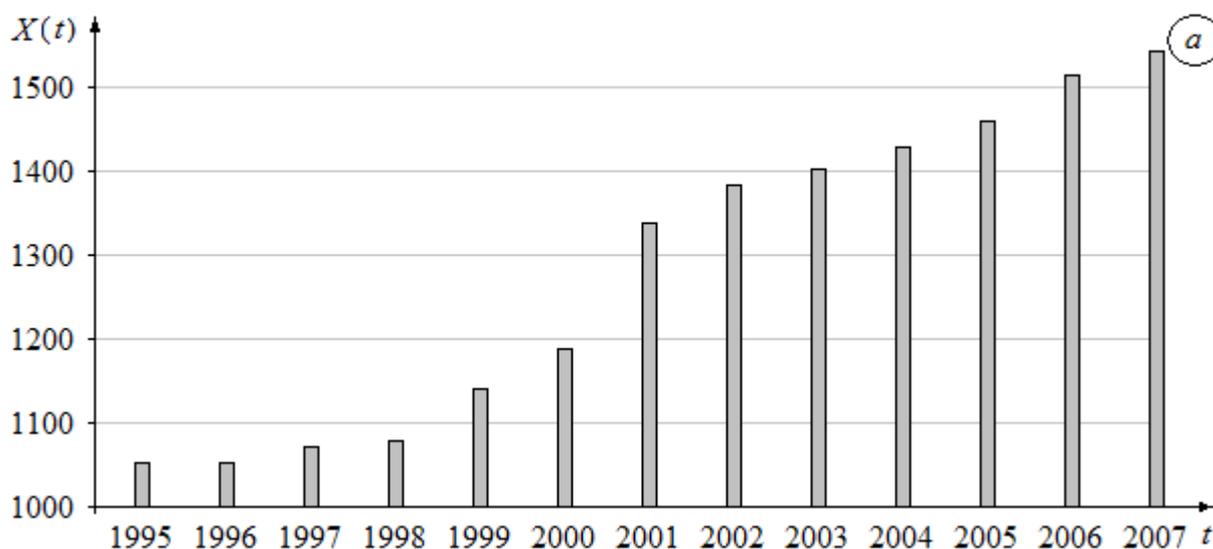
где амплитуда  $A$ , период  $T$  и начальная фаза  $\theta$  подбираются при моделировании,  $B = \overline{M}(X)$ , а  $t$  изменяется с шагом один год. При этом в качестве  $\overline{M}(X)$

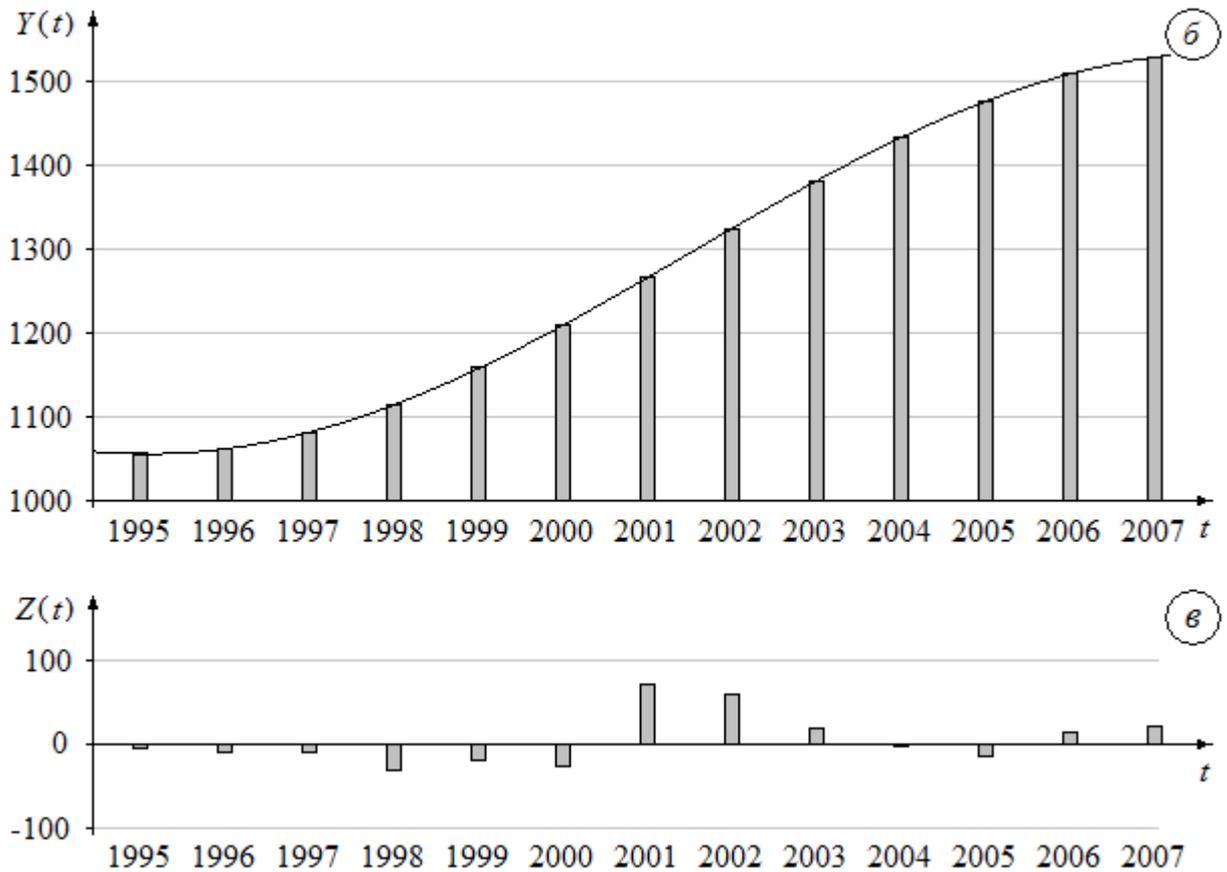
принимается среднее значение элементов  $X(t)$  ряда за наблюдаемое число шагов. Что касается сезонных колебаний ПЗ, то, как уже указывалось, на годовые значения ПЗ они не влияют. Поэтому предлагаемая модель их не учитывает. Вместе с тем синусоидальная модель может оказаться удобной для моделирования сезонных составляющих временных рядов ПЗ.

Выбирая модель  $Y(t)$ , фактически одновременно получают и статистическую последовательность  $Z(t)$  значений случайной составляющей ряда, так как  $X(t)=Y(t)+Z(t)$ . При этом можно предложить следующий необходимый критерий соответствия тренда  $Y(t)$  исходному ряду  $Z(t)$ : *необходимо, чтобы на моделируемом участке временного ряда среднее значение  $Z(t)$  было бы нулевым, т.е.*

$$\overline{M}[Z(t)] = 0. \quad (3.6)$$

Отметим, что используя различные модели выделения трендов и различные параметры этих моделей, можно получить бесконечное число трендов  $Y(t)$ , удовлетворяющих рекомендуемому критерию (3.6) для их выбора. Тем не менее этот критерий, рассматриваемый далее только как необходимый, значительно конкретизирует задачу построения трендов.





**Рис. 3.6.** Временной ряд показателя общей заболеваемости населения РФ (а), тренд ряда и его огибающая (б), случайная составляющая временного ряда (в)

На рис. 3.6 приведён пример декомпозиции участка временного ряда показателя общей заболеваемости  $X(t)$  населения РФ за 12 лет. Модель участка тренда  $Y(t)$  этого ряда при  $t = 1995, 1996, \dots, 2007$  представляет собой детерминированную последовательность с огибающей вида (3.5). При этом  $A = 235$  на 1000 человек населения,  $B = \bar{M}(Y) = 1287$  на 1000 человек населения,  $T = 26$  лет,  $\Theta = -0,55\pi$ ,  $Z(t) = X(t) - Y(t)$ ,  $\bar{M}(Z) = 0$ ,  $\sigma(Z) = 27,72$  на 1000 человек населения.

В рассмотренном примере интервал изменения тренда, равный  $2A$ , составил 30,8% от его среднего значения  $B$ . По-видимому, при малом значении указанного отношения, например, менее 5%, без ущерба для результатов моделирования можно принять, что тренд имеет постоянное значение.

Как уже указывалось, за признак правильности проведения декомпозиции временного ряда можно принять нулевое значение оценки математиче-

ского ожидания значений его случайной составляющей. Для последовательности  $Z(t)$ , приведённой на рис. 3.6.в, это условие выполняется. Однако в общем случае решение задачи декомпозиции временного ряда является неоднозначным, а указанный признак нельзя считать необходимым или достаточным. Отметим также, что для рядов с большим числом элементов часто не удаётся получить тренд на основе только одной модели, “достаточно хорошо” отражающий тенденции изменения ряда на всей его длине. В этом случае для разных участков ряда можно использовать разные модели тренда, то есть использовать сплайновую аппроксимацию тренда [44], сглаживая соответствующие участки ряда алгебраическими полиномами, синусоидами, экспонентами и т.д.

Таким образом, декомпозиция моделируемого нестационарного временного ряда на тренд и случайную составляющую, представляющую собой последовательность независимых случайных величин с соответствующим одномерным законом распределения, позволяет моделировать их отдельно с последующим сложением (композицией). При этом практически отпадает необходимость в решении крайне сложной задачи – моделировании временного ряда с заданными корреляционной функцией и одномерным законом распределения. Стационарные временные ряды имеют постоянное значение математического ожидания и постоянный тренд, равный математическому ожиданию. Поэтому моделирование таких рядов значительно упрощается.

### **3.5.2. Моделирование случайной составляющей временных рядов показателей здоровья и показателей работы медицинских учреждений**

Для разработки модели случайной составляющей ряда некоторого ПЗ необходимо проанализировать статистику этого ПЗ и получить статистический ряд значений случайной составляющей, получающийся после вычитания модели тренда из основного ряда. Во всех приведённых ниже примерах для трендов использовалась удобная для представления нестационарных уча-

стков рядов модель (3.6), параметры которой выбирались на основе метода наименьших квадратов [10, 25].

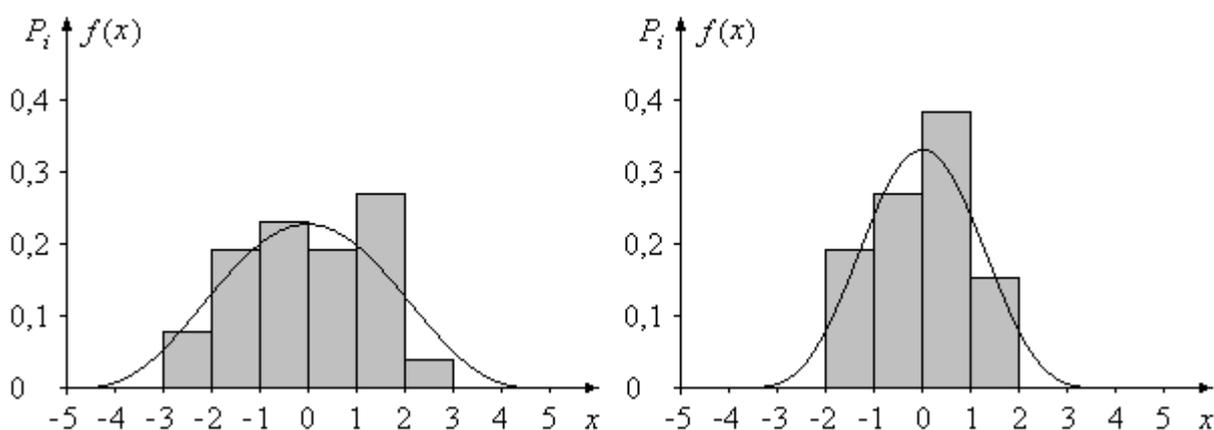
Построение модели случайной составляющей временного ряда заключается в следующем:

- на основе ряда случайной составляющей строится статистическое распределение (ряд распределения или многоугольник распределения) этой составляющей;
- производится аппроксимация статистического распределения некоторой функцией, принимаемой за функцию плотности;
- принимается, что последовательные реализации случайной составляющей являются независимыми.

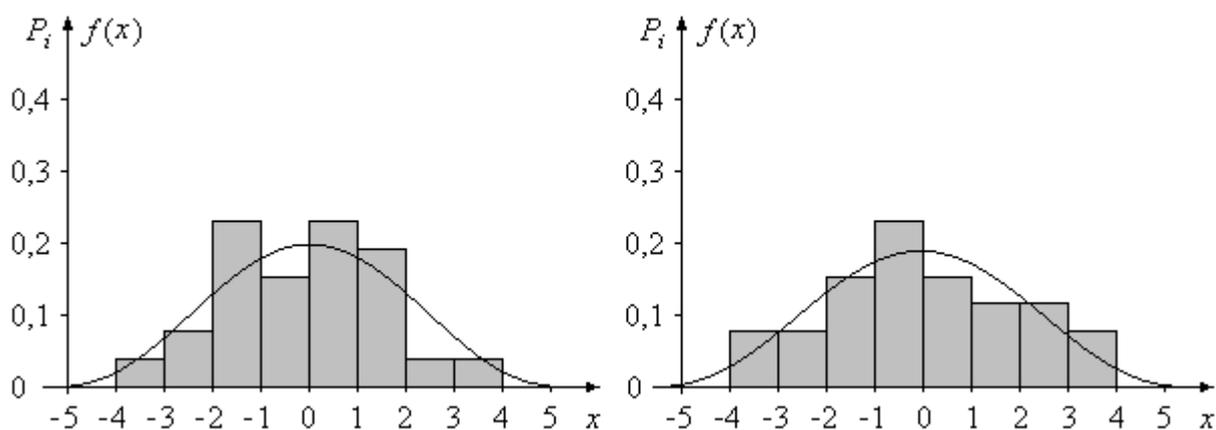
При аппроксимации статистического распределения обычно стараются подогнать его под какое-либо классическое распределение. Независимость значений случайной составляющей временного ряда обычно является некоторым допущением, удобным для моделирования, так как ЭВМ генерирует независимые, равновероятные случайные величины, на основе которых получают случайные величины с другими распределениями. Кроме того, при декомпозиции коротких участков временного ряда достоверность определения корреляционной функции весьма мала ввиду малости объема выборки. Как показывает практика, если среднее значение модуля случайной составляющей не превышает 10 процентов среднего значения модуля тренда, введение в модель независимости значений случайной составляющей вполне допустимо.

Рассмотрим примеры построения моделей случайной составляющей временного ряда, т.е. последовательностей независимых значений этой составляющей. На рис. 3.7 и рис. 3.8 приведены статистические распределения показателей рождаемости и смертности для регионов, сильно отличающихся по численности населения – для РФ в целом и для Новгородской области. Они построены на основе статистических данных за 1980 ÷ 2007 годы. Меньший разброс значений показателя смертности для РФ по сравнению с

Новгородской областью отражает общую тенденцию повышения стабильности ПЗ с увеличением численности населения региона. Непрерывными кривыми показаны возможные модели функций плотности  $f(x)$  этих случайных величин  $X$ . Для них характерным является наличие одного максимума, местоположение которого обычно несколько не совпадает с математическим ожиданием, равным нулю. В обе стороны от максимума  $f(x)$  монотонно уменьшается до нуля к концам отрезка  $[a, b]$ , т.е.  $f(a) = f(b) = 0$ . Такие  $f(x)$  имеют колоколообразную форму.



**Рис. 3.7.** Статистические распределения случайной составляющей временных рядов коэффициентов общей рождаемости (слева) и общей смертности и графики аппроксимирующих их функций плотности для РФ



**Рис. 3.8.** Статистические распределения случайной составляющей временных рядов коэффициентов общей рождаемости (слева) и общей смертности и графики аппроксимирующих их функций плотности для Новгородской области

Приведённые на рисунках гистограммы распределения соответствующих показателей получены в результате статистического анализа рядов их случайной составляющей, которые получались путём вычитания трендов, полученных для указанного временного интервала, из исходных рядов.

Для аппроксимации гистограмм случайной составляющей временных рядов показателей средней продолжительности предстоящей жизни, общей заболеваемости и инвалидности для РФ, всех её федеральных округов и Новгородской области наиболее подходящими также оказались колоколообразные функции. При равенстве единице интеграла от такой функции на отрезке  $[a, b]$  эту функцию можно считать функцией плотности рассматриваемой случайной величины.

Проведённые исследования позволяют утверждать, что колоколообразные распределения хорошо аппроксимируют и гистограммы случайной составляющей временных рядов показателей учреждений здравоохранения.

### **3.5.3. Моделирование случайных величин с колоколообразными распределениями**

Колоколообразные распределения некоторой случайной величины  $X$  можно получать по крайней мере на основе одного из трёх классических распределений. Для распределений, симметричных относительно середины отрезка  $[a, b]$  изменения ПЗ, такими распределениями могут быть усечённое снизу нормальное распределение [57] или распределение, представляющее собой взвешенную сумму нескольких равномерно распределённых случайных величин [130]. Для реализации несимметричных колоколообразных распределений можно использовать бета-распределение [78].

Функция плотности нормального распределения имеет два параметра и определяется выражением

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-M)^2}{2\sigma^2}}, \quad (3.7)$$

в котором параметры  $M$  и  $\sigma$  – соответственно математическое ожидание и среднее квадратическое отклонение случайной величины  $X$ , изменяющейся от  $-\infty$  до  $\infty$ .

Рис. 3.9 поясняет метод усечения функции плотности нормального распределения с  $M(x) = 0,5(a + b)$  снизу. После усечения  $f(x)$  остаётся функ-

ция  $f_*(x)$ , отличная от нуля только в  $[a, b]$ , где она равна  $f(x) - f(x)$ . Умножение этой функции на коэффициент

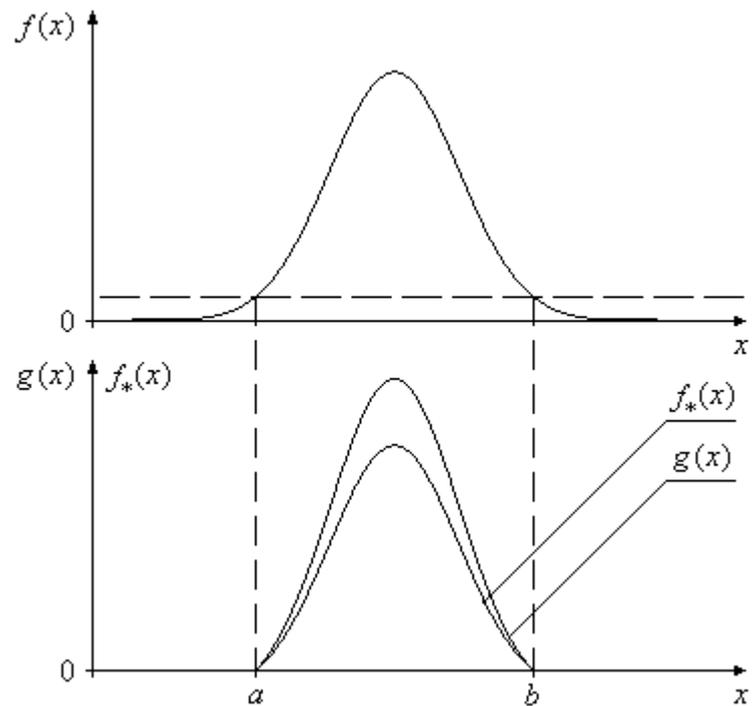
$$k = 1/[2\Phi(b/\sigma) - 1 - 2bf(a)]$$

превращает её в функцию плотности  $g(x) = kf_*(x)$ , так как интеграл от  $g(x)$  на  $[a, b]$ , равен единице.

Распределение  $g(x)$  имеет форму колокола, его  $M(x)$  равно  $0,5(a + b)$ , а значение среднего квадратического отклонения

моделируемой  $X$  можно подобрать путём варьирования значением  $\sigma$  распределения  $f(x)$ . При этом для моделирования случайной составляющей ряда принимается  $a = -b$ .

Моделирование последовательности случайных величин с колоколообразной функцией плотности  $g(x)$  просто выполняется методом исключения [38, 68, 103]. В соответствии с этим методом моделируется последовательность случайных точек  $(x, y)$ , равномерно распределённых в прямоугольнике



**Рис. 3.9.** Реализация колоколообразного распределения с помощью усечения снизу функции плотности нормально распределённой случайной величины  $X$

с основанием  $ab$  и с высотой  $h \geq f_*(x)_{\max}$ . При этом на каждом шаге в случае  $y \leq f_*(x)_{\max}$  значение  $x$  принимается за значение моделируемого показателя. В противном случае данный шаг пропускается, т.е. полученное значение  $x$  из ряда исключается.

Кроме рассмотренного метода, моделирование случайных величин с колоколообразным распределением возможно на основе взвешенной суммы (сложения с весами) нескольких равномерно распределённых случайных величин  $C_i$ . Оно может быть реализовано с помощью алгоритма

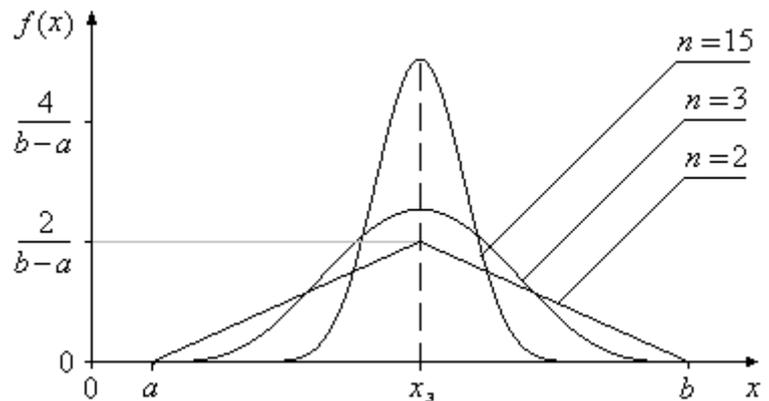
$$x = \frac{1}{n} \sum_{i=1}^n C_i,$$

в котором все  $C_i$  распределены равномерно в промежутке  $[a, b]$ . С увеличением числа слагаемых  $n$ , являющимся параметром распределения, график функции  $f(x)$  вытягивается вверх (рис. 3.10), а значение  $\sigma(x)$  уменьшается в соответствии с выражением

$$\sigma(x) = \frac{b-a}{\sqrt{12n}}.$$

При неограниченном увеличении  $n$   $f(x)$  стремится к  $\delta$ -функции [57, 107].

Рассмотренные два метода позволяют получать симметричные относительно середины отрезка  $[a, b]$ , распределения. При этом значение  $M(x)$  совпадает со значением  $x_3$ , при котором имеет место экстремум функции плотности. Они равны  $0,5(b-a)$ . Однако, как уже отмечалось, большинство распределений рассмат-



**Рис. 3.10.** Графики распределений ПЗ на основе взвешенной суммы  $n$  равновероятных случайных величин при разных  $n$

риваемых показателей несимметричны. Такие распределения удобно аппроксимировать, например, бета-распределением [57].

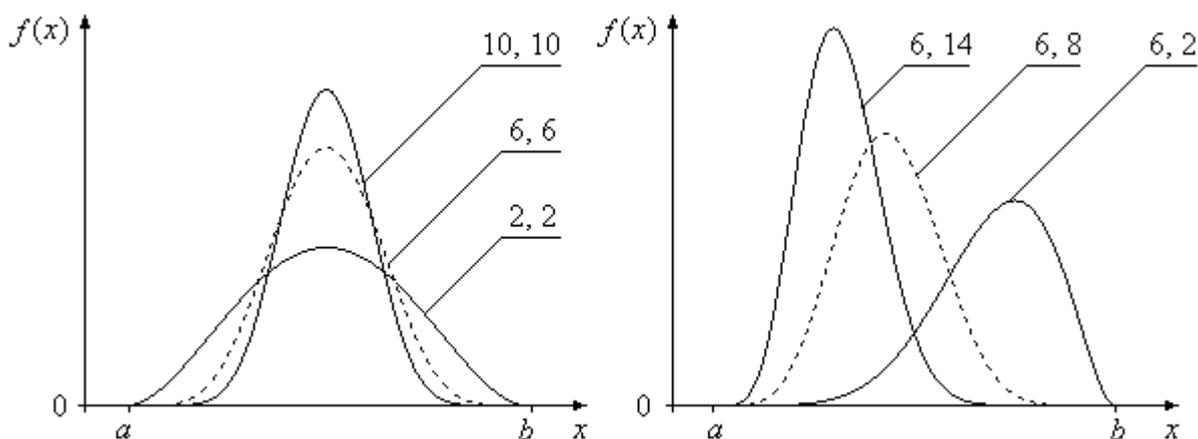
Бета-распределение имеет два параметра, варьирование которыми позволяет моделировать распределения большинства ПЗ здоровья. Для целых значений этих параметров ( $m > 0$  и  $n > 0$ ) функция плотности бета-распределения имеет вид:

$$f(x) = \begin{cases} \frac{m+n+1}{(b-a)^{m+n+1}} C_{m+n}^n (x-a)^m (b-x)^n & \text{при } x \in [a, b], \\ 0 & \text{в противном случае.} \end{cases} \quad (3.8)$$

Математическое ожидание, среднее квадратическое отклонение и соответствующее экстремуму значение  $x_0$  определяются выражениями:

$$M(x) = \frac{a(n+1)+b(m+1)}{m+n+2}, \quad \sigma(x) = \frac{b-a}{m+n} \sqrt{\frac{mn}{m+n+1}}, \quad x_0 = \frac{mb+na}{m+n}.$$

Рис. 3.11 иллюстрирует графики функции (3.8), обозначаемой часто как бета( $m, n$ ). Не трудно заключить, что этим колоколообразным распределением можно варьировать в довольно широких пределах и оно может служить моделью для большинства распределений ПЗ. Однако аналитический расчёт подходящих значений параметров  $m$  и  $n$  для распределения бета( $m, n$ ) по полученным статистическим данным довольно неудобен. Поэтому их обычно подбирают путём моделирования функции (3.8) при различных значениях



**Рис. 3.11.** Графики бета-распределения при разных значениях  $m$  и  $n$

параметров и сравнения её со статистическим распределением моделируемого показателя.

Моделирование ПЗ и случайной составляющей временного ряда ПЗ с распределением бета( $m, n$ ) просто выполняется методом исключения [38, 78, 103], согласно которому на каждом шаге модели моделируются случайные величины  $C_1 \in [a, b]$  и  $C_2 \in [0, x_s]$ , вычисляется значение  $f(c_1)$  и в случае выполнения условия  $c_2 < f(c_1)$  величина  $c_1$  принимается за значение ПЗ. В противном случае указанное значение  $c_1$  игнорируется.

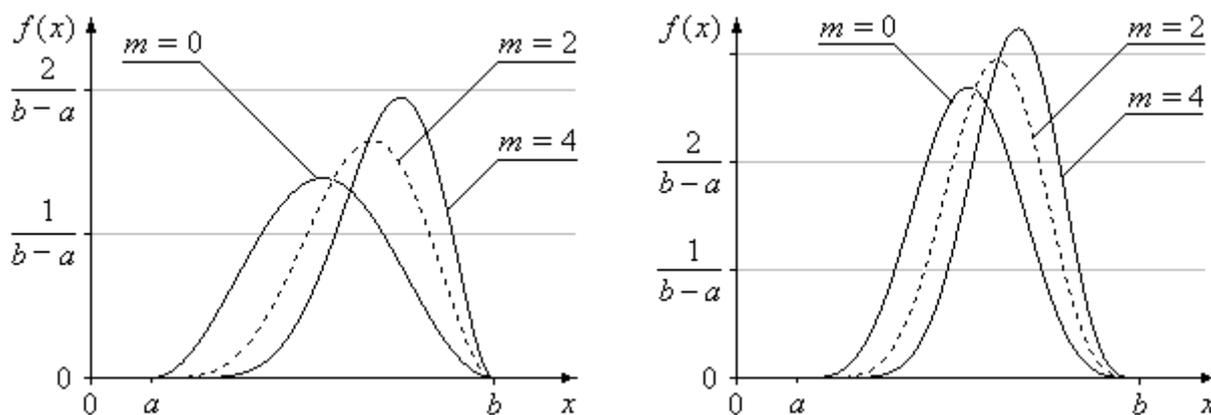
Удобным распределением для моделирования показателей здоровья и показателей работы учреждений здравоохранения, а также для моделирования случайной составляющей их временных рядов является предложенное в [53] множество колоколообразных распределений типа  $x \sin x$ , функция плотности которых в промежутке  $[a, b]$  имеет вид

$$f(x) = M \cdot \left( \frac{x-a}{b-a} \right)^m \cdot \left[ 1 + \sin \left( \frac{2\pi(x-a)}{b-a} - \frac{\pi}{2} \right) \right]^n, \quad (3.9)$$

где  $m$  и  $n$  – параметры (положительные вещественные числа), а  $M$  – масштабный коэффициент, обеспечивающий равенство единице интеграла от  $f(x)$  на  $[a, b]$ . За пределами этого промежутка  $f(x) = 0$ . Значение  $M$  можно просто найти как величину, обратную интегралу от  $f(x)$  в пределах от  $a$  до  $b$  при  $M = 1$ .

Удобство использования распределений (3.9) состоит в том, что значениями  $m$  и  $n$  можно независимо друг от друга изменять асимметрию (параметр  $m$ ) и крутость (параметр  $n$ ) функций плотности. Рис. 3.12 поясняет влияние параметров  $m$  и  $n$  на вид функций  $f(x)$ . Симметричные распределения реализуются при  $m = 0$ . Распределения с отрицательной асимметрией приведены на рисунке при  $m = 2$  и  $m = 4$ . Для получения распределений с по-

ложительной асимметрией в выражении (3.9) нужно заменить величину  $x - a$  на  $b - x$ .



**Рис. 3.12.** Графики распределений типа  $x \sin x$  при  $m \in \{0, 2, 4\}$  для  $n = 1$  (слева) и  $n = 2$

При аппроксимации гистограммы моделируемых показателей подбор подходящих значений параметров  $a$  и  $b$  распределения (3.9), а при необходимости и промежутка  $[a, b]$ , наиболее просто выполняется путём моделирования графиков  $f(x)$  на графике аппроксимируемой ею гистограммы распределения рассматриваемого показателя методом проб и ошибок. В рассматриваемых в дальнейшем примерах выбор параметров распределения (3.9) производился именно таким способом.

Для аппроксимации распределений случайной составляющей временных рядов (ВР) показателей рождаемости и смертности, представленных на рис. 3.7 и рис. 3.8, использовались все четыре рассмотренные выше распределения. При этом распределения ПЗ, полученные на основе бета-распределения и нормального распределения, практически совпали. Распределения, ПЗ, полученные на основе распределения взвешенной суммы равномерно распределенных случайных величин (СВ), оказались очень близкими к двум первым распределениям (на рис. 3.7 и 3.8 они не приведены). Во всех случаях принималось  $M(X) = 0$ . В табл. 3.3 приведены параметры всех распределений, аппроксимирующих статистические распределения указанных ПЗ. Зна-

чения  $a$  и  $b$  этих распределений определялись по значениям  $\sigma$  статистического распределения при нулевом значении его математического ожидания. Математические ожидания  $M(X)$  при всех приведённых распределениях также равны нулю.

Т а б л и ц а 3.3. Параметры распределений, использованных для моделирования случайной составляющей временных рядов показателей рождаемости и смертности (рис. 3.7 и 3.8).

Аппроксимируемое распределение	Моделирующее распределение	Параметры моделирующих распределений
Случайной составляющей ВР показателя рождаемости для РФ ( $\sigma = 1,819$ )	Бета-распределение	$a = -4,812, b = 4,812, m = 3, n = 3, \sigma = 1,819$
	Усечённое снизу нормальное	$a = -4,812, b = 4,812, \sigma = 1,819, k = 1,0775$
	Суммы равномерно распределённых СВ	$a = -5,100, b = 5,100, m = 0, \sigma = 1,700$
	Типа $x \sin x$	$a = -4,812, b = 4,812, m = 0, n = 1,297, M = 0,0938, \sigma = 1,819$
Случайной составляющей ВР показателя смертности для РФ ( $\sigma = 1,314$ )	Бета-распределение	$a = -3,941, b = 3,941, m = 5, n = 5, \sigma = 1,314$
	Усечённое снизу нормальное	$a = -3,941, b = 3,941, \sigma = 1,314, k = 1,0304$
	Суммы равномерно распределённых СВ	$a = -3,941, b = 3,941, n = 3, \sigma = 1,314$
	Типа $x \sin x$	$a = -3,941, b = 3,941, m = 0, n = 1,279, M = 0,1153, \sigma = 1,314$
Случайной состав-	Бета-распределение	$a = -5,516, b = 5,516, m = 3, n = 3, \sigma = 2,085$



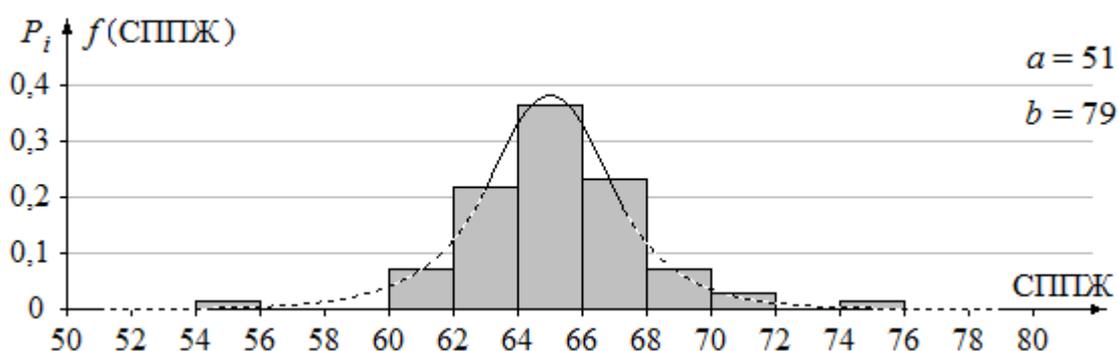
$M(\Pi_3)$	12,125	12,141	13,067	14,149	10,951	10,950	15,854	16,257
$\sigma(\Pi_3)$	3,365	3,168	2,152	2,362	3,264	3,060	5,496	5,445
$R(1)$	0,747	0,687	0,886	0,925	0,848	0,726	0,644	0,540
$R(2)$	0,469	0,455	0,730	0,772	0,699	0,590	0,300	0,298
$R(3)$	0,222	0,187	0,519	0,572	0,455	0,329	-0,059	0,082
$R(4)$	0,062	0,023	0,270	0,349	0,237	0,175	-0,348	-0,017
$R(5)$	-0,127	-0,181	0,044	0,157	0,023	-0,050	-0,534	-0,190
$R(6)$	-0,269	-0,285	-0,157	-0,027	-0,242	-0,289	-0,543	-0,329
$R(7)$	-0,294	-0,259	-0,226	-0,136	-0,377	-0,313	-0,371	-0,283
$R(8)$	-0,266	-0,134	-0,188	-0,169	-0,362	-0,110	-0,073	-0,119

Сравнивая приведённые в табл. 3.4 характеристики исходного и моделируемого временных рядов можно заключить, что они достаточно близки. По-видимому, допустимую степень их отклонения можно принимать за критерий, которому должны удовлетворять как моделируемый ряд в целом, так и его тренд и случайная составляющая. Например: максимально допустимое отклонение математического ожидания – не более 5%, среднего квадратического отклонения – до 10%,  $R(1)$  – до 15% и  $R(2)$  – до 20%.

Отдельное моделирование тренда и случайной составляющей позволило обойтись моделированием независимой случайной составляющей. В этом заключается существенное достоинство декомпозиции рядов. Не применяя такую декомпозицию ряда, пришлось бы решать сложную задачу моделирования коррелированной последовательности с заданным законом распределения.

В том случае, когда с помощью одной  $f(x)$  достаточно точно аппроксимировать гистограмму не удаётся, можно воспользоваться сплайновой аппроксимацией [44], при использовании которой промежуток  $[a, b]$  разбивается на отрезки  $[a_i, b_i]$ , и гистограммы на эти отрезках аппроксимируются функциями  $f_1(x), f_2(x), \dots, f_n(x)$ , являющимися звеньями сплайна. При моделировании каждого значения  $x$  вначале разыгрывается номер  $i$  звена сплайна,

причём вероятности выбора каждого отрезка должны быть равны значениям интегралов от  $f_i(x)$  на этих отрезках. Затем моделируется  $x$  по выбранной  $f_i(x)$ . Такая аппроксимация и моделирование с её использованием в целом являются более сложными. Вместе с тем такой метод позволяет интерполировать и соответственно моделировать не только колоколообразные, но и любые другие распределения показателей здоровья и работы учреждений здравоохранения. Наиболее распространённой является алгебраическая сплайновая интерполяция, в частности, полиномиальная, когда каждая  $f_i(x)$  является полиномом соответствующей степени [44, 60]. Если параметры каждой  $f_i(x)$  не зависят от параметров функций других звеньев, то такой сплайн называется локальным. Выбор параметров звеньев такого сплайна значительно проще [60] чем для сплайнов общего вида.



**Рис. 3.13.** Гистограмма распределения показателя средней продолжительности предстоящей жизни при рождении для населения РФ 2005 года рождения и её сплайновая аппроксимация трёхзвенным сплайном

На рис. 3.13 иллюстрируется результат аппроксимации гистограммы показателя СППЖ населения 70 административных единиц РФ, родившегося в 2005 г., трёхзвенным локальным сплайном. Колоколообразное распределение  $f(\text{СППЖ})$  на отрезке  $[51, 79]$  имеет математическое ожидание 65,3 года – такое же как и гистограмма. Оно разбито на три участка:  $[51, 63]$ ,  $[63, 67,3]$  и  $[67,3, 79]$ , на которых  $f_i(x)$  представляют собой следующие экспоненциальные функции:

$$f_1(x)=0,000303 \cdot [\exp(0,5467 \cdot (x-51))-1], \quad f_2(x)=0,382 \cdot \exp\left[-\frac{(x-65)^2}{6,919}\right],$$

$$f_3(x) = 0,17786 \cdot \exp[(67,3 - x) \cdot 0,598].$$

Незначительное отличие производных этих функций в точках стыковки звеньев сплайна, не заметные на рисунке, в рассматриваемых задачах вполне допустимы.

При моделировании каждого значения СППЖ или других случайных величин с распределением, представляемым сплайновой функцией, вначале разыгрывается номер звена сплайна. Вероятность выбора каждого  $i$ -го звена должна быть равна интегралу от  $f_i(x)$  на отрезке  $[a_i, b_i]$ , а сумма таких вероятностей должна быть равна единице. Затем моделируется  $x_i \in [a_i, b_i]$  - одним из методов, изложенных, например, в [38, 78, 103].

Рассмотренные колоколообразные распределения являются основными распределениями для моделирования показателей здоровья населения и случайной составляющей временных рядов этих показателей. Случайные величины, в том числе показатели здоровья, с другими распределениями могут быть реализованы с помощью методов, излагаемых в литературе по математическому моделированию [38, 78, 103 и др.].

## **ГЛАВА 4. МОДЕЛИ ИНТЕГРАЛЬНОГО ПОКАЗАТЕЛЯ ЗДОРОВЬЯ НАСЕЛЕНИЯ НА ОСНОВЕ “СТАНДАРТНЫХ” ПОКАЗАТЕЛЕЙ ЗДОРОВЬЯ**

### **4.1. Интегральная оценка здоровья населения, интегральные показатели**

Среди математических моделей в области здравоохранения значительное место занимают модели, связанные с интегральной оценкой здоровья населения. Попытки создания моделей интегральных показателей (ИП) здоровья населения, характеризующих в целом состояние здоровья населения различных регионов, ведутся уже в течение нескольких лет.

Для пояснения понятия “интегральная оценка здоровья населения” отметим ещё раз, что при оценке здоровья принято выделять 3 его уровня [75, 80, 95, 157]:

- здоровье отдельного человека (индивидуальное здоровье);
- групповое здоровье (здоровье социальных, этнических групп, населения административных территорий);
- общественное здоровье (здоровье общества, субпопуляции в целом).

В литературе характеристики группового и общественного здоровья в статике и в динамике рассматриваются как интегральные понятия индивидуального здоровья рассматриваемого множества индивидуумов, например, жителей некоторого региона. Причем эти характеристики являются не просто суммой соответствующих характеристик указанного множества индивидуумов, а совокупностью взаимосвязанных данных, выраженных количественными и качественными показателями соответственно построенной модели ИП здоровья. Математические описания зависимостей ИП группового или общественного здоровья от статистических показателей здоровья и задают модели соответствующих ИП.

На необходимость создания научно-обоснованных моделей ИП здоровья указывали академики Н.М. Амосов [3], Ю.П. Лисицин [74 – 77] и ряд дру-

гих авторитетных учёных [95, 102, 116, 158 и др.]. Так, по мнению Ю.П. Лищицына, “в различные периоды развития общества и, в частности, современной системы здравоохранения были и остаются актуальными вопросы оценки состояния здоровья, анализа обоснованности и силы влияния тех или иных факторов на формирование и сохранение определенного уровня общественного здоровья” Указанные показатели необходимы не только для удобного сравнения состояния здоровья населения различных регионов, но также и для принятия надлежащих мер по охране и укреплению здоровья. Отметим, что разработка ИП ведётся не только в медицине [84 и др.].

Вместе с тем анализ литературы по математическим моделям ИП здоровья [1, 11, 39, 42, 45, 66, 73, 102, 112, 116, 145, 151] показал, что еще нет модели или моделей, принятых большинством специалистов в качестве основных, наиболее подходящих для интегральных характеристик здоровья населения. В отдельных моделях имеются математические неувязки. Следовательно, проблема разработки и совершенствования указанных моделей оставалась актуальной. Поэтому авторами была продолжена работа по созданию и исследованию математических моделей ИП общественного здоровья населения.

При разработке новых моделей авторы старались устранить недостатки, присущие известным моделям, и исходили из требований Всемирной организации здравоохранения (ВОЗ) к ИП здоровья населения, сформулированным ещё в 1971-м году\*. Указанных требований – десять:

- Доступность данных. Должна существовать возможность для определения ИП без “сложных специальных исследований”.
- Полнота охвата. ИП должен быть получен из данных, охватывающих все население, для которого он предназначен.

---

\*) EUR/RC 47/II от 23.07.1997 г.// Всемирная организация здравоохранения.

- **Качество.** Национальные или территориальные данные не должны изменяться во времени и пространстве таким образом, чтобы на ИП оказывалось значительное влияние.
- **Универсальность.** ИП по возможности должен быть отражением группы факторов, которые определены и влияют на уровень здоровья.
- **Вычислимость.** ИП должен рассчитываться как можно более простым способом, расчет не должен быть дорогостоящим.
- **Приемлемость и интерпретируемость.** ИП должен быть приемлем и, несомненно, должны существовать приемлемые методы для расчета ИП и его интерпретации.
- **Воспроизводимость.** При использовании ИП здоровья разными специалистами в различных условиях и в различное время результаты должны быть идентичными.
- **Специфичность.** ИП должен отражать изменения только в тех явлениях, выражением которых она служит.
- **Чувствительность.** ИП здоровья должен быть чувствительным к изменениям соответствующих явлений.
- **Валидность.** ИП должен быть истинным выражением фактов, мерой которых он является.

По-видимому, некоторые из перечисленных требований уже не являются актуальными. Так, в настоящее время в медицинских учреждениях достаточно широко используется современная компьютерная техника, которая обладает очень высоким быстродействием и весьма большой памятью. Поэтому понятие сложности вычислений интегральных показателей уже существенно отличается от этого понятия в 1971-м году, запрограммированные вычисления ЭВМ выполняет за ничтожные доли секунды. Однако в целом перечисленные требования по-прежнему являются важной составляющей методологии разработки интегральных показателей здоровья населения.

В 1980-м году С.П. Ермаковым было предложено дополнить приведенный перечень требований ВОЗ еще тремя требованиями, которые повышают обоснованность использования разработанных ИП:

- Репрезентативность. ИП должен отражать изменения в здоровье отдельных возрастно-половых и других контингентов населения, выделенных для целей изучения.
- Иерархичность. ИП должен конструироваться по единому принципу для разных иерархических уровней, выделяемых в изучаемой совокупности населения для учитываемых заболеваний, их стадий и последствий. Должна существовать возможность его унифицированной свертки и развертки по составляющим компонентам.
- Целевая состоятельность. ИП должен адекватно отражать цели сохранения и развития (улучшения) здоровья и стимулировать общество к поиску наиболее эффективных путей достижения этих целей.

Эти требования в предлагаемых авторами моделях также учитываются. Кроме того, ввиду отсутствия единого понятия общественного здоровья и общепринятой методологии его оценки авторы исходили из общей концепции общественного здоровья, воплощённой в Европейской стратегии “Здоровье для всех на двадцать первое столетие”<sup>\*)</sup>. В этой концепции здоровье рассматривается как состояние, позволяющее вести активную в социальном и в экономическом плане жизнь. Однако рекомендации о том, что следует считать мерой указанной активности и как ИП общественного здоровья должен зависеть от соответствующих параметров, т.е. показателей здоровья, в документах ВОЗ и в указанной концепции отсутствуют.

При разработке моделей ИП, авторы выдвинули ещё одно требование к ИП:

- Удобство для использования. Значения ИП должны изменяться в нормированном промежутке, например от нуля до единицы.

---

<sup>\*)</sup> EUR/RC 47/II от 23.07.1997 г.// Всемирная организация здравоохранения.

Указанный промежуток и был принят авторами. При этом случай ИП = 0 соответствует наихудшему состоянию здоровья населения, например, – по головной заболеваемости и смерти всех индивидуумов рассматриваемого региона. Случай ИП = 1 соответствует максимально достигаемому состоянию здоровья. ИП является безразмерной величиной.

#### 4.2. Однопараметрические модели

Наиболее простыми моделями ИП являются однопараметрические модели. При этом, по мнению авторов, из всех известных ПЗ для использования в однопараметрических моделях ИП группового или общественного здоровья наибольшего внимания заслуживают следующие показатели: общий коэффициент смертности (ОКС) [73], средняя продолжительность жизни (СПЖ) [62] и особенно средняя продолжительность предстоящей жизни при рождении (СППЖ). Первый и последний из этих показателей входят в перечень медицинских ПЗ, публикуемых ежегодно ГОСКОМСТАТОм РФ.

Показатель ОКС для каждого региона определяется путём деления числа умерших жителей региона от всех причин за рассматриваемый период (обычно – год) на численность населения региона и умножения на 1000, то есть он указывает число смертных случаев на 1000 жителей региона. По этому показателю в определённой степени можно судить о состоянии санитарного обслуживания населения, о специфике природно-климатических, социально-экономических и экологических условий жизни, о региональной и возрастной патологии отдельных общественных групп населения [73]. При более углублённом изучении причин смертности рассматривают отдельно смертность от различных болезней, смертность по различным возрастным группам населения, смертность от несчастных случаев, смертность, связанную с социально-экономическими условиями жизни.

Значение показателя СПЖ определяется как средний возраст умерших в рассматриваемом году, т.е. путём деления суммы возрастов всех умерших

соответствующего региона в этом году на число умерших. СПЖ также зависит от указанных факторов, которые в общем случае оказывают на неё противоположное влияние, что свидетельствует о некоторой отрицательной корреляции между смертностью и СПЖ.

Если для построения ИП здоровья населения использовать показатель общей смертности, то ИП можно определять согласно выражению:  $ИП = 1 - ОКС/1000$ . В этом случае ИП будут зависеть линейно от ОКС, причём при  $ОКС = 0$   $ИП = 1$  (максимально возможное значение), а при  $ОКС = 1$   $ИП = 0$  (минимально возможное значение ИП). Однако в этом случае большинство значений ИП будет близким к 0,98, что не удобно для практического использования. Поэтому для определения ИП на основе показателя ОКС можно предложить выражение

$$ИП = 1 - K \cdot ОКС/1000, \quad (4.2)$$

в котором  $K$  – весовой коэффициент, выбираемый из [25, 35]. Тогда значения этого ИП будут близки к значениям, получаемым согласно рассматриваемым далее моделям. Однако при этом теоретически значения данного ИП при очень высокой смертности могут быть и отрицательными.

Для однопараметрической модели ИП на основе СПЖ авторы предлагают принять:

$$ИП = K \cdot СПЖ / СПЖ_{\max}, \quad (4.3)$$

где в качестве максимально возможного значения СПЖ принимается соответствующая величина, например, 100 или 115 лет, а  $K \in [0,95, 1,05]$ . В этом случае также  $ИП \in [0, 1]$ .

Средняя продолжительность предстоящей жизни при рождении представляет собой математическое ожидание числа лет, которое предстоит прожить поколению родившихся и не умерших в рассматриваемом году при условии, что на всём протяжении жизни индивидуумов этого поколения смертность в каждой возрастной группе будет такой же, какой она была в этом го-

ду. Обычно значения показателя СППЖ несколько отличаются от значений показателя СПЖ.

Известно несколько моделей (вариантов) расчета СППЖ [28, 103, 102]. В последние годы для расчёта СППЖ рассматриваемое поколение делится на возрастные группы  $[l_{i-1}, l_i)$ ,  $[l_i, l_{i+1})$ , ..., в качестве которых используются интервалы  $[0, 2)$ ,  $[2, 3)$ ,  $[3, 4)$ ,  $[4, 5)$ ,  $[5, 10)$ ,  $[10, 15)$  ... лет. Затем для каждой возрастной группы вычисляются повозрастные показатели смертности, равные отношению числа умерших из  $i$ -й группы в рассматриваемом году к численности населения данного возраста на начало этого года. По указанным показателям определяются статистические вероятности  $q_i$  дожить родившемуся до возраста  $l_i$  и умереть в возрасте, соответствующем  $i$ -й возрастной группе. Упорядоченное множество вероятностей  $q_i$  в здравоохранении принято называть таблицей смертности (дожития). Используя указанные вероятности и полагая, что в каждой из указанных групп средний возраст умирающих равен  $0,5(l_i + l_{i+1})$ , т.е. соответствует середине  $i$ -й группы, получаем:

$$\text{СППЖ} = 0,5 \sum q_i (l_i + l_{i+1}) \quad (4.4)$$

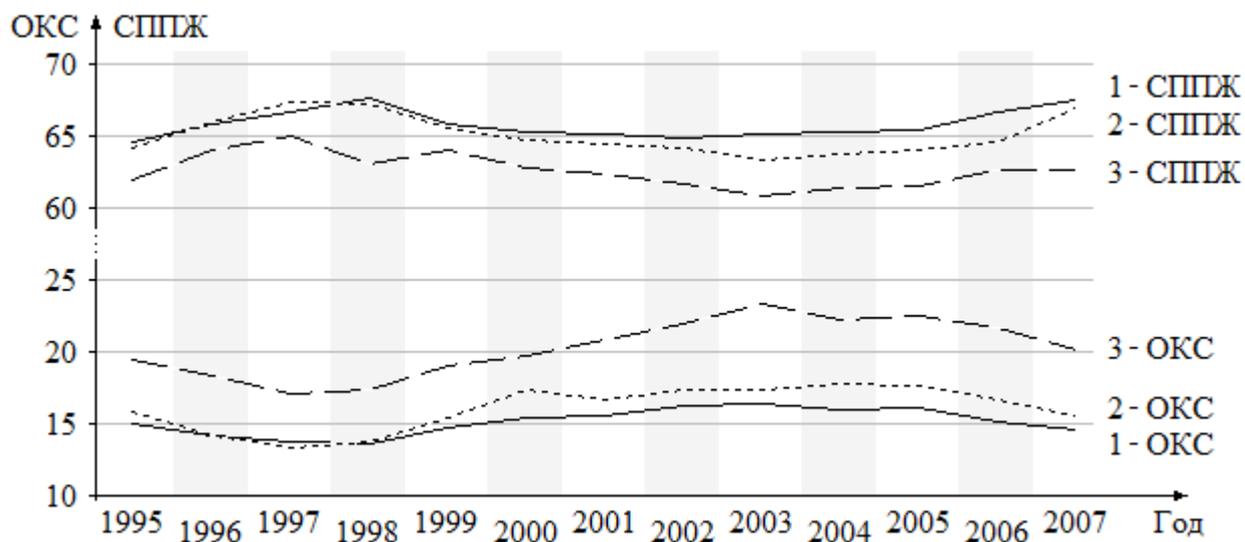
ИП на основе СППЖ можно определять аналогично выражению (4.3), а именно:

$$\text{ИП} = K \cdot \text{СППЖ} / \text{СППЖ}_{\text{макс}}, \quad (4.5)$$

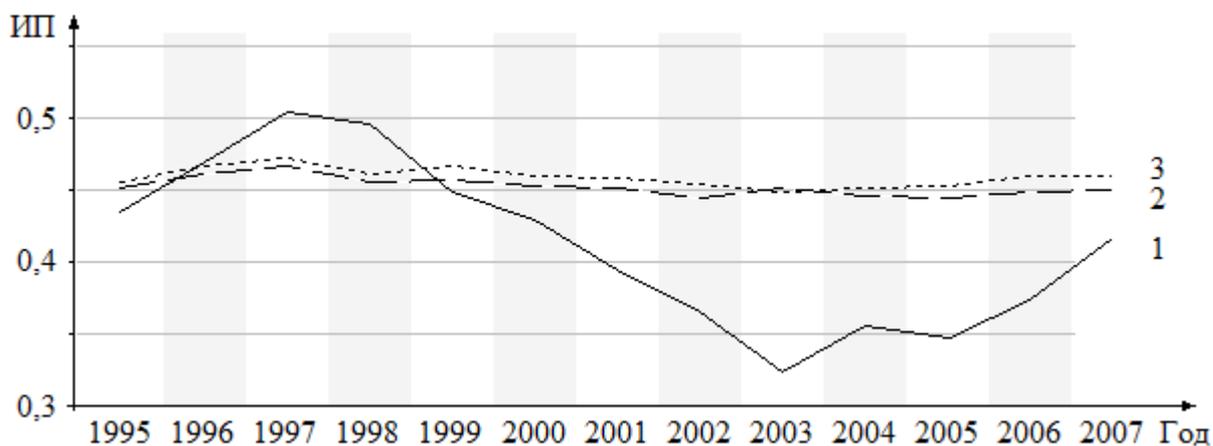
где в качестве  $\text{СППЖ}_{\text{макс}}$  принимается 100 или 115 лет, а  $K \in [0,6, 0,7]$ . Значения промежутков для выбора весового коэффициента  $K$  рекомендуются на основе моделирования ИП здоровья. При использовании рассматриваемых моделей для сравнения здоровья населения разных регионов следует использовать одну и ту же модель, т.е. значения  $K$  в этой модели должны быть одинаковыми.

На рис. 4.1 приведены графики изменения используемых ПЗ, а на рис. 4.2 – графики изменения ИП, полученных согласно выражениям (4.2) при  $K = 29$ , (4.3) при  $K = 0,65$  и  $\text{СППЖ}_{\text{макс}} = 100$  лет, (4.5) при  $K=0,65$  и  $\text{СППЖ}_{\text{макс}} =$

100 лет. Из графиков следует, что наибольший разброс значений имеет ИП на основе показателя ОКС, а наименьший – на основе СППЖ.



**Рис. 4.1.** Динамика статистических показателей общей смертности на 1000 человек населения и средней продолжительности предстоящей жизни при рождении для РФ (1), Северо-Западного федерального округа РФ (2) и Новгородской области (3)



**Рис. 4.2.** Динамика интегральных показателей здоровья населения Новгородской области на основе статистических показателей общей смертности на 1000 человек населения (1), средней продолжительности жизни (2) и средней продолжительности предстоящей жизни при рождении (3)

Средняя продолжительность жизни является итоговыми показателями жизни человека, она может неоднозначно зависеть от таких промежуточных факторов в его жизни как болезни, обращаемость к медицинской помощи, бытовые условия и др. Для интегральной оценки здоровья важно учитывать не столько среднюю продолжительность жизни, сколько продолжительность

здоровой жизни. Поэтому указанный показатель можно рекомендовать для построения ИП на его основе лишь для быстрой, прикидочной интегральной оценки здоровья населения, не требующей высокой точности определения этого показателя. Для обеспечения достаточно адекватной характеристики здоровья населения с помощью ИП модель ИП следует строить на основе нескольких, наиболее значимых статистических показателей общественного здоровья.

### **4.3 Структура многопараметрической модели интегрального показателя общественного здоровья населения**

В общем случае ИП здоровья населения является некоторой функцией от нескольких ПЗ. Следовательно, его тоже можно отнести к статистическим показателям. Указанная функция, которая может быть линейной или нелинейной, представляет собой аналитическую модель ИП.

Анализ причинно-следственных связей основных факторов, влияющих на показатели здоровья населения, показывает, что большинство этих показателей в значительной степени зависит от качества окружающей среды в рассматриваемом регионе, от эффективности работы системы здравоохранения, образа жизни населения, социально-экономических условий жизни и наследственности. При этом на состояние окружающей среды значительное влияние оказывает образ жизни населения, а на состояние системы здравоохранения больше всего влияют социально-экономические условия жизни. Кроме того, как будет показано в § 4.3, при построении моделей ИП из большого количества ПЗ следует учитывать только ПЗ, оказывающие на значения ИП существенное влияние. Так, для интегральной характеристики здоровья населения не имеют существенного значения разновидности заболеваний населения. Поэтому заболеваемость достаточно разделять лишь на общую по обращениям и на исчерпанную, а также по возрастному принципу.

Исходя из изложенного при построении моделей ИП общественного здоровья населения авторы предлагают исходить из предлагаемой ими и приводимой на рис. 4.3 структуры выявленных причинно-следственных связей основных природных и общественных факторов, влияющих на ПЗ, и основных ПЗ, определяющих ИП общественного здоровья населения [63].



**Рис. 4. 3.** Причинно-следственные связи основных природных и общественных факторов, влияющих на показатели здоровья. Структура многопараметрической модели интегрального показателя общественного здоровья населения

Приведённые на рис. 4.3 показатели, характеризующие состояние здоровья населения, являются параметрами ИП. Все или часть этих показателей используются соответствующим алгоритмом для вычисления значения ИП. Максимальное число таких ПЗ в приведённой на этом рисунке равно 9. Однако число параметров модели может быть сокращено или увеличено. Так, можно учитывать по возрастной смертность, а заболеваемость ещё разделять на острую и хроническую. Для определения значений истощенной заболеваемости и значения индекса физического состояния необходимо проводить

обследование населения. Но эти показатели, важные для интегральной оценки здоровья населения, в государственной статистике не приводятся.

Во всех случаях, когда по каким-либо причинам перечень ПЗ, учитываемых интегральным показателем, необходимо сократить или увеличить, принцип построения предложенной структуры модели ИП не изменяется. Если при этом не учитывать показатели физического развития и исчерпанную заболеваемость, то получим структурную схему модели ИП, которая использует только ПЗ, публикуемые в настоящее время ГОСКОМСТАТОм РФ (государственная статистика).

Перечень учитываемых моделями ПЗ с одной стороны определяется стремлением учесть большинство наиболее существенных показателей, влияющих на интегральную оценку здоровья населения, а с другой – не усложнять модели и обеспечить достаточное влияние каждого из использованных в модели ПЗ на значения ИП. Так, в силу указанного в рекомендуемые модели не включены показатели заболеваемости по различным видам болезней, а используются только показатели общей заболеваемости.

Заметим, что структурная схема модели ещё не является моделью. Для того, чтобы она стала таковой, необходимо раскрыть характер влияния каждого входящего в неё показателя на ИП, то есть построить или принять конкретный алгоритм, преобразующий значения ПЗ в значение ИП.

Рассматриваемые математические модели ИП общественного здоровья населения являются аналитическими, поскольку значения ИП на каждом шаге работы модели определяются значениями ПЗ, полученными на том же шаге. По виду зависимостей, характеризующих влияния на ИП учитываемых им показателей, математические модели ИП можно разделить на линейные и нелинейные.

#### **4.4. Линейные многопараметрические модели**

Согласно линейным моделям ИП его значение определяется алгебраической суммой учитываемых моделью показателей здоровья, умноженных на

соответствующие коэффициенты, называемые весовыми или весами. Достоинством линейных моделей является то, что в них приращения ИП пропорциональны приращениям каждого из указанных ПЗ. Величина приращения ИП при получении любым из этих ПЗ некоторого приращения не зависит от конкретного значения этого ПЗ, а зависит только от его приращения.

В обобщённом виде любую из предлагаемых моделей ИП, имеющую  $m$  параметров, можно задать выражением

$$\text{ИП} = A - B + C = \sum_{i=1}^j K_i \text{ПЗ}_i - \sum_{i=j+1}^m K_i \text{ПЗ}_i + C, \quad (4.6)$$

в котором  $A$  и  $B$ , являются обобщёнными составляющими модели, а именно – суммами произведений  $K_i \text{ПЗ}_i$  где  $K_i$  – весовые коэффициенты. При этом в произведения суммы  $A$  входят те  $\text{ПЗ}_i$ , увеличение которых приводит к увеличению значения ИП (индекс физического состояния, рождаемость, средняя продолжительность предстоящей жизни при рождении), а в произведения суммы  $B$  – те  $\text{ПЗ}_i$ , увеличение которых приводит к уменьшению значения ИП (заболеваемость, инвалидность, смертность). Величина  $C$  в выражении (4.6) является константой (постоянным числом). Она определяет среднее значение ИП в модели выбранного вида.

Весовые коэффициенты (веса) показателей здоровья в предлагаемых моделях выбираются таким образом, чтобы  $\text{ИП} \in [0, 1]$ . Для обеспечения этого условия учтём, что каждое произведение  $K_i \text{ПЗ}_i$  является вкладом  $V_{k_i}$   $i$ -го ПЗ в сумму  $A$  (при  $i \leq j$ ) или в сумму  $B$  (при  $i > j$ ), т.е. положительным или отрицательным приращением  $\Delta \text{ИП}_i$ , получаемым интегральным показателем при изменении только  $i$ -го ПЗ. В дальнейшем средние значения таких вкладов используются в разработанной методике определения значений весовых коэффициентов.

Разумно варьируя перечнем используемых в модели показателей здоровья и весовыми коэффициентами, на основе выражения (4.6) можно получить много моделей ИП общественного здоровья, соответствующих рассмот-

ренной структуре ИП (рис. 4.3). Некоторые из таких моделей [63, 85, 88] использовались в системе здравоохранения Новгородской области.

Рассмотрим четыре наиболее характерные линейные модели ИП, в двух из которых используются только ПЗ, ежегодно публикуемые ГОСКОМСТАТОМ РФ (табл. 4.1).

Т а б л и ц а 4.1. Линейные модели интегрального показателя общественного здоровья населения

№	Реализуемые зависимости
1	$\text{ИП} = K_{\text{ОКР}} \text{ОКР} + K_{\text{СППЖ}} \text{СППЖ} - K_{\text{ОЗОД}} \text{ОЗОД} - K_{\text{ОЗОВ}} \text{ОЗОВ} - K_{\text{ОКСД}} \text{ОКСД} - K_{\text{ОКСВ}} \text{ОКСВ} - K_{\text{ПИНВ}} \text{ПИНВ} + C_1$
2	$\text{ИП} = K_{\text{ОКР}} \text{ОКР} + K_{\text{СППЖ}} \text{СППЖ} - K_{\text{ОЗО}} \text{ОЗО} - K_{\text{ОКС}} \text{ОКС} - K_{\text{ПИНВ}} \text{ПИНВ} + C_2$
3	$\text{ИП} = K_{\text{ОКР}} \text{ОКР} + K_{\text{СППЖ}} \text{СППЖ} + K_{\text{ИФС}} \text{ИФС} - K_{\text{ОЗИД}} \text{ОЗОД} - K_{\text{ОЗИВ}} \text{ОЗОВ} - K_{\text{ОКСД}} \text{ОКСД} - K_{\text{ОКСВ}} \text{ОКСВ} - K_{\text{ПИНВ}} \text{ПИНВ} + C_3$
4	$\text{ИП} = K_{\text{ОКР}} \text{ОКР} + K_{\text{СППЖ}} \text{СППЖ} + K_{\text{ИФС}} \text{ИФС} - K_{\text{ОЗИ}} \text{ОЗИ} - K_{\text{ОКС}} \text{ОКС} - K_{\text{ПИНВ}} \text{ПИНВ} + C_4$

В приведенных моделях используются следующие обозначения:

- ОКР – общий коэффициент рождаемости (младенцев, родившихся живыми);
- СППЖ – средняя продолжительность предстоящей жизни (при рождении);
- ОЗО – общая заболеваемость населения по обращаемости в учреждения здравоохранения;
- ОЗОД – общая заболеваемость детского населения (от одного года до 17 лет включительно) по обращаемости;
- ОЗОВ – общая заболеваемость взрослого населения (с 18 лет) по обращаемости;
- ОЗИ – общая заболеваемость населения исчерпанная;
- ОЗИД – общая заболеваемость детского населения исчерпанная;
- ОЗИВ – общая заболеваемость взрослого населения исчерпанная;
- ОКС – общий коэффициент смертности населения;

- ОКСД – общий коэффициент смертности детского населения;  
 ОКСВ – общий коэффициент смертности взрослого населения;  
 ПИНВ – первичная инвалидность (общая);  
 ИФС – индекс физического состояния (см. § 3.2);  
 $K_i$  – весовой коэффициент (вес) показателя здоровья, указанного в индексе  $i$ ;  
 $C_i$  – число, определяющее среднее значение ИП для  $i$ -й модели.

В государственной статистике значения всех используемых в приведённых моделях ПЗ, кроме показателей СППЖ и ИФС, приводятся в расчёте на 1000, на 10000 или на 100000 человек соответствующего населения (всего, детского или взрослого). Для вычисления значения ИП это не имеет значения, так как при этом соответственно изменяются значения весовых коэффициентов ПЗ, учитываемых используемой моделью.

Первые две модели ИП общественного здоровья использует только ПЗ, приводимые в государственной статистике. Вторая модель является упрощённым вариантом первой модели. В этой модели в качестве показателей заболеваемости и смертности используются только показатели общей заболеваемости и общей смертности. Пример динамики ПЗ, входящих во 2-ю модель, приводится на рис. 4.4.

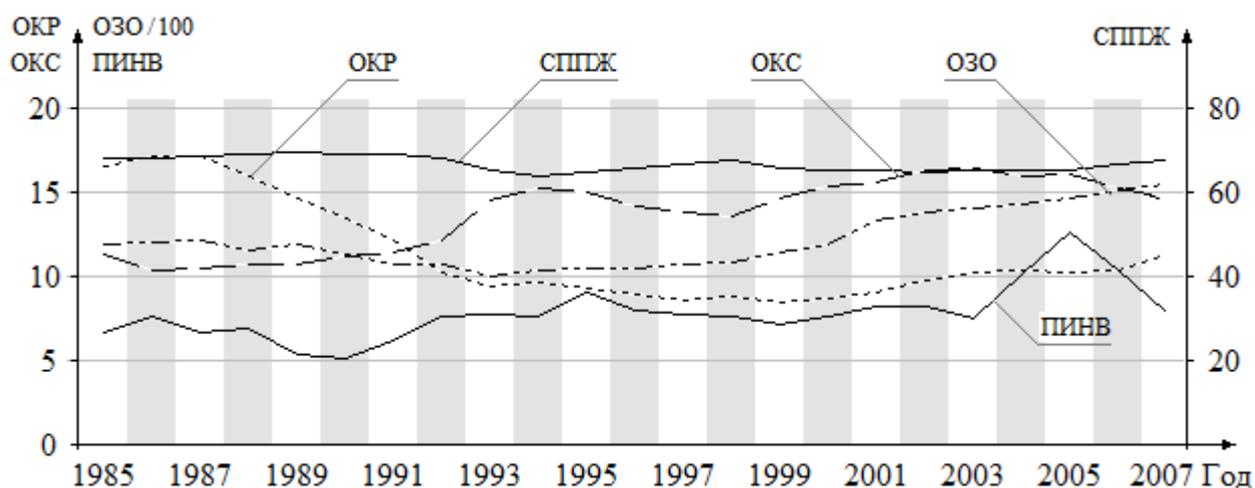


Рис. 4.4. Динамика показателей здоровья населения Российской Федерации

Третья модель является наиболее полной из моделей, приведённых в табл. 4.1. Она использует важные для интегральной оценки здоровья населения дополнительные ПЗ, которые ГОСКОМСТАТом РФ не публикуются. К таким ПЗ относится индекс физического состояния населения, являющийся интегральной характеристикой физического развития населения (п. 3.2). Кроме того, для улучшения интегральной оценки заболеваемости предлагается использовать в модели не показатели заболеваемости по обращаемости, а показатели исчерпанной (истинной) заболеваемости.

Дело в том, что значения всех показателей заболеваемости, вносимых в государственную базу данных, определяются по обращаемости за медицинской помощью, т.е. они не учитывают контингент больных, не обратившихся к врачу, и поэтому не могут достаточно адекватно отражать положение с заболеваемостью населения. Обычно процент не обратившихся к врачу больных устанавливается путём обследования соответствующей выборки населения рассматриваемого региона. Общую заболеваемость по обращениям, дополненную общей заболеваемостью не обратившихся к врачу больных, называют *исчерпанной* или *истинной* заболеваемостью [68, 95]. Согласно проводившимся в 1995-м и в 2005 годах исследованиям истинной заболеваемости городского населения Новгородской области [68] эта заболеваемость по различным классам болезней превысила заболеваемость по обращению в среднем в 2,5 раза в 1985 году и в 1,8 раза в 2005 году, т.е. в целом большинство болевших за врачебной помощью не обращались.

Четвёртая модель является упрощённым вариантом третьей модели. Пока две последние модели являются моделями на перспективу, так как статистический анализ всех входящих в них показателей здоровья проводится только в Новгородской области.

По-видимому, для регионов, не проводящих исследования исчерпанной заболеваемости и физического развития населения, для вычисления значений ИП с помощью хотя бы 4-й модели можно воспользоваться результатами

указанного исследования в других регионах. Так, для получения исчерпанной заболеваемости значения показателей заболеваемости этих регионов можно умножить на средний коэффициент, показывающий отношение исчерпанной заболеваемости к заболеваемости по обращениям для регионов с известной исчерпанной заболеваемостью.

Выбор весовых коэффициентов для рассматриваемых моделей ИП может быть достаточно просто произведён согласно предложенной в [52, 53] методике, основанной на выборе (на распределении) вкладов  $V_k$  показателей здоровья в ИП на усреднённой статистике ПЗ, в качестве которой предлагается использовать средние значения ПЗ населения РФ за несколько лет. Такие статистические данные приводятся в табл. 4.2, в которой значения показателей заболеваемости, смертности, инвалидности и рождаемости указаны в расчете на 1000 человек соответствующего населения.

Т а б л и ц а 4.2. Средние значения показателей здоровья населения РФ за 1985-2005г

Название показателя	Значения	Название показателя	Значения
Общий коэффициент рождаемости	9,309	Общая заболеваемость по обращаемости	1183,784
Общий коэффициент смертности населения	15,191	Общая заболеваемость по обращаемости детского населения	1526,697
Общий коэффициент смертности младенцев	15,936	Общая заболеваемость по обращаемости взрослого населения	1040,357
Общий коэффициент смертности детского населения	16,309	Первичная инвалидность	7,476
Общий коэффициент смертности взрослого населения	14,481	Средняя продолжительность предстоящей жизни	66,691

Выбор вкладов показателей здоровья в значение ИП в моделях нормированного ИП качества систем произвольного назначения должен быть произведён таким образом, чтобы с одной стороны ИП реагировал на минимальное, отличное от нуля приращение ПЗ, а с другой стороны не выходил бы из

промежутка  $[0, 1]$ . С увеличением числа подсистем рассматриваемой системы и с увеличением степени разброса значений параметров системы увеличивается и возможный разброс суммы вкладов  $V_{k_i}$ , равной  $A + B$ , и, соответственно, значений ИП. Согласно [49, 55, 58] для большинства систем среднее значение суммы  $A + B$ , гарантирующее изменение ИП в  $[0, 1]$ , может быть выбрано согласно эмпирическому выражению

$$A + B = 1,1/(1+n)^k, \quad (4.7)$$

в котором  $n$  – число подсистем в системе, не имеющих своих подсистем, а  $k$  – число из промежутка  $[0,2, 0,3]$ , зависящее от степени разброса параметров системы. Применительно к моделям ИП общественного здоровья населения системой является РФ, а подсистемами, не имеющими своих подсистем, являются административные единицы федеральных округов, т.е. края, республики, области, автономные округа (всего 71 административная единица). При этом для моделей рассматриваемого ИП рекомендуется [73] принять  $k = 0,24$ . Следовательно, для ИП общественного здоровья населения РФ в целом, федеральных округов и их административных образований выражение (4.7) принимает вид:

$$A + B = 1,1/72^{0,24}, \quad (4.8)$$

Если в приведённых выражениях принять  $n = 1$ , то такая модель будет предназначена только для РФ в целом. В этом случае  $(A + B)_{\text{макс}} = 1,1$ , и расчётное значение  $A + B$  выбирается принадлежащим  $[-0,05, 1,05]$  при ограничении ИП по значениям 0 и 1, что несколько увеличивает чувствительность ИП по сравнению со случаем выбора  $A + B \in [0, 1]$ . Если же разрабатываемую модель необходимо применять и для районов, значение  $n$  должно быть соответственно увеличено. В дальнейшем будут рассматриваться только модели ИП здоровья населения, предназначенные для РФ в целом, федеральных округов и их административных единиц на уровне краёв, областей, автономных республик и округов.

Предлагаемая методика расчёта весовых коэффициентов моделей ИП предусматривает следующее:

- Принимается, что среднее значение суммы  $A+B$  модулей обобщённых составляющих  $A$  и  $B$  в моделях ИП (4.6), получаемых на основе усреднённых значений ПЗ для РФ, определяется согласно выражению (4.8) и, следовательно, равно 0,4. Это гарантирует достаточно широкий "рабочий интервал" изменения ИП в  $[0, 1]$ , расчётные значения которого только теоретически с крайне малой вероятностью могут выйти за пределы указанного отрезка (случаи, когда, например, нет ни больных, ни инвалидов, а рождаемость и средняя продолжительность предстоящей жизни весьма велики, или когда, например, при отсутствии рождаемости почти всё население переболело и вымерло). Если какой-либо из указанных случаев всё же произойдёт, то принимается, что ИП равен соответственно 1 или 0. Путём уменьшения указанной суммы можно добиться гарантированного изменения ИП только в  $[0, 1]$ . Однако при этом существенно сужается интервал, в котором практически может изменяться ИП, что неудобно для его использования.
- С участием экспертов устанавливаются соотношения вкладов ПЗ в значение ИП для выбранной модели и по ним определяются значения этих вкладов и величины  $A$  и  $B$ . Так как для 2-й модели экспертами рекомендовано принять вклады  $V_{\text{ОКР}}$ ,  $V_{\text{ОЗО}}$  и  $V_{\text{ПИНВ}}$  равными 25% от  $A+B$ , а  $V_{\text{СППЖ}}$  и  $V_{\text{ОКС}}$  – равными 12,5% от  $A+B$ , т.е.:  $V_{\text{ОКР}} = V_{\text{ОЗО}} = V_{\text{ПИНВ}} = 0,1$  и  $V_{\text{СППЖ}} = V_{\text{ОКС}} = 0,05$ . Следовательно,  $A = V_{\text{ОКР}} + V_{\text{СППЖ}} = 0,15$  и  $B = V_{\text{ОЗО}} + V_{\text{ПИНВ}} + V_{\text{ОКС}} = 0,25$ . В табл. 4.3 приведены значения вкладов  $V_{k_i}$  для каждой из рассмотренных математических моделей ИП (табл. 4.1), а также значения обобщённых составляющих  $A$  и  $B$  этих моделей.

Т а б л и ц а 4.3. Рекомендуемые значения вкладов показателей здоровья в ИП общественного здоровья

Модель	Значения вкладов $V_{k_i}$	$A$	$B$
1	0,1 (ОКР), 0,05 (СППЖ), 0,025 (ОЗОД), 0,075 (ОЗОВ), 0,017 (ОКСД), 0,033 (ОКСВ), 0,1 (ПИНВ)	0,15	0,25
2	0,1 (ОКР), 0,05 (СППЖ), 0,1 (ОЗО), 0,05 (ОКС), 0,1 (ПИНВ)	0,15	0,25
3	0,08 (ОКР), 0,04 (СППЖ), 0,08 (ИФС), 0,02 (ОЗИД), 0,06 (ОЗИВ), 0,0135(ОКСД), 0,0265 (ОКСВ), 0,08 (ПИНВ)	0,2	0,2
4	0,08 (ОКР), 0,04 (СППЖ), 0,08 (ИФС), 0,08 (ОЗИ), 0,04 (ОКС), 0,08 (ПИНВ)	0,2	0,2

- С помощью полученных значений  $V_{k_i}$  и средних значений показателей здоровья определяются значения всех весовых коэффициентов для рассматриваемой модели. Для этого используется выражение  $V_{k_i} = K_i \cdot ПЗ_{i\text{ ср}}$ , из которого следует:

$$K_i = V_{k_i} / ПЗ_{i\text{ ср}}, \quad (4.9)$$

где  $ПЗ_{i\text{ ср}}$  – среднее значение соответствующего ПЗ из табл. 4.2.

- Вычисляются значения чисел  $C$ , определяющих местоположение интервала изменения ИП. Для этого рекомендуется принять, что для каждой модели середина интервала  $[A, B]$  соответствует значению ИП, равному 0,5, т.е.  $0,5A - 0,5B + C = 0,5$ . Тогда значение  $C$  определяется согласно выражению

$$C = 0,5 + 0,5(B - A).$$

В дальнейшем такие значения чисел  $C$  используются во всех рассматриваемых линейных моделях. Для первой и второй моделей из табл. 4.1 получаем  $C = 0,55$ , а для третьей и четвертой –  $C = 0,5$ .

Значения весовых коэффициентов для моделей ИП, приведённых в табл. 4.1, рассчитанные с соответствии с предлагаемой методикой их расчёта,

позволяют представить рассматриваемые модели с конкретными значениями весовых коэффициентов (табл. 4.4). Подставляемые в эти модели ПЗ должны быть рассчитаны на то же число жителей, что и ПЗ<sub>i</sub> в выражении (4.9). Для определения весовых коэффициентов показателей физического развития и исчерпанной заболеваемости необходимо знать средние значения этих показателей для Российской Федерации.

Т а б л и ц а 4.4. Линейные модели ИП с числовыми значениями весовых коэффициентов

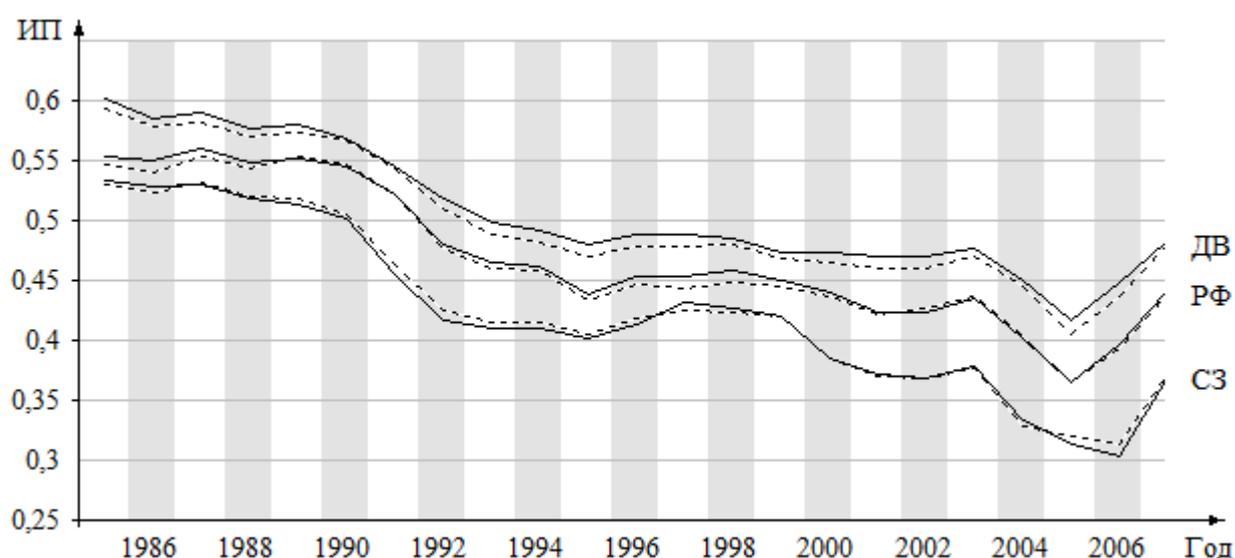
№	Реализуемые зависимости
1	$\text{ИП} = (10,742 \text{ ОКР} + 0,750 \text{ СППЖ} - 0,016 \text{ ОЗОД} - 0,072 \text{ ОЗОВ} - 0,920 \text{ ОКСД} - 2,417 \text{ ОКСВ} - 13,376 \text{ ПИНВ})/1000 + 0,55$
2	$\text{ИП} = (10,742 \text{ ОКР} + 0,750 \text{ СППЖ} - 0,084 \text{ ОЗО} - 3,291 \text{ ОСМ} - 13,376 \text{ ПИНВ})/1000 + 0,55$
3	$\text{ИП} = (8,594 \text{ ОКР} + 0,600 \text{ СППЖ} + K_{\text{ИФС}} \text{ ИФС} - K_{\text{Озид}} \text{ ОЗИД} - K_{\text{Озив}} \text{ ОЗИВ} - 0,828 \text{ ОКСД} - 1,830 \text{ ОКСВ} - 10,701 \text{ ПИНВ})/1000 + 0,5$
4	$\text{ИП} = (8,594 \text{ ОКР} + 0,600 \text{ СППЖ} + K_{\text{ИФС}} \text{ ИФС} - K_{\text{Ози}} \text{ ОЗИ} - 2,633 \text{ ОКС} - 10,701 \text{ ПИНВ})/1000 + 0,5$

Отметим, что нет необходимости ежегодно определять значения масштабных коэффициентов. Принятый для использования вариант модели с рассчитанными коэффициентами может быть рабочим в течение многих лет. Примеры динамики ИП общественного здоровья, рассчитанных согласно рекомендуемым линейным моделям, приведены на рис. 4.5 и рис. 4.6.

Рис. 4.5 иллюстрирует влияние на ИП разделения показателей ОЗО и ОСМ на две группы: детей и взрослых, что, по мнению авторов, позволяет более объективно оценивать здоровье населения с помощью ИП. По графикам для Дальневосточного и Северо-Западного федеральных округов, имеющих соответственно максимальное и минимальное средние значения ИП за 2000÷2005 г., можно также заключить, что ИП здоровья населения в

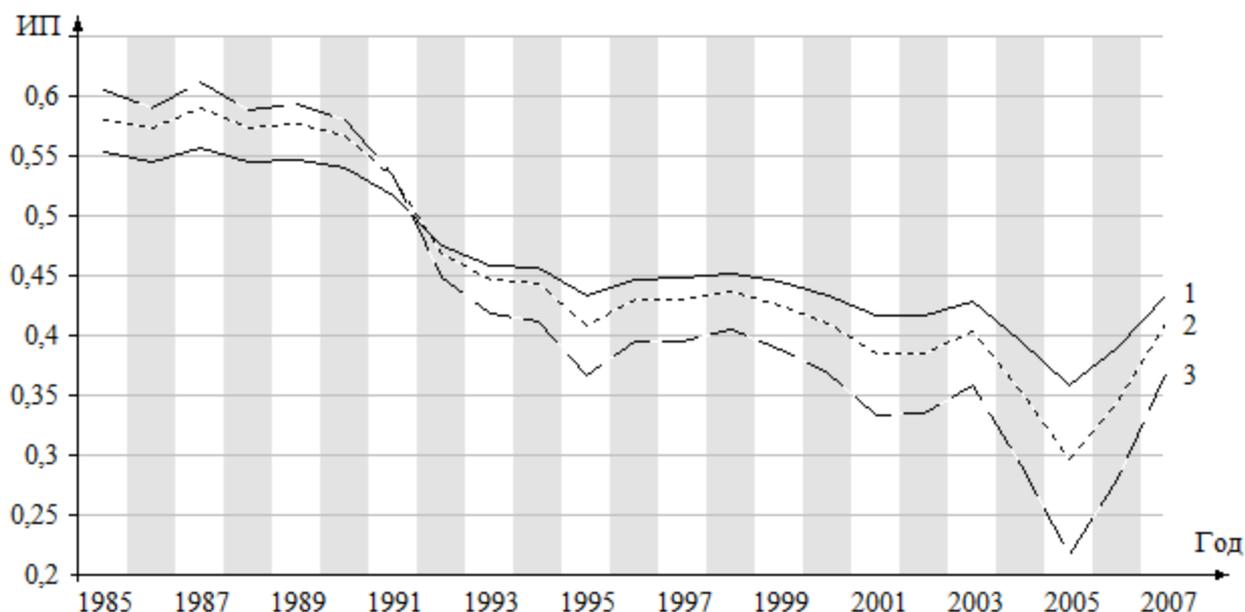
1985÷2006 г. изменялся в промежутке  $[0,34, 0,62]^*$ . На рис. 4.6 демонстрируется влияние на ИП, получаемый на основе второй модели, значения суммы  $A + B$  обобщённых составляющих модели ИП и разделения её на вклады ПЗ, соотношение которых не изменяется.

По приведённым на рис. 4.5 и рис. 4.6 графикам можно заключить, что при применении различных моделей ИП общественного здоровья имеют место одинаковые тенденции изменения ИП. Следовательно, временные ряды ИП, порождаемые этими моделями, имеют значительную положительную взаимную корреляцию. Так, для первой и второй моделей (табл. 4.5) на выборке ИП здоровья населения России за 1985 ÷ 2007 годы коэффициент взаимной корреляции получаемых ИП равен 0,972. Отметим также, что при увеличении суммы модулей положительной ( $A$ ) и отрицательной ( $B$ ) составляющих ИП пропорционально указанному увеличению расширяется и интервал изменения ИП (рис. 4.6). При этом увеличивается вероятность достижения интегральным показателем граничных значений интервала  $[0, 1]$  его изменения.



**Рис. 4.5.** Динамика интегральных показателей здоровья населения Российской Федерации, Дальневосточного и Северо-Западного ФО, полученных на основе моделей 1 (точечные линии) и 2 (сплошные линии)

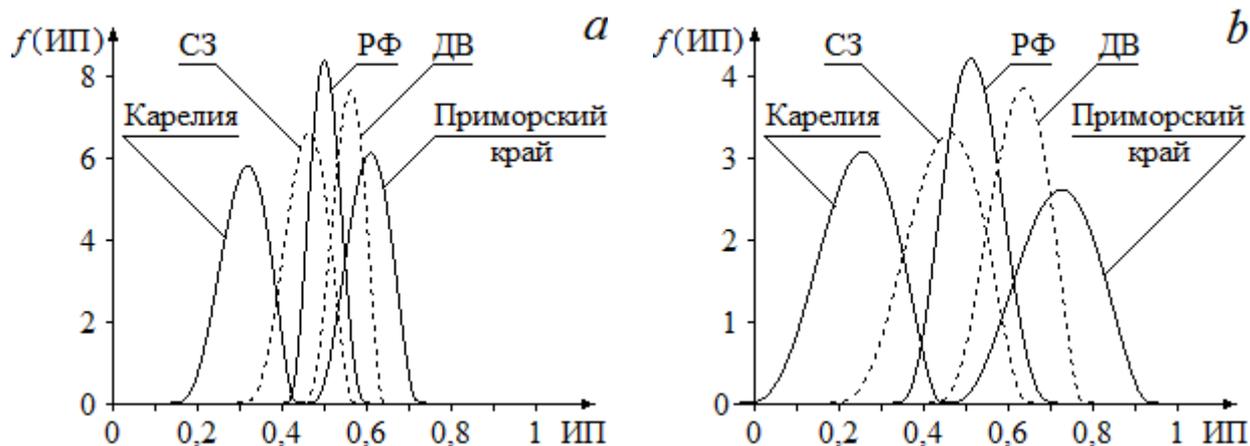
\* Так как Федеральные округа были образованы в 2000г., то ПЗ их населения за 1985÷1999 г. определялись по ПЗ населения включённых в них административных единиц.



**Рис. 4.6.** Динамика интегральных показателей здоровья населения Российской Федерации, полученных на основе 2-й модели при значениях  $A + B$ , равных 0,4 (1), 0,6 (2) и 0,8 (3)

достижения интегральным показателем граничных значений интервала его изменения.

Графики функций плотности  $f(\text{ИП})$ , приведённые на рис. 4.7 слева и полученные с использованием статистических ПЗ за 1996 ÷ 2006 годы, свидетельствуют о том, что при рекомендованной для расчёта весовых коэффициентов модели ИП сумме обобщённых составляющих модели  $A+B = 0,4$  предложенная методика расчёта этих коэффициентов обеспечила определённый запас для изменения ИП до границ промежутка  $[0, 1]$ . Получаемые границы изменения ИП для всех регионов России определялись по минимальному значению ИП здоровья населения Республики Карелия, имевшей это значение ИП, и по максимальному значению ИП здоровья населения Приморского края, по которому был зарегистрирован указанный максимум. При  $A+B = 0,4$  ИП здоровья населения всех регионов России изменялся от 0,15 до 0,72.



**Рис. 4.7.** Аппроксимированные статистические распределения ИП здоровья населения России, Северо-Западного и Дальневосточного федеральных округов и двух их административных единиц, полученные на основе 2-й модели при  $A+B = 0,4$  (a) и  $A+B = 0,8$  (b)

Графики функций плотности  $f(\text{ИП})$ , приведённые на рис. 4.7 справа, отражают случай, когда выбрано  $A + B = 0,8$ . В этом случае ИП здоровья населения каждого региона изменялся в более широких пределах. Минимальное значение ИП здоровья населения всех регионов России оказалось равным  $-0,01$  (ограничение этого значения даёт 0), а максимальное  $-0,93$ . Поэтому при использовании ИП, в модели которого принято  $A+B = 0,8$ , время от времени может иметь место выход расчётных значений ИП за пределы промежутка  $[0, 1]$ .

В рассмотренных моделях в число ПЗ, использованных для расчета ИП, входил показатель заболеваемости по обращениям. При замене этой заболеваемости на исчерпанную заболеваемость, которая может быть больше примерно в 2 раза, естественно произойдёт определённое уточнение значений ИП. Однако в этом случае существенных изменений ИП не произойдёт, так как в выражении (4.9) при неизменном значении принятого вклада заболеваемости в значение ИП с увеличением показателя заболеваемости соответственно уменьшится весовой коэффициент  $K_{\text{ози}}$ . Поэтому в регионах, в которых при переходе на учёт исчерпанной заболеваемости возрастание забо-

леваемости произойдёт пропорционально её возрастанию в РФ, а ИП здоровья населения не изменится. В противном случае произойдёт некоторое изменение ИП в ту или иную сторону.

Как уже указывалось, предложенная методика определения весовых коэффициентов с близкой к нулю вероятностью ограничения ИП допускает достижение интегральным показателем значений 0 и 1. Для оценки значений ПЗ, при которых указанное событие всё же может произойти, воспользуемся 2-й моделью. Рассмотрим случаи, когда значения величин  $A$  и  $B$  в выражении (4.6) изменятся по сравнению с приведёнными в табл. 4.3 в противоположные стороны, в  $K_0$  или в  $K_1$  раз, и при этом ИП достигнет соответственно значения 0, если  $A$  уменьшится, а  $B$  увеличится в  $K_0$  раз, или значения 1, если  $A$  увеличится, а  $B$  уменьшится в  $K_1$  раз. Определим соответствующие этим событиям значения  $K_0$  и  $K_1$ .

Из выражения (4.6) для 2-й модели с учётом данных табл. 4.3 следует:

$$\text{1-й случай: ИП} = A/K_0 - BK_0 + C = 0,15/K_0 - 0,25K_0 + 0,55 = 0, \text{ т.е. } K_0 \approx 2,43.$$

$$\text{2-й случай: ИП} = AK_1 - B/K_1 + C = 0,15K_1 - 0,25/K_1 + 0,55 = 1, \text{ т.е. } K_1 \approx 2,76.$$

Очевидно, указанное одновременное изменение обобщённых составляющих модели ИП имеет крайне низкую вероятность.

Таким образом, достижение таких значений ПЗ, при которых ИП станет равным нулю или единице, является почти невозможным событием. Вместе с тем можно показать, что предложенная методика определения весовых коэффициентов позволяет в 3÷5 раз расширить рабочий интервал изменения ИП по сравнению с методикой, использовавшейся в [85, 88, 152] и предполагающей равенство единице суммы этих коэффициентов.

#### 4.5. Нелинейные многопараметрические модели

Целью введения нелинейности в зависимость ИП общественного здоровья от соответствующих показателей является возможность использования

простого способа обеспечения нормированности ИП. При этом в отличие от линейных моделей можно обеспечить изменение ИП в промежутке  $[0, 1]$  при любых положительных значениях весовых коэффициентов, что является достоинством предлагаемых нелинейных моделей.

В обобщённом виде предлагаемые нелинейные модели можно представить в виде

$$\text{ИП} = \frac{A}{A + K_M B}, \quad (4.10)$$

где  $A$  и  $B$  – обобщённые параметры модели [88], определяемые согласно суммам в выражении (4.6). Как и в линейных моделях, параметр  $A$  является суммой произведений ПЗ, с ростом которых ИП увеличивается, на соответствующие им весовые коэффициенты, а параметр  $B$  определяется как сумма произведений ПЗ, с ростом которых ИП уменьшается, на их весовые коэффициенты. Причём эта сумма умножается ещё на масштабный коэффициент  $K_M$ , от выбора которого изменяется влияние параметра  $B$  на ИП, т.е. ИП по-разному изменяется в  $[0, 1]$  при одном и том же изменении ПЗ.

Если в выражении (4.10)  $A = 0$ , а  $B > 0$ , то  $\text{ИП} = 0$ . Если же  $A > 0$ , а  $B = 0$ , то  $\text{ИП} = 1$ . Особым случаем, имеющим чисто теоретическое значение, является случай, когда  $A = B = 0$ . Практически это может иметь место только в административных единицах с крайне малой численностью населения (например,  $10 \div 20$  человек). В данном случае согласно выражению (4.9) получаем:  $\text{ИП} = 0/0$  (неопределённость). Предлагается принять, что при этом  $\text{ИП} = 0.5$ . Такое значение ИП можно рассматривать как предельное, если  $A$  и  $B$  одновременно стремятся к нулю при выполнении условия  $K_M B/A = 1$ . Поэтому в рассматриваемом случае получаем:

$$\lim_{A, B \rightarrow 0} \text{ИП} = \lim_{A, B \rightarrow 0} \frac{1}{1 + K_M B/A} = \frac{1}{1+1} = 0,5.$$

Отметим, что значение ИП = 0,5 при  $A = B = 0$  можно получить, прибавив к числителю и знаменателю выражения (4.9) малые постоянные величины с отношением 1/2, например, 0,0001 и 0,0002. Ввиду малости эти величины не будут влиять на ИП. Такую возможность следует иметь в виду при определении значений ИП для административных единиц с малой численностью населения.

Как и в случае линейных моделей, на основе общей структуры модели ИП (рис. 4.3) можно привести большое количество нелинейных моделей, соответствующих выражению (4.10). В табл. 4.5 приведены 4 нелинейные модели ИП, которые по используемым параметрам аналогичны соответствующим линейным моделям ИП, приведённым в табл. 4.1 (одинаковые ПЗ используют следующие пары моделей: 1-я и 5-я, 2-я и 6-я, 3-я и 7-я, 4-я и 8-я).

Т а б л и ц а 4.5. Нелинейные модели ИП общественного здоровья

№	Алгоритм (модель)
5	$A = K_{\text{окр}} \text{ОКР} + K_{\text{СППЖ}} \text{СППЖ},$ $B = K_{\text{ОЗОД}} \text{ОЗОД} + K_{\text{ОЗОВ}} \text{ОЗОВ} + K_{\text{ОКСД}} \text{ОКСД} + K_{\text{ОКСВ}} \text{ОКСВ} + K_{\text{ПИНВ}} \text{ПИНВ}$ $\text{ИП} = A / (A + K_M B)$
6	$A = K_{\text{окр}} \text{ОКР} + K_{\text{СППЖ}} \text{СППЖ},$ $B = K_{\text{ОЗО}} \text{ОЗО} + K_{\text{ОКС}} \text{ОКС} + K_{\text{ПИНВ}} \text{ПИНВ}$ $\text{ИП} = A / (A + K_M B)$
7	$A = K_{\text{окр}} \text{ОКР} + K_{\text{СППЖ}} \text{СППЖ} + K_{\text{ИФС}} \text{ИФС}$ $B = K_{\text{ОЗИД}} \text{ОЗИД} + K_{\text{ОЗИВ}} \text{ОЗИВ} + K_{\text{ОКСД}} \text{ОКСД} + K_{\text{ОКСВ}} \text{ОКСВ} + K_{\text{ПИНВ}} \text{ПИНВ}$ $\text{ИП} = A / (A + K_M B)$
8	$A = K_{\text{окр}} \text{ОКР} + K_{\text{СППЖ}} \text{СППЖ} + K_{\text{ИФС}} \text{ИФС},$ $B = K_{\text{ОЗИ}} \text{ОЗИ} + K_{\text{ОКС}} \text{ОКС} + K_{\text{ПИНВ}} \text{ПИНВ},$ $\text{ИП} = A / (A + K_M B)$

Модели 5 и 6 используют только показатели здоровья, входящие в публикуемые ГОСКОМСТАТОМ РФ статистические данные. В модели 7 и 8 по-

мимо этих показателей входят и дополнительные показатели (индекс физического состояния и исчерпанная общая заболеваемость).

Для рассматриваемых моделей удобно сохранить методику определения весовых коэффициентов, разработанную для линейных моделей. В этом случае можно пользоваться рекомендуемыми значениями вкладов  $PZ_i$ , приведёнными в табл. 4.3. Значение  $K_m$  можно выбрать, например, исходя из того, чтобы при средних значениях ПЗ (табл. 4.2) как и для соответствующих линейных моделей ИП был бы равен 0,5. В этом случае из выражения (4.10) получаем, что для нелинейных моделей 5 и 6 (табл. 4.5) следует принять  $K_m = A/B = 0,6$ , а для нелинейных моделей 7 и 8 –  $K_m = 1$ . Однако определение  $K_m$  возможно и на основе других соображений. Так, поскольку более близкие результаты 6-я и 2-я модели дают при  $K_m = 0,75$ , то в дальнейшем это значение  $K_m$  и используется.

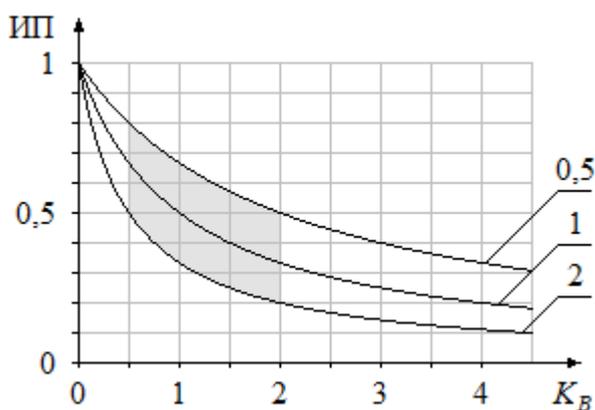
Т а б л и ц а 4.6. Нелинейные модели ИП с числовыми значениями коэффициентов

№	Реализуемые зависимости
5	$A = 10,742 \text{ ОКР} + 0,750 \text{ СППЖ},$ $B = 0,016 \text{ ОЗОД} + 0,072 \text{ ОЗОВ} + 0,920 \text{ ОСД} + 2,417 \text{ ОСВ} + 13,376 \text{ ПИНВ},$ $\text{ИП} = A/(A + 0,75B)$
6	$A = 10,742 \text{ ОКР} + 0,750 \text{ СППЖ},$ $B = 0,084 \text{ ОЗО} + 3,291 \text{ ОКС} + 13,376 \text{ ПИНВ},$ $\text{ИП} = A/(A + 0,75B)$
7	$A = 8,594 \text{ ОКР} + 0,600 \text{ СППЖ} + K_{\text{ИФС}} \text{ ИФС},$ $B = K_{\text{ОЗИД}} \text{ ОЗИД} + K_{\text{ОЗИВ}} \text{ ОЗИВ} + 0,828 \text{ ОКСД} + 1,830 \text{ ОКСВ} + 10,701 \text{ ПИНВ},$ $\text{ИП} = A/(A + 0,75B)$
8	$A = 8,594 \text{ ОКР} + 0,600 \text{ СППЖ} + K_{\text{ИФС}} \text{ ИФС},$ $B = K_{\text{ОЗИ}} \text{ ОЗИ} + 2,633 \text{ ОКС} + 10,701 \text{ ПИНВ},$ $\text{ИП} = A/(A + 0,75B)$

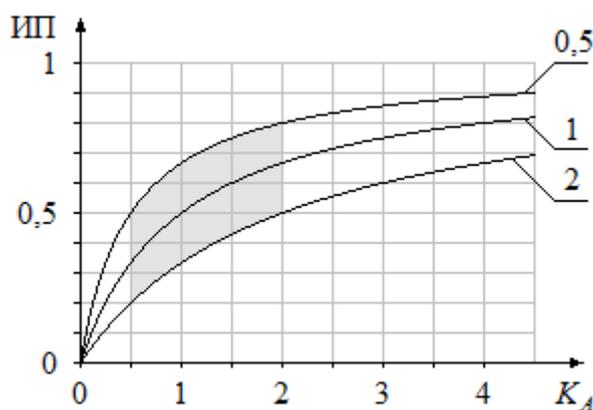
Продemonстрируем характер нелинейной зависимости ИП от соответствующих ПЗ на двух примерах для 6-й модели при изменении только составляющей  $A$  или только составляющей  $B$ . Пусть значения всех показателей этой модели равны произведениям их средних значений, указанных в табл. 4.2, на величину  $K_A$  для показателей группы  $A$  (ОКР и СППЖ) и на величину  $K_B$  для показателей группы  $B$  (ОЗО, ОКС и ПИНВ). В этом случае выражение для ИП принимает вид:

$$\text{ИП} = \frac{0,15K_A}{0,15K_A + 0,75 \cdot 0,25K_B} = \frac{K_A}{K_A + 1,25K_B}. \quad (4.11)$$

На рис. 4.8 приведены графики зависимости (4.11) от  $K_B$  при трёх значениях  $K_A$ , а на рис. 4.9 – от  $K_A$  при трёх значениях  $K_B$ . Затемнённая область на рисунках соответствует “рабочей области” значений ИП, за которую он обычно не выходит, изменяясь в  $[0,2, 0,8]$ .

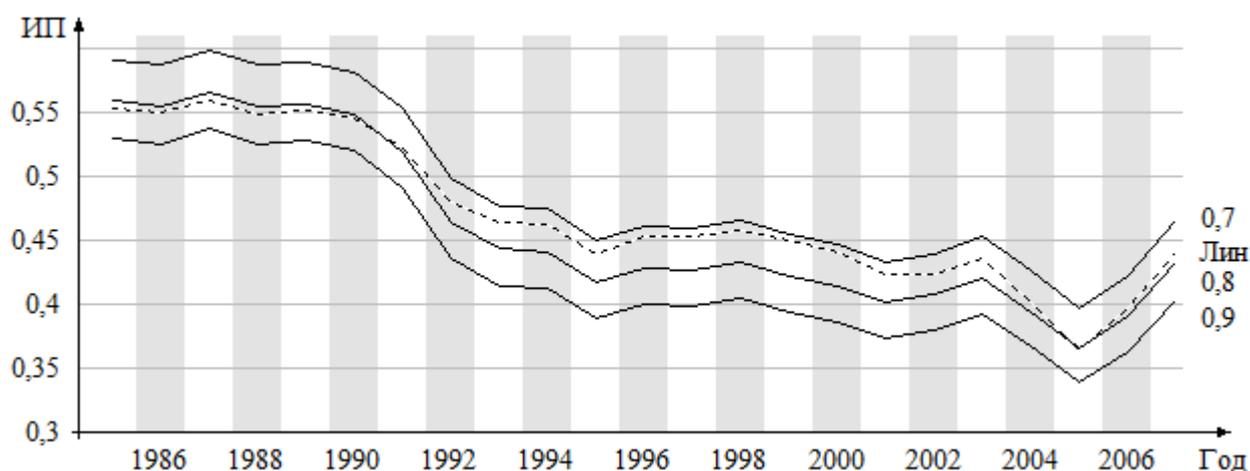


**Рис. 4.8.** Графики зависимости ИП от значения  $K_B$  при  $K_A \in \{0,5, 1, 2\}$



**Рис. 4.9.** Графики зависимости ИП от значения  $K_A$  при  $K_B \in \{0,5, 1, 2\}$

На рис. 4.10 для сравнения иллюстрируется динамика ИП, полученных с помощью линейной и нелинейной моделей на основе статистических данных по РФ. По графикам можно заключить, что даже “родственные” линейные и нелинейные модели с одинаковыми параметрами (2-я и 6-я) обычно дают несколько отличающиеся графики изменения ИП.



**Рис. 4.10.** Динамика интегральных показателей здоровья населения РФ, полученных на основе 2-й линейной модели (пунктирная линия) и на основе 6-й нелинейной модели при разных  $K_M$  (0,7, 0,8 и 0,9)

Достоинством линейных моделей ИП является то, что изменения ИП пропорциональны изменениям учитываемых моделями ПЗ. Вместе с тем при использовании этих моделей есть крайне малая, но всё же отличная от нуля вероятность выхода определяемых по линейной зависимости ИП за пределы отрезка  $[0, 1]$ . Для нелинейных моделей, имеющих на разных участках изменения ПЗ различные динамические свойства, выход ИП их промежутка  $[0, 1]$  исключён даже теоретически.

Учитывая одинаковую динамичность линейных моделей во всём промежутке  $[0, 1]$ , позволяющую сравнивать приращения ИП на любых участках указанного промежутка, авторы всё же отдают предпочтение этим моделям.

Отметим, что для сравнения здоровья населения различных регионов с помощью ИП здоровья необходимо, чтобы во всех регионах ИП определялись на основе одной и той же модели.

#### **4.6. Чувствительность интегральных показателей здоровья населения**

С чувствительностью ИП качества или функционирования любой многопараметрической системы связаны вопросы рационального выбора числа параметров, учитываемых моделью ИП, необходимой точности определения

значений этих параметров и реализации точного представления значений ИП. Поэтому рассмотрим понятие чувствительности и указанные вопросы применительно к ИП здоровья населения.

Итак, пусть модель ИП определена выражением  $ИП = f(ПЗ_1, ПЗ_2, \dots, ПЗ_n)$ , т.е. в общем случае зависит от  $n$  параметров. Если каждый из  $ПЗ_i$  изменится на величину  $\Delta ПЗ_i$ , то, полагая, что функция  $f$  является дифференцируемой, можно просто найти [70] приращение  $\Delta ИП_i$ , которое в этом случае получит ИП:

$$\Delta ИП_i \approx \sum_{i=1}^n \frac{\partial ИП}{\partial ПЗ_i} \Delta ПЗ_i. \quad (4.12)$$

Причём для случая линейной зависимости выходной величины от каждого параметра выражение (4.12) является точным [70].

В ряде прикладных задач (системный анализ, измерительная техника и др.) каждая частная производная от выходной величины системы по одной из нескольких независимых входных величин этой системы называется чувствительностью выходной величины по соответствующей входной. Распространим это понятие и на модели ИП здоровья.

Учитывая указанное, под чувствительностью ИП по  $i$ -му ПЗ, учитываемому моделью ИП, будем понимать частную производную  $\partial ИП / \partial ПЗ_i$ . Обозначая её  $Ч_i$ , получаем:

$$\Delta ИП_i \approx \sum_{i=1}^n \Delta ИП_i = \sum_{i=1}^n Ч_i \Delta ПЗ_i, \quad (4.13)$$

где величина  $\Delta ИП_i$  является реакцией ИП на изменение значения  $ПЗ_i$  при постоянстве значений всех остальных ПЗ, учитываемых моделью.

Приведём примеры выражений для  $Ч_i$ , получаемых путём дифференцирования функций ИП по соответствующим аргументам  $ПЗ_i$ . При этом ограничимся рассмотрением 2-й и 6-й моделей. Для 2-й модели (линейной) получим:  $Ч_{ОКР} = K_{ОКР} = 0,008398$ ,  $Ч_{СППЖ} = K_{СППЖ} = 0,000750$ ,  $Ч_{ОЗО} = K_{ОЗО} = -0,0000814$ ,  $Ч_{ОКС} = K_{ОКС} = -0,003682$ ,  $Ч_{ПИНВ} = K_{ПИНВ} = -0,013377$ . Знак ми-

нус в этих выражениях свидетельствует о том, что с увеличением рассматриваемого ПЗ значение ИП уменьшается. Для нелинейных моделей выражения для чувствительности ИП имеют более сложный вид. Для 6-й модели для чувствительности ИП по общему коэффициенту рождаемости получаем:

$$\chi_{\text{ОКР}} = \frac{0,75K_{\text{ОКР}}(K_{\text{ОЗО}}\text{ОЗО} + K_{\text{ОСМ}}\text{ОСМ} + K_{\text{ПИНВ}}\text{ПИНВ})}{[K_{\text{ОКР}}\text{ОКР} + K_{\text{СППЖ}}\text{СППЖ} + 0,75(K_{\text{ОЗО}}\text{ОЗО} + K_{\text{ОСМ}}\text{ОСМ} + K_{\text{ПИНВ}}\text{ПИНВ})]^2},$$

т.е. в отличие от линейных моделей  $\chi_{\text{ОКР}}$  зависит от значений нескольких ПЗ.

Используя выражение (4.13), можно рационально выбирать число разрядов для записи значений ПЗ и ИП, а также показать, что стремление значительно увеличить число ПЗ, являющихся аргументами функции ИП, оказывается неоправданным. Остановимся на этих вопросах.

В таблицах ГОСКОМСТАТа РФ до 2000 г. значения ПЗ приводились в расчёте на 1000 человек населения с числом разрядов после запятой от одного до четырёх, причём последний из этих разрядов обычно не заполнялся. В таблицах с 2000 г. указанные значения приводятся в расчёте как на 1000, так и на 10000 или на 100000 человек с одним разрядом после запятой. Последнее по точности представления эквивалентно указанию значений ПЗ на 1000 человек с 3 разрядами после запятой. Таким образом, можно считать, что в таблицах ГОСКОМСТАТа РФ значения ПЗ представляются с абсолютной погрешностью 0,0005. При этом в оптимальном случае изменение некоторого  $\text{ПЗ}_i$  на 0,001, т.е. на одну единицу младшего разряда, должно приводить к изменению расчётного значения интегрального показателя  $\Delta\text{ПЗ}_i$  не менее, чем на одну единицу младшего разряда ИП. В противном случае ИП не отреагирует на указанное изменение  $\text{ПЗ}_i$  т.е. ”не почувствует“ этого изменения  $\text{ПЗ}_i$ .

Рассмотрим пример. Пусть ИП определяется на основе пятипараметрической 2-й модели, в которой из весовых коэффициентов минимальное значение, равное 0,000081, имеет  $K_{\text{ОЗО}}$ . В этом случае при изменении ОЗО на минимально возможное, отличное от нуля значение 0,001, допускаемое таблицами ГОСКОМСТАТа РФ, согласно выражению (4.13) получаем:  $\Delta\text{ОЗО} =$

0,001. Следовательно, для реагирования нормированного ИП на указанное изменение ОЗО запись ИП должна иметь 9 разрядов после запятой.

Проанализируем теперь как влияет на работоспособность модели нормированного ИП увеличение числа ПЗ, по которым определяется значение ИП. Для этого рассмотрим модель ИП с  $n$  параметрами, полагая, что хотя бы по одному из ПЗ минимальное значение модуля чувствительности ИП равно  $0,0004/n$ , что соответствует рассмотренному примеру, в котором при  $n = 5$   $K_{OZO} = 0,000081$ . Будем предполагать, что минимальное, отличное от нуля изменение каждого ПЗ по-прежнему равно 0,001. Так как при увеличении  $n$  сумма вкладов показателей здоровья в ИП остаётся постоянной, то в этом случае при изменении ПЗ, имеющего указанный весовой коэффициент, на единицу младшего разряда значение ИП линейной модели должно будет измениться на  $0,0000004/n$ . Поэтому при неограниченном увеличении значения  $n$  чувствительность ИП по ПЗ с минимальными  $K_i$  будет стремиться к нулю, т.е. число разрядов после запятой в представлении ИП, необходимое для получения отличного от нуля  $\Delta$ ИП при изменении ПЗ с указанным  $K_i$  на 0,001, будет стремиться к бесконечности. Причём с увеличением числа указанных ПЗ и соответствующим уменьшением практически всех  $K_i$  уменьшение чувствительности ИП произойдёт по всем ПЗ.

К аналогичному выводу можно придти и для нелинейных моделей. Следовательно, существует определённый компромисс между сложностью модели ИП и чувствительностью ИП. Поэтому в предложенных моделях ИП используется от 5 до 10 ПЗ, являющихся наиболее существенными для оценки здоровья населения.

С усложнением моделей может потребоваться дальнейшее увеличение числа указанных разрядов ввиду уменьшения чувствительности ИП по соответствующим ПЗ, что неудобно для практического использования ИП. По-видимому, нет смысла учитывать в рассматриваемых моделях такие ПЗ, весовые коэффициенты которых меньше 0,0001 или вклады которых в ИП

меньше 0,001. Это определяет и максимальное число параметров моделей интегральных показателей общественного здоровья населения.

#### 4.7. Сравнение интегральных оценок здоровья населения регионов Российской Федерации

Разработанные модели ИП позволяют сравнивать здоровье населения различных регионов. В табл. 4.7 и 4.8 приводятся значения ИП здоровья населения РФ и её федеральных округов (ФО), полученные с помощью 2-й (линейной) и 6-й (нелинейной) моделей. Названия федеральных округов в первом столбце этих таблиц следуют в порядке уменьшения среднего значения их ИП здоровья за 8 рассмотренных лет.

Т а б л и ц а 4.7. Значения интегрального показателя общественного здоровья населения РФ и федеральных округов (на основе 2-й модели)

РФ и ФО	Годы	За 8 лет	2000	2001	2002	2003	2004	2005	2006	2007
Россия		0,4156	0,440	0,423	0,424	0,436	0,402	0,364	0,396	0,440
Уральский		0,4617	0,468	0,468	0,470	0,477	0,456	0,425	0,443	0,486
Южный		0,4606	0,482	0,456	0,455	0,473	0,443	0,423	0,454	0,500
Дальневосточный		0,4605	0,474	0,469	0,470	0,476	0,450	0,417	0,448	0,481
Сибирский		0,4373	0,455	0,448	0,446	0,455	0,426	0,392	0,415	0,461
Приволжский		0,4072	0,423	0,420	0,424	0,431	0,398	0,357	0,379	0,425
Центральный		0,3753	0,393	0,380	0,381	0,398	0,360	0,319	0,360	0,411
Северо-Западный		0,3524	0,384	0,372	0,368	0,379	0,333	0,314	0,304	0,366

Т а б л и ц а 4.8. Значения интегрального показателя общественного здоровья населения РФ и федеральных округов (на основе 6-й модели)

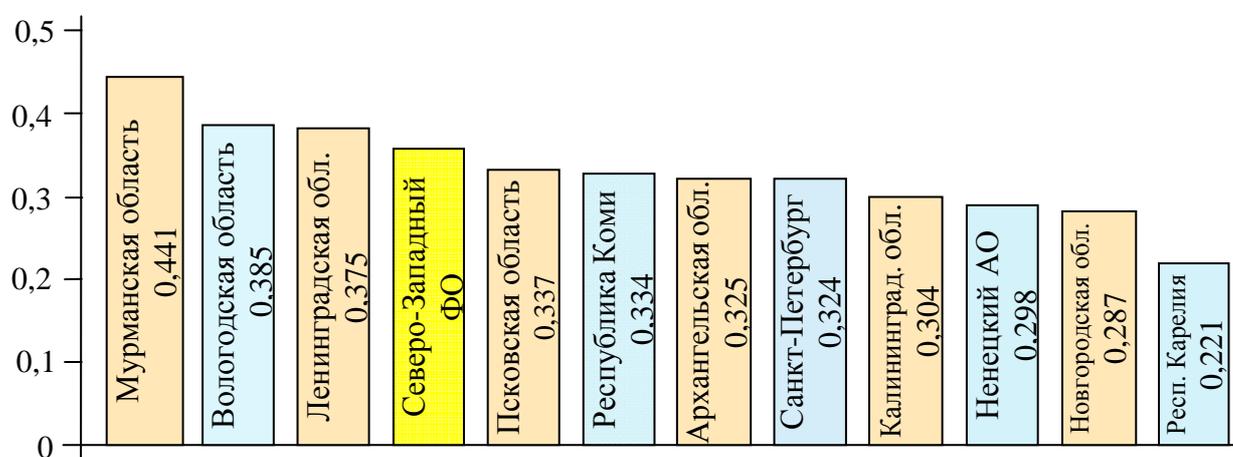
РФ и ФО	Годы	За 8 лет	2000	2001	2002	2003	2004	2005	2006	2007
Россия		0,4187	0,429	0,417	0,423	0,436	0,410	0,381	0,406	0,448
Южный		0,4687	0,485	0,457	0,460	0,481	0,454	0,435	0,463	0,514
Уральский		0,4659	0,464	0,465	0,472	0,481	0,462	0,433	0,451	0,497
Дальневосточный		0,4651	0,468	0,467	0,471	0,482	0,459	0,427	0,456	0,491

Сибирский	0,4441	0,450	0,446	0,449	0,461	0,436	0,408	0,428	0,474
Приволжский	0,4099	0,412	0,413	0,421	0,431	0,405	0,372	0,391	0,435
Центральный	0,3761	0,378	0,371	0,376	0,393	0,367	0,339	0,367	0,417
Северо-Западный	0,3619	0,370	0,365	0,369	0,383	0,354	0,342	0,335	0,378

Из приведенных таблиц следует, что применение нелинейной модели ИП может приводить к некоторому изменению значений ИП по сравнению с результатами, получаемыми с помощью линейной модели ИП. Однако при этом тенденции изменения ИП сохраняются.

Недостатком нелинейных моделей является то, что они не позволяют объективно сравнивать значения прироста (приращения), увеличения прироста и темпа роста ИП показателей здоровья населения различных регионов. Линейные же модели дают приращения ИП, пропорциональные приращениям ПЗ при любых значениях ПЗ. Поэтому в дальнейшем используются только линейные модели ИП общественного здоровья населения.

С помощью предложенных моделей можно получить и интегральные оценки здоровья населения административных образований федеральных округов. Пример распределения таких единиц по величине ИП приведён на рис 4.11. При этом следует иметь в виду, что в общем случае с уменьшением численности населения административных единиц увеличивается разброс значений их ИП здоровья их населения.



**Рис. 4.11.** Диаграмма ИП общественного здоровья населения административных образований Северо-Западного ФО в 2006 г., полученная с помощью 2-й модели

Таким образом, разработанные модели интегрального показателя общественного здоровья населения на основе показателей здоровья, публикуемых ГОСКОМСТАТОм РФ, могут использоваться для оценки и сравнения состояния здоровья населения различных регионов. Для более адекватного отражения получаемыми интегральными показателями состояния здоровья населения целесообразно учитывать в моделях интегрального показателя и приводить в публикуемой государственной статистике показателей здоровья показатели физического развития, инвалидности и исчерпанной заболеваемости населения. Так, расчёт ИП здоровья населения Новгородской области за 2005-й год согласно 4-й модели, т.е. с учётом индекса физического состояния (0,541) и исчерпанной заболеваемости (3255,5), даёт значение ИП, равное 0,283 вместо значения 0,202, получаемого на основе 2-й модели (в 2005 году ИП здоровья населения Новгородской области был минимальным).

## ГЛАВА 5. КОРРЕЛЯЦИЯ И РЕГРЕССИЯ. МОДЕЛИ ЗАВИСИМОСТЕЙ

### 5.1. Типы зависимостей

Любое множество объектов без указания зависимостей между ними оказывается хаотичным и не информативным. Поэтому изучение связи между различного рода показателями типа «фактор(ы)-отклик», «доза-эффект» является одной из важнейших задач статистики. Установление типа зависимости и непосредственного количественного соотношения для значений показателей позволяет:

- выявить объективно существующие причинно-следственные связи,
- исследовать «механизм» взаимодействия рассматриваемых процессов,
- использовать возможности влияния на отклик через факторы, его формирующие;
- составить математически обоснованный прогноз на будущее.

Многие прикладные задачи, в частности задачи медицинской статистики, требуют установления вида связи (зависимости) между факторами риска и заболеваниями, и между самими факторами риска, между значениями одного и того же показателя, но характеризующего разные группы объектов. Во всех случаях факторы выступают как случайные величины.

Как известно, случайные величины  $X$  и  $Y$  могут быть либо независимыми, либо зависимыми. Зависимость случайных величин также подразделяется на функциональную и статистическую. В математике функциональной зависимостью переменной  $Y$  от переменной  $X$  называют зависимость вида  $Y = f(X)$ , где каждому допустимому значению  $X$  ставится в соответствие по определенному правилу единственно возможное значение  $Y$ .

В реальных задачах исследования факторов  $X$  и  $Y$  – это случайные величины и между ними может существовать зависимость иного рода, называемая стохастической (статистической) зависимостью. При этом каждому

значению  $X$  может соответствовать не одно значение  $Y$ , как при функциональной зависимости, а целое множество значений. Среди этого множества значений  $Y$  можно найти среднее  $M(Y|X = x)$ , которое для каждого значения  $X$  свое, т.е. условное математическое ожидание  $M(Y|X = x)$  – это функция от  $x$ . Множество возможных значений  $(x, y)$  на графике образуют линию зависимости  $y$  от  $x$ :

$$y = M(Y | X = x), \quad (5.1)$$

вид которой может быть самым разнообразным (прямая, парабола, экспонента и т.д.) и определяется случайными величинами  $X$  и  $Y$ . Соответствующий график называют линией регрессии  $Y$  на  $X$ .

Зависимость случайных величин называют стохастической (статистической), если изменение одной из них приводит к изменению закона распределения другой. Изменения закона распределения могут проявляться по-разному. В частности, если изменение одной из случайных величин влечет изменение среднего другой случайной величины, то стохастическую зависимость называют корреляционной. Сами случайные величины, связанные корреляционной зависимостью, оказываются коррелированными.

Корреляционной будет зависимость заболеваемости от воздействия внешних факторов, например запыленности, уровня радиации, солнечной активности, уровня холестерина, эстрогенов в крови и т.д. Имеется корреляция между дозой ионизирующего излучения и числом мутаций, между пигментом волос человека и цветом глаз, между уровнем жизни населения и показателями уровня смертности. Именно корреляционные зависимости наиболее часто встречаются в природе в силу взаимовлияния и тесного переплетения огромного множества самых различных факторов, определяющих значения изучаемых показателей.

Напомним, что независимые случайные величины – это частный случай некоррелированных. Некоррелированные случайные величины могут оказаться как независимыми, так и зависимыми.

При корреляционной зависимости  $Y$  и  $X$  можно наблюдать тенденцию роста: с увеличением значений  $X$  среднее значение  $Y$  возрастает или с увеличением значений  $X$  среднее значение  $Y$  уменьшается. В этих случаях говорят соответственно о положительной и отрицательной корреляции.

## 5.2. Выборочный коэффициент корреляции

Как известно степень зависимости случайных величин  $X$  и  $Y$  (двух признаков) характеризуется значением коэффициента корреляции:

$$r = r_{xy} = \frac{K(X, Y)}{\sqrt{D(X)}\sqrt{D(Y)}},$$

где  $K(X, Y)$ - корреляционный момент (ковариация) случайных величин  $X$  и  $Y$ ,  $D(X)$  и  $D(Y)$  – соответствующие дисперсии.

При этом коэффициент корреляции, как и всякая другая теоретическая характеристика, вычисляется, исходя из всех возможных значений  $X$  и  $Y$ . На практике же мы не имеем возможности охватить наблюдениями все означенное множество, а используем лишь ограниченное число наблюдений: двумерную выборку значений  $(x, y)$ . Полученные числа можно занести в таблицу (табл. 5.1).

Т а б л и ц а 5.1. Запись двумерной выборки

<b><math>X</math></b>	$x_1$	$x_2$	...	$x_n$
<b><math>y</math></b>	$y_1$	$y_2$	...	$y_n$

По данным наблюдений можно вычислить значение коэффициента корреляции так же, как и в случае системы дискретных случайных величин, с той лишь разницей, что вместо известных вероятностей для каждой пары

возможных значений будем использовать соответствующий аналог: относительную частоту  $1/n$ .

Формула для вычисления выборочного коэффициента корреляции  $r_B$  случайных величин  $X$  и  $Y$  выглядит так:

$$r_B = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \cdot \sum_{j=1}^n (y_j - \bar{y})^2}}. \quad (5.2)$$

Если число наблюдений достаточно велико и особенно, если наблюдения объединяются поинтервально, т.е. все значения, попавшие в интервал, округляются до значения середины интервала, то каждая из наблюдаемых пар значений может встретиться неоднократно. В этом случае обычно данные заносят в таблицу с учетом частот встречаемости. Такую таблицу, сгруппированными данными называют корреляционной.

**Пример 5.1.** Получена корреляционная таблица (табл. 5.2.), составленная по выборке 150 студентов возраста 20 - 22 лет. Конкретные данные: случайная величина  $X$  – стаж курильщика (количество лет), случайная величина  $Y$  – жизненная емкость легких (ЖЕЛ) в мл. Значения  $X$  – середины соответствующих интервалов (0 - 2), (2 - 4), (4 - 6), (6 - 8), (8 - 10). Для некурящих полагаем  $X = 0$ . Значения  $Y$  также рассматриваются поинтервально: (3000 - 3500), (3500 - 4000), (4000 - 4500), (4500 - 5000), (5000 - 5500).

Таблица 5.2. Корреляционная таблица соотношения стажа курильщика и показателя ЖЕЛ

$X \backslash Y$	3250	3750	4250	4750	5250	$n_{x_i}$
0	1	5	25	17	7	55
1	2	-	8	8	3	21
3	2	4	9	5	4	24
5	2	6	6	4	1	19
7	3	12	2	1	-	18
9	4	7	-	2	-	13
$n_{y_j}$	14	34	50	37	15	$n = 150$

В первом столбце и первой строке таблицы указаны наблюдаемые значения соответственно случайных величин  $X$  и  $Y$ . На пересечении фиксированных строки

и столбца записано число – частота встречаемости этой пары в данной выборке. Например, пара значений  $X = 3, Y = 4250$  (запишем в виде  $(3; 4250)$ ) встречается 9 раз, а комбинация  $(9; 5250)$  не встречается ни разу. В общем случае частоту пары  $(X; Y)$  в выборке обозначим  $n_{xy}$ , или более конкретно: пара  $(x_i; y_j)$  имеет частоту  $n_{ij}$ . В последнем столбце представлены суммарные значения частот встречаемости в выборке пар с каждым из наблюдаемых значений  $X$ , например,  $X = 0$  встречается с различными значениями  $Y$  суммарно 55 раз ( $1 + 5 + 25 + 17 + 7 = 55$ ),  $X = 1$  появляется 21 раз и т.д. Аналогично, в последней строке корреляционной таблицы выписаны суммарные частоты встречаемости в наблюдениях каждого из значений  $Y$ . Сумма всех частот, записанных в таблице для каждой пары значений  $X, Y$ , равна  $n$  и, естественно,

$$\sum n_{ij} = \sum n_{x_i} = \sum n_{y_j} = n.$$

В рассматриваемой таблице  $n = 150$ .

Формула для выборочного коэффициента корреляции  $r_B$ , вычисляемого по корреляционной таблице, т. е. с учетом частот встречаемости в наблюдениях каждой пары  $(x, y)$ , принимает вид

$$r_B = \frac{\sum n_{xy} (x - \bar{x})(y - \bar{y})}{\sqrt{\sum n_x (x - \bar{x})^2 \cdot \sum n_y (y - \bar{y})^2}}, \quad (5.3)$$

где  $\bar{x}$  и  $\bar{y}$  соответствующие выборочные средние

$$\bar{x} = \frac{1}{n} \sum n_x x, \quad \bar{y} = \frac{1}{n} \sum n_y y,$$

а суммирование распространяется в знаменателе на все возможные  $x$  и  $y$  соответственно и в числителе – на все возможные пары  $(x, y)$ .

Если в числителе выражения (5.3) под знаком суммы выполнить умножение, то данная формула преобразуется к виду

$$r_B = \frac{\sum n_{xy} xy - n \bar{x} \bar{y}}{\sqrt{\sum n_x (x - \bar{x})^2 \sum n_y (y - \bar{y})^2}}. \quad (5.4)$$

Заметим, что в знаменателе выражение под знаком корня равно  $n \bar{s}_x^2 \cdot n \bar{s}_y^2$ , где  $\bar{s}_x^2, \bar{s}_y^2$  – выборочные дисперсии.

Проведя вычисления непосредственно по числовым данным таблицы 12, получаем

$$\begin{aligned} \bar{x} &= 2,87333; & \bar{y} &= 4266,66667; \\ r_B &= \frac{\sum n_x xy - n \bar{x} \bar{y}}{\sqrt{\sum n_x (x - \bar{x})^2 \cdot \sum n_y (y - \bar{y})^2}} = \end{aligned}$$

$$= \frac{1716250 - 150 \cdot 2,87333 \cdot 4266,66667}{\sqrt{1408,593 \cdot 46708333}} = -0,47829.$$

Как известно, коэффициент корреляции случайных величин (генеральных совокупностей)  $X$  и  $Y$ , изменяясь по модулю в пределах от 0 до 1, характеризует тесноту линейной связи: от полной независимости (или, по крайней мере, некоррелированности) случайных величин при  $r = 0$  до линейного соотношения,  $Y = aX + b$ , при  $|r| = 1$ . Причем при  $r > 0$  возрастание одной переменной влечет и рост другой (положительная корреляция), а при  $r < 0$  возрастание одной переменной влечет убывание другой (отрицательная корреляция). Отметим, что в данном случае речь идет не об однозначном возрастании или убывании переменной по всем значениям, а об общей тенденции. В приведенном примере отрицательное значение коэффициента корреляции свидетельствует о тенденции уменьшения значений показателя ЖЕЛ при увеличении стажа курильщика. Вместе с тем следует отметить невысокое значение коэффициента корреляции, что вызвано качеством данных, а именно: показатель ЖЕЛ, как известно, существенно зависит от роста и массы тела индивидуума (и других характеристик, не связанных с курением), а эти факторы при расчете не учтены. Объективно вне зависимости от фактора курения индивидуум невысокого роста с малой массой тела имеет меньшую емкость легких, чем индивидуум высокого роста с большой массой тела. Неоднородность рассматриваемой совокупности не позволила вычленить именно изучаемую зависимость, а полученное значение коэффициента корреляции оказалось сформированным под воздействием различных не учитываемых при расчетах разнонаправленных факторов. Исправить ситуацию можно двумя способами:

- переходом от абсолютных к относительным показателям (в рассмотренном примере вместо показателя– отклика ЖЕЛ следует использовать показатели типа

$$\frac{\text{ЖЕЛ}}{\text{Рост}}, \frac{\text{ЖЕЛ}}{\text{Масса тела}}, \frac{\text{ЖЕЛ}}{\text{Массо–ростовой коэффициент}}, \frac{\text{ЖЕЛ}}{\text{Индекс массы тела}});$$

- формируя совокупность, однородную по другим факторам (например, включать в рассматриваемую совокупность лишь индивидуумов, близких по значениям роста и массы тела); разумеется, выборочная совокупность должна быть однородна и по половозрастному признаку.

В частности, в приведенном примере, формируя выборку на основе имеющихся наблюдений при росте 165 - 175 см и нормальной массе тела (согласно градации массо-ростового коэффициента), получим корреляционную таблицу, (табл. 5.3). Всего в новой выборке оказались показатели 77 студентов из первоначально отобранных 150.

Т а б л и ц а 5.3. Корреляционная таблица соотношения стажа курильщика и показателя ЖЕЛ, составленная по новым данным.

$X \backslash Y$	3250	3750	4250	4750	5250	$n_{x_i}$
0	-	1	8	12	1	22
1	-	-	8	6	-	14
3	-	2	9	3	-	14
5	1	3	4	1	-	9
7	3	5	2	-	-	10
9	4	4	-	-	-	8
$n_{y_j}$	8	15	31	22	1	$n = 77$

Для рассматриваемых значений (т.е. по более однородной выборке) коэффициент корреляции признаков вновь можно найти по формуле (5.4) и он оказывается существенно выше:  $r_B = -0,7535$ . Выборочный коэффициент корреляции  $r_B$  – оценка коэффициента корреляции  $r$ , рассчитанного по всей генеральной совокупности, т. е.  $r_B \cong r$ . Следовательно, рассчитав  $r_B$ , мы также можем судить о силе линейной связи. В случае если выборка имеет достаточно большой объем  $n$ , порядка сотен, то можно воспользоваться  $r_B$  как точечной оценкой коэффициента корреляции  $r$ .

### 5.3. Проверка независимости признаков

Рассмотрим два признака. Первый из них характеризуется рядом значений  $x_1, x_2, \dots, x_n$ , а второй – соответствующими значениями  $y_1, y_2, \dots, y_n$ .

Исходя из имеющихся данных, требуется установить зависимы ли эти признаки. Полагаем, что наблюдаемые значения  $x_i$  представляют генеральную совокупность  $\mathbf{X}$ , а значения  $y_i$  – генеральную совокупность  $\mathbf{Y}$ , и тогда задача сводится к проверке независимости случайных величин  $X$  и  $Y$ . В этом нам поможет коэффициент корреляции. Предположим, что случайные величины  $X$  и  $Y$  имеют нормальное распределение. Как известно (гл.2), условие некоррелированности нормально распределенных  $X$  и  $Y$  равносильно их независимости. Т.е. условие  $r_{xy} = 0$  гарантирует независимость случайных величин (а значит, и независимость признаков), а условие  $r_{xy} \neq 0$  означает зависимость  $X$  и  $Y$ , причем величина коэффициента  $r_{xy}$  свидетельствует о степени имеющейся зависимости.

Заметим, что значение выборочного коэффициента корреляции является лишь оценкой «истинного» теоретического значения  $r_{xy}$  и отличается от него в силу различных случайных причин (несмотря на соблюдение всех возможных условий репрезентативности отбора, просто, не повезло с выборкой). Даже при очевидной независимости признаков, скорее всего, окажется  $r_B \neq 0$ . Поэтому следует установить, отличие  $r_B$  от нуля вызвано случайными причинами, связанными с выборкой (незначимо), или же оно принципиально, т.е. объясняется именно зависимостью признаков (значимо).

Итак, в качестве нулевой гипотезы полагаем  $H_0: r_{xy} = 0$ , тогда конкурирующая гипотеза  $H_1: r_{xy} \neq 0$ .

Для проверки гипотезы  $H_0$  требуется, исходя из выборочных данных, подобрать статистику-критерий, которая бы через выборочный коэффициент корреляции  $R_B$  характеризовала  $r_{xy}$  и при этом, распределение которой было бы хорошо известно. Таким критерием является статистика

$$T = \frac{R_B \sqrt{n-2}}{\sqrt{1-R_B^2}}, \quad (5.5)$$

имеющая распределение Стьюдента с  $(n - 2)$  степенями свободы при справедливости нулевой гипотезы.

Следовательно, для проверки  $H_0$  используется критерий Стьюдента, т.е. по выбранному уровню значимости  $\alpha$  и числу степеней свободы  $(n - 2)$  необходимо по соответствующей таблице критических точек распределения Стьюдента (табл. П.3) найти  $t_{кр}$ , затем определить  $T_{набл}$  по формуле

$$T_{набл} = \frac{r_B \sqrt{n-2}}{\sqrt{1-r_B^2}} \quad (5.6)$$

и сравнить полученные значения  $t_{кр}$  и  $T_{набл}$ . Если  $|T_{набл}| < t_{кр}$ , то наблюдаемое отклонение  $r_B$  от нуля незначимо, и нулевая гипотеза принимается: случайные величины независимы. Если же  $|T_{набл}| > t_{кр}$ , то принимается гипотеза  $H_1$ : случайные величины  $X$  и  $Y$  зависимы.

Замечание. Вместо использования таблицы критических значений (квантилей) распределения Стьюдента можно воспользоваться специальной таблицей «наибольших случайных значений коэффициента корреляции» (табл. П.4). В этой таблице при заданном уровне значимости  $\alpha$  ( $\alpha = 0,001; 0,01; 0,027; 0,05$ ) и известном числе степеней свободы  $m = n - 2$  представлены максимальные значения, которые может иметь величина  $|r_B|$  при справедливости нулевой гипотезы.

**Пример 5.2.** По данным примера 5.1, при 77 наблюдениях, требуется определить, зависит ли в генеральных совокупностях значение показателя ЖЕЛ ( $Y$ ) от стажа курильщика ( $X$ ). Распределение случайных величин  $X$  и  $Y$  предполагается нормальным.

Используем критерий Стьюдента для проверки гипотезы  $H_0 : r_{xy} = 0$ .

В нашем примере число наблюдений  $n = 77$ , выборочный коэффициент корреляции  $r_B = -0,7535$ . Найдем  $T_{набл}$ :

$$T_{набл} = \frac{-0,7535 \sqrt{77-2}}{\sqrt{1-(-0,7535)^2}} = -9,9255.$$

Выберем уровень значимости  $\alpha = 0,01$  и по таблице критических значений распределения Стьюдента находим  $t_{кр} = t_{кр}(0,01; 75) = 2,643$ .

Так как  $|T_{набл}| > t_{кр}$ , то при выбранном уровне значимости нулевую гипотезу отвергаем; следовательно, случайные величины  $X$  и  $Y$  зависимы. К данному выводу можно было прийти, используя таблицу, наибольших случайных значений коэффициента корреляции (табл. П.4). В этой таблице при уровне значимости  $\alpha = 0,01$  и числе степеней свободы  $m = 77 - 2 = 75$  наибольшее возможное значение  $|r_B|$ , вызванное случайными причинами (а на самом деле  $r_{xy} = 0$ ), может быть 0,29. Наше наблюдаемое значение  $|r_B| = 0,7535$ , следовательно, наблюдаемое отклонение от нуля коэффициента корреляции вызвано не только случайными причинами, а и зависимостью  $X$  и  $Y$ . Таким образом,  $Y$  зависит от  $X$ , т.е. значение показателя ЖЕЛ зависит от стажа курильщика.

#### 5.4. Проверка гипотезы о силе линейной связи двух признаков

Допустим, что гипотеза о независимости двух исследуемых признаков отвергнута, как в приведенном выше примере. Следовательно, статистически установлена зависимость признаков. Идем дальше: установим, сколь сильна эта зависимость. Для решения вопроса о линейной зависимости нужно знать величину коэффициента корреляции. Мы же располагаем лишь его полномочным представителем, выборочным коэффициентом корреляции, который в силу разных случайных обстоятельств несколько отличается от своего теоретического оригинала  $r_{xy}$ . Вновь строим гипотезы: в нулевой гипотезе предполагаем равенство  $r_{xy}$  конкретному числу, отличному от нуля

$$H_0 : r_{xy} = a,$$

тогда конкурирующая гипотеза

$$H_1 : r_{xy} \neq a.$$

Далее для проверки нулевой гипотезы необходимо подобрать подходящий критерий. Фишером установлено, что статистика

$$W = \frac{1}{2} \ln \frac{1+R_B}{1-R_B}, \quad (5.7)$$

построенная по выборкам из  $X$  и  $Y$  достаточно большого объема  $n$  ( $n > 50$ ) имеет приближенно нормальное распределение. В случае справедливости

нулевой гипотезы,  $r_{xy} = a$ , параметры этого распределения можно найти по формулам:

$$m_W = M(W) = \frac{1}{2} \ln \frac{1+a}{1-a} + \frac{a}{2(n-1)}, \quad (5.8)$$

$$\sigma_W^2 = D(W) = \frac{1}{n-3}. \quad (5.9)$$

Для использования в качестве критерия более удобна нормированная величина, т. е. с нулевым средним и единичной дисперсией. Полагаем

$$Z = \frac{W - m_W}{\sigma_W}, \text{ причем } Z \sim N(0; 1) \quad (5.10)$$

Этот критерий и используется для проверки нулевой гипотезы.

Далее схема простая: вычисляем  $Z_{набл}$ , при выбранном уровне значимости  $\alpha$  по таблице функции Лапласа (табл. П.1) находим  $z_{кр}$  и сравниваем эти величины. В зависимости от полученных значений  $Z_{набл}$  и  $z_{кр}$  нулевая гипотеза либо принимается, либо отвергается.

**Пример 5.3.** Нами установлено, что значение показателя ЖЕЛ зависит от стажа курильщика. Найдем степень этой корреляционной зависимости.

Вычислено значение  $r_B = -0,7535$ . Проверим нулевую гипотезу  $H_0$ :  $r_{xy} = -0,8$  при уровне значимости  $\alpha = 0,05$  (число  $a = -0,8$  выбрано, исходя из численного значения  $r_B = -0,7535$ ). Конкурирующая гипотеза  $H_1$ :  $r_{xy} \neq -0,8$ .

Для использования критерия  $Z$  вначале вычислим  $m_W$ ,  $\sigma_W$ :

$$m_W = \frac{1}{2} \ln \frac{1+(-0,8)}{1-(-0,8)} + \frac{-0,8}{2(77-1)} \approx -1,1039,$$

$$\sigma_W = \sqrt{\frac{1}{77-3}} \approx 0,1162.$$

Тогда

$$Z = \frac{W - (-1,1039)}{0,1162}.$$

Находим наблюдаемое значение критерия:

$$Z_{набл} = \frac{\frac{1}{2} \ln \frac{1+(-0,7535)}{1-(-0,7535)} - (-1,1039)}{0,1162} \approx 1,0576.$$

Теперь при уровне значимости  $\alpha = 0,05$  критическое значение критерия  $z_{кр}$  находится из таблицы значений функции Лапласа:  $z_{кр} = 1,96$ . Сравнивая  $Z_{набл}$  с найденным критическим значением, видим, что наблюдаемое значение критерия находится в области принятия гипотезы:  $|Z_{набл}| < z_{кр}$ .

Следовательно, нулевая гипотеза принимается: при уровне значимости  $\alpha = 0,05$  (т.е. с вероятностью ошибиться не более чем в 5% реальных случаев) коэффициент корреляции равен 0,8. Эта зависимость достаточно сильная (при  $r = 1$  зависимость уже функциональная  $Y = aX + b$ , где  $a, b$  – некоторые числа).

### 5.5. Выборочная регрессия

Для зависимых признаков можно построить регрессионную модель, численно отражающую эту корреляционную зависимость, т.е. найти функциональную зависимость, *приблизенно* связывающую значения исследуемых показателей. Поскольку приближения всегда берутся с той или иной степенью точности и данные связаны не однозначно (одному значению  $X$  могут соответствовать несколько значений  $Y$ ), то наряду с самим соотношением (моделью) вводятся характеристики, численно отражающие качество полученной модели, степень ее соответствия реальным данным. При этом вычисления обычно проводятся с помощью стандартного программного обеспечения. В исходных данных соответствующие выборочные значения должны быть указаны попарно.

Например, в рассмотренном выше примере (см. табл.5.3.) линейная модель зависимости показателя ЖЕЛ (отклик  $Y$ ) от стажа курильщика (фактор  $X$ ) имеет вид

$$y = 4577326 - 118124x.$$

Качественные характеристики этой модели: коэффициент детерминации  $R^2$  равен 0,56777, стандартная ошибка - 322,89.

Таким образом, в отличие от коэффициента корреляции, характеризующего зависимость показателей одним числом, регрессионная зависимость

является функциональной, указывающей связь между всеми теоретически возможными значениями показателей.

Обсудим данную теорию подробнее.

Пусть для системы случайных величин  $X$  и  $Y$  наблюдаемые пары значений  $(x, y)$  оформлены в виде корреляционной таблицы 5.4. В этой таблице в отличие от табл. 5.3 данные представлены в общем виде.

Т а б л и ц а 5.4. Корреляционная таблица

$X \backslash Y$	$y_1$	$y_2$	$\dots$	$y_j$	$\dots$	$y_l$	$n_{x_i}$
$x_1$	$m_{11}$	$m_{12}$	$\dots$	$m_{1j}$	$\dots$	$m_{1l}$	$n_{x_1}$
$x_2$	$m_{21}$	$m_{22}$	$\dots$	$m_{2j}$	$\dots$	$m_{2l}$	$n_{x_2}$
$\dots$	$\dots$	$\dots$	$\dots$	$\dots$	$\dots$	$\dots$	$\dots$
$x_i$	$m_{i1}$	$m_{i2}$	$\dots$	$m_{ij}$	$\dots$	$m_{il}$	$n_{x_i}$
$\dots$	$\dots$	$\dots$	$\dots$	$\dots$	$\dots$	$\dots$	$\dots$
$x_k$	$m_{k1}$	$m_{k2}$	$\dots$	$m_{kj}$	$\dots$	$m_{kl}$	$n_{x_k}$
$n_{y_j}$	$n_{y_1}$	$n_{y_2}$	$\dots$	$n_{y_j}$	$\dots$	$n_{y_l}$	$n$

Случайная величина  $X$  принимает  $k$  различных значений, а случайная величина  $Y$  –  $l$  различных значений. На пересечениях строк и столбцов имеющих значения указана соответствующая частота, например, пара  $(x_i, y_j)$  встречается в нашей выборке  $m_{ij}$  раз (какие-то  $m_{ij}$  могут равняться нулю). В последнем столбце и строке выписаны суммы соответствующих частот для значений  $X$  и  $Y$ . Например,

$$n_{x_i} = m_{i1} + m_{i2} + \dots + m_{il},$$

$$n_{y_2} = m_{12} + m_{22} + \dots + m_{k2}.$$

Сумма всех возможных частот  $m_{ij}$  равна  $n$  и, конечно же, равна сумме этих частот, рассматриваемых отдельно по строкам и столбцам таблицы:

$$n = \sum_{i=1}^k \sum_{j=1}^l m_{ij} = \sum_{i=1}^k n_{x_i} = \sum_{j=1}^l n_{y_j}. \quad (5.11)$$

Каждому числу  $x_i$  соответствует целый набор значений  $y_1, y_2, \dots, y_l$  с конкретными частотами  $m_{i1}, m_{i2}, \dots, m_{il}$  (вновь отметим, что какие-то из час-

тот могут быть равны нулю). Вычислим среднее этих значений, которое обозначим  $\bar{y}_{x_i}$  :

$$\bar{y}_{x_i} = \frac{1}{n_{x_i}}(y_1 m_{i1} + y_2 m_{i2} + \dots + y_l m_{il}). \quad (5.12)$$

Среднее значение  $\bar{y}_{x_i}$  (условное среднее значение  $y$  при условии, что  $X = x_i$ ) можно вычислить для каждого  $x_i$ . Занесем полученные данные в табл. 5.5.

Таблица 5.5. Условные средние значения

$x$	$x_1$	$x_2$	...	$x_k$
$\bar{y}_x$	$\bar{y}_{x_1}$	$\bar{y}_{x_2}$	...	$\bar{y}_{x_k}$
$n_x$	$n_{x_1}$	$n_{x_2}$	...	$n_{x_k}$

Из табл. 5.5 легко прослеживается зависимость (соответствие) средних значений  $\bar{y}_x$  от значений  $X$ , т. е.

$$\bar{y}_x = \varphi^*(x). \quad (5.13)$$

Напомним, что в корреляционной таблице 5.4 представлены выборки случайных величин  $X$  и  $Y$ , значения которых рассматриваются попарно. Если бы мы имели возможность рассмотреть все возможные значения  $X$  и  $Y$ , то наше среднее  $\bar{y}_x$  оказалось бы не чем иным, как условным математическим ожиданием  $M(Y|X = x)$ , которое при разных  $x$  представляет собой функцию,  $\varphi(x)$ , называемую регрессией. Равенство вида

$$M(Y|X = x) = \varphi(x) \quad (5.14)$$

- это уравнение регрессии  $Y$  на  $X$  (теоретическое выражение).

В нашем случае, когда генеральные совокупности  $X$  и  $Y$  представлены конкретными наборами выборочных значений,  $\bar{y}_x$  является оценкой теоретической величины  $M(Y|X = x)$ , а уравнение (5.13) – выборочный аналог уравнения регрессии. Уравнение (5.13) называют *выборочным уравнением рег-*

рессии  $Y$  на  $X$ . Функция  $\varphi^*(x)$  - это выборочная регрессия  $Y$  на  $X$ , а график функции  $\varphi^*(x)$  - выборочная линия регрессии  $Y$  на  $X$ .

Совершенно аналогично выборочным уравнением регрессии  $X$  на  $Y$  является уравнение

$$\bar{x}_y = \psi^*, \quad (5.15)$$

где выборочные средние  $\bar{x}_y$  при различных значениях  $Y$  находятся из корреляционной таблицы 5.4, а уравнение (5.15) является выборочным аналогом теоретического уравнения регрессии  $M(X|Y = y) = \psi(y)$ .

Аналогично табл. 5.5 для выборочных значений  $\bar{x}_y$  можно, пользуясь корреляционной таблицей 5.4, построить табл. 5.6:

Т а б л и ц а 5.6. Условные средние значения  $x$

$y$	$y_1$	$y_2$	...	$y_l$
$\bar{x}_y$	$\bar{x}_{y_1}$	$\bar{x}_{y_2}$	...	$\bar{x}_{y_l}$
$n_y$	$n_{y_1}$	$n_{y_2}$	...	$n_{y_l}$

Из выражения (5.13) видно, что  $X$  и  $\bar{y}_x$  связаны некоторой функциональной зависимостью. Аналогично согласно выражению (5.15) –  $Y$  и  $\bar{x}_y$  также связаны соответствующей функциональной зависимостью. В то же время непосредственно  $X$  и  $Y$  имеют зависимость лишь корреляционную (если  $\varphi(x)$  и  $\psi(y)$  отличны от  $const$ ). Задачей исследователя является нахождение такой функциональной зависимости, которая бы наиболее адекватно отвечала реальным данным.

С помощью таблиц 5.5 и 5.6, можно определить, какой вид функциональной зависимости связывает  $X$  и  $\bar{y}_x$  или  $Y$  и  $\bar{x}_y$ .

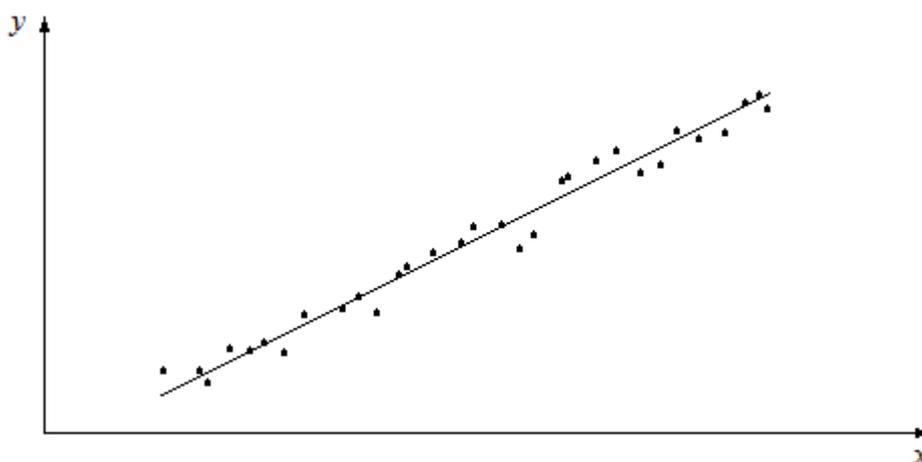
Далее следует подобрать коэффициенты функций  $\varphi^*(x)$ ,  $\psi^*(y)$  так, чтобы зависимость была «наилучшей» в смысле соответствия наблюдаемым значениям.

Для построения приемлемого уравнения регрессии также требуется знать степень корреляционной связи между случайными величинами  $X$  и  $Y$ , т. е. величину рассеяния значений относительно среднего:  $y$  относительно  $\bar{y}_x$  или  $x$  относительно  $\bar{x}_y$ .

Так как средние с ростом числа наблюдений обладают свойством стабилизироваться, нивелировать случайные отклонения, то регрессия представляет собой «истинные» значения без влияния различного рода случайных факторов. Таким образом, задача регрессионного анализа – определить приближенное уравнение регрессии и оценить допускаемую ошибку.

Важнейшим является вопрос выбора вида функции регрессии  $\varphi^*(x)$  (или  $\psi^*(x)$ ), например, линейной  $y = a + b x$  или нелинейной (показательной, логарифмической и т. д.).

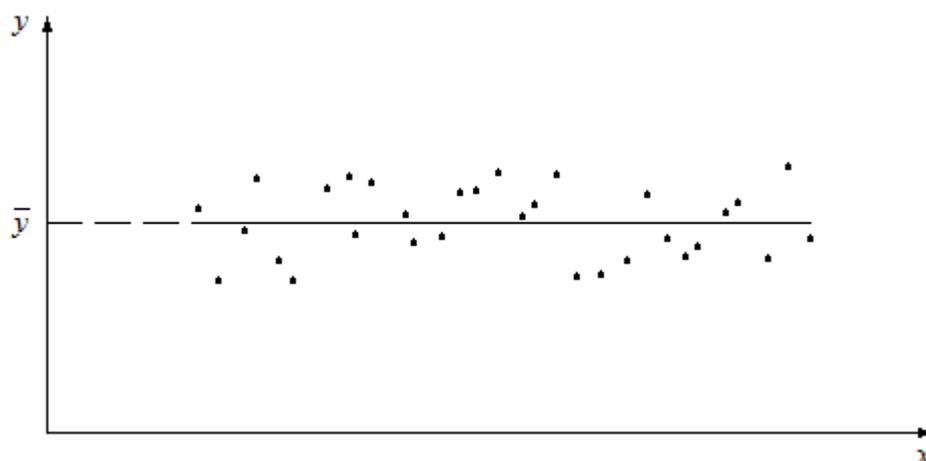
На практике вид функции регрессии можно определить, построив на координатной плоскости множество точек, соответствующих всем имеющимся парам наблюдений  $(x, y)$ . Например, на рис. 5.2 отчетливо видна тенденция роста значений  $y$  с ростом  $x$ , при этом средние значения  $y$  визуальнo располагаются на прямой. Поэтому целесообразно использовать линейную модель\* зависимости  $y$  от  $x$ .



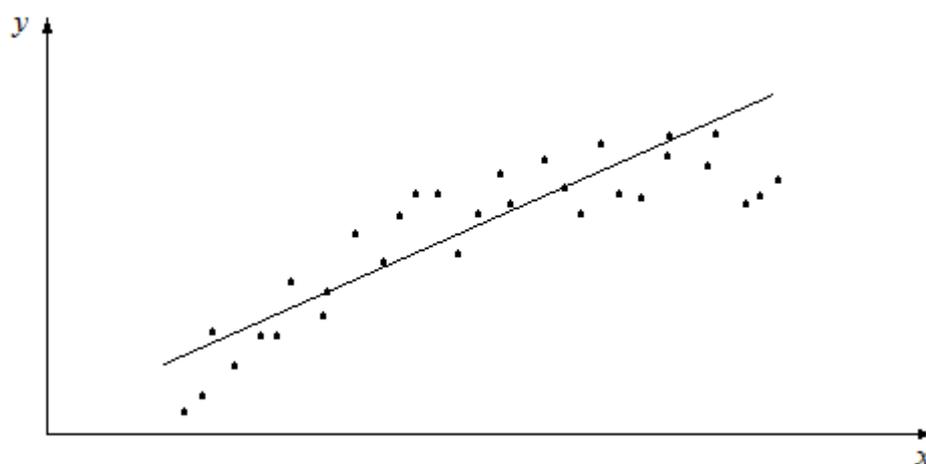
**Рис. 5.2.** Линейная регрессия значима. Модель  $Y = a + bX$

\* Отметим, что вид зависимости  $Y$  от  $X$  (или  $X$  от  $Y$ ) принято называть моделью. Далее будем придерживаться этой терминологии.

На рис. 5.3 средние значения  $y$  не зависят от  $x$ , следовательно, линейная регрессия незначима (функция регрессии постоянна и равна  $\bar{y}$ ).



**Рис. 5.3.** Линейная регрессия незначима. Модель  $Y = \bar{Y}$



**Рис. 5.4.** Линейная регрессия значима. Но желательно проверить нелинейную модель  $Y = aX^2 + bX + c$ .

### 5.6. Параметры выборочного уравнения регрессии при линейной зависимости

Выбрав вид функции регрессии, т. е. вид рассматриваемой модели зависимости  $Y$  от  $X$  (или  $X$  от  $Y$ ), например, линейную модель  $Y = bX + a$ , необходимо определить конкретные значения коэффициентов модели. Обратимся к линейной модели. При различных значениях  $a$  и  $b$  можно построить бесконечное число зависимостей вида  $Y = bX + a$  (т. е. на координатной плоскости

имеется бесконечное количество прямых), нам же необходима такая зависимость, которая соответствует наблюдаемым значениям наилучшим образом.

Как известно из [27, 29, 107 и др.], в случае модели линейной зависимости  $X$  и  $Y$  функция регрессии имеет вид

$$Y = m_y + r \frac{\sigma_y}{\sigma_x} (X - m_x)$$

или, что то же самое,

$$Y = \left( m_y - r \frac{\sigma_y}{\sigma_x} m_x \right) + r \frac{\sigma_y}{\sigma_x} X. \quad (5.16)$$

Мы же ищем линейную функцию  $a + bX$  исходя лишь из некоторого количества имеющихся наблюдений. Поэтому для нахождения функции с наилучшим соответствием наблюдаемым значениям используем метод наименьших квадратов, согласно которому находим коэффициенты  $a$  и  $b$  так, чтобы сумма квадратов отклонений наблюдаемых значений от значений на прямой линии регрессии оказалась наименьшей:

$$\sum_i (y_i - (a + bx_i))^2 - \min.$$

Чтобы отличать найденные наилучшие значения коэффициентов  $b$  и  $a$  от всех других обозначим их соответственно  $\beta$  и  $\alpha$ . Чтобы отличать значения  $y$ , относящиеся к уравнению регрессии, пометим их: будем обозначать  $\tilde{y}$ . Тогда выборочное уравнение линейной регрессии примет вид

$$\tilde{y} = \alpha + \beta x$$

или в другой записи

$$\tilde{y} = \bar{y} + \rho_{yx} (x - \bar{x}), \text{ где } \rho_{yx} = r_B \frac{s_y}{s_x}. \quad (5.17)$$

Коэффициент  $\rho_{yx}$  указывает, на какую величину изменится значение  $\tilde{y}$  в уравнении регрессии при изменении  $x$  на единицу.

Согласно выражению (5.17), каждому наблюдаемому значению  $x_i$  соответствует не только наблюдаемое  $y_i$ , но и значение  $\tilde{y}_i$ , удовлетворяющее уравнению регрессии  $\tilde{y} = \bar{y} + \rho_{yx}(x - \bar{x})$ .

Коэффициенты  $\alpha$  и  $\beta$ , полученные по выборочным данным, естественно, отличаются от соответствующих коэффициентов в уравнении регрессии генеральных совокупностей  $Y$  на  $X$ , т. е. являются их оценками. По этим оценкам можно проверить гипотезы о равенстве коэффициентов  $a$  и  $b$  в уравнении регрессии конкретным числам.

Замечание. Совершенно аналогично можно найти и уравнение регрессии  $X$  на  $Y$  того же вида, но, поменяв местами  $x$  и  $y$ , т. е. положив  $y$  в качестве независимой переменной, а  $x$  – функция от  $y$ . Тогда для уравнения регрессии получи:

$$\tilde{x} = \bar{x} + \rho_{xy}(y - \bar{y}), \text{ где } \rho_{yx} = r_B \frac{s_y}{s_x}. \quad (5.18)$$

В математической постановке задачи безразлично, какую из переменных выбрать независимой, а какую зависимой от нее, поскольку речь идет лишь о двух наборах чисел. Заметим, что уравнения (5.17) и (5.18), как правило, не совпадают. В практических задачах обычно имеет смысл лишь односторонняя зависимость. Например, для показателей возраста и заболеваемости населения логично рассматривать заболеваемость в зависимости от возраста. Обратное же соотношение, возраст в зависимости от заболеваемости, вообще говоря, бессмысленно, хотя математически и может быть найдено.

Итак, линия регрессии представляет множество средних значений. При этом легко вычислить характеристику разброса значений  $y$  относительно линии регрессии. Определим эту характеристику выражением

$$\tilde{s}^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \tilde{y}_i)^2 \quad (5.19)$$

Величина  $\tilde{s}^2$  называется *остаточной дисперсией* или *дисперсией относительно линии регрессии*. Эта величина – одна из оценок дисперсии слу-

чайной величины  $Y$ . Выбор коэффициентов  $\alpha$  и  $\beta$  в уравнении линейной регрессии с помощью метода наименьших квадратов гарантирует минимальное значение  $\tilde{s}^2$ .

Значение стандартного отклонения  $\tilde{s}^2 = \sqrt{\tilde{s}^2}$  называется *стандартной ошибкой оценки регрессии*. Очевидно, чем меньше ошибка, тем лучше регрессионная модель описывает реальную зависимость.

Поскольку коэффициенты  $\alpha$  и  $\beta$  находятся по конкретной выборке, для различных выборок они будут разными, т.е.  $\alpha$  и  $\beta$  являются случайными величинами, оценками истинных коэффициентов в уравнении регрессии. Естественно, что эти случайные величины имеют свои средние (истинные значения коэффициентов) и свои стандартные отклонения. Стандартные отклонения коэффициентов  $\sigma_\alpha$  и  $\sigma_\beta$  называются *стандартными ошибками коэффициентов регрессии*. Стандартные ошибки коэффициентов регрессии также используют при построении доверительных интервалов для оценки коэффициентов регрессии, для характеристики качества построенной модели. Повторимся: чем меньше значения стандартных ошибок (для оценки регрессии, коэффициентов регрессии), тем лучше модель. О «малости» найденных значений судят по их уровням значимости. В многочисленных пакетах компьютерных программ, используемых для вычисления регрессии, все эти характеристики обычно представлены.

### **5.7. Использование линейной регрессии в случае нелинейной зависимости**

Пусть имеются числовые данные парных наблюдений  $(x, y)$ , всего  $n$  пар. Исследуется корреляционная зависимость  $Y$  от  $X$ . Построение линейной модели  $y = a + bx$  оказывается неэффективным, например, малое значение коэффициента корреляции (слабая зависимость между  $X$  и  $Y$ ) обесценивает построенную модель. В этом случае можно попытаться построить нелинейную модель, которая «внутренне линейна», а именно: следует преобразовать исходные данные так, чтобы построенная по новым данным линейная модель

оказалась лучше предыдущей. Поясним сказанное конкретными преобразованиями.

Предположим, что исходные данные записаны в виде таблицы

$x$	$x_1$	$x_2$	...	$x_i$	...	$x_n$
$y$	$y_1$	$y_2$	...	$y_i$	...	$y_n$

Вначале преобразуем данные независимой переменной  $X$ .

Обратное преобразование. В строке значений  $x$  запишем величины  $1/x_i$  (естественно, в случае, когда все  $x_i$  отличны от нуля), тогда преобразованные данные имеют вид

$1/x$	$1/x_1$	$1/x_2$	...	$1/x_i$	...	$1/x_n$
$y$	$y_1$	$y_2$	...	$y_i$	...	$y_n$

Соответствующая линейная модель, построенная по этим данным,  $y=a+b/x$ , оказывается относительно  $x$  нелинейной, но относительно величины  $1/x$  эта модель, конечно же, линейна, что и позволяет для нахождения коэффициентов  $a$  и  $b$  использовать приведенную выше теорию линейной регрессии. Этим и объясняется термин «внутренне линейна», характеризующий регрессионную модель.

Возможно, эта модель будет более удачной, чем предыдущая.

Логарифмическое преобразование. Если все наблюдаемые значения  $x_i$  положительны, то вместо исходных данных  $x_i$  можно взять величины  $\ln x_i$ , тогда преобразованные данные запишутся в виде

$\ln x$	$\ln x_1$	$\ln x_2$	...	$\ln x_i$	...	$\ln x_n$
$y$	$y_1$	$y_2$	...	$y_i$	...	$y_n$

и по ним можно построить линейную модель  $y = a + b \ln x$ , которая линейна относительно  $\ln x$ , но, конечно же, нелинейна относительно  $x$ .

Преобразование типа степени. Если все  $x_i$  положительны, то вместо исходных данных  $x_i$  можно взять значения  $x_i^c$ , где  $c$  – некоторая постоянная (например, при  $c = 3$  все  $x_i$  возводятся в куб, а при  $c = 0,5$  из  $x_i$  извлекается квадратный корень и т. д.). Тогда преобразованные данные запишутся в виде

$x^c$	$x_1^c$	$x_2^c$	...	$x_i^c$	...	$x_n^c$
$y$	$y_1$	$y_2$	...	$y_i$	...	$y_n$

а соответствующая линейная модель окажется такой:  $y = a + bx^c$  ( $c$  - заранее заданная постоянная).

Преобразование отрицательных данных. Если среди наблюдаемых значений  $x_i$  имеются нули или отрицательные значения, то операция логарифмирования или извлечения квадратного корня для этих величин невозможна. В этом случае все  $x_i$  можно сместить на одно и то же число, а затем выполнить выбранную операцию логарифмирования, извлечения корня и т. д. Например, вместо  $(x_i : -1, 0, 1, 2)$  можно исследовать, добавив 2, последовательность  $(2 + x_i : 1, 2, 3, 4)$ , значения которой уже допускают, например, логарифмирование. Отметим, что смещение значений на одно и то же число равносильно добавлению к случайной величине неслучайного слагаемого, а эта операция, как известно, не изменяет ни дисперсию, ни корреляцию. Таким образом, построенная модель, линейная относительно новых данных, имеет вид

$$y = a + b \ln(2 + x).$$

Можно строить модели и преобразовывая зависимую переменную  $y$ . Например, модель  $y = ax^b$  легко сводится к линейной модели в результате операции логарифмирования обеих частей равенства:

$$\ln y = \ln a + b \ln x.$$

Следовательно, вместо значений  $x_i$  можно взять преобразованные величины  $\ln x_i$ , а вместо  $y_i$  – значения  $\ln y_i$ .

Замена  $y_i$  на  $1/y_i$  приводим к модели

$$y = \frac{1}{a + bx}.$$

Таким образом, метод построения линейной модели может быть использован и для построения нелинейных моделей. Разумеется, множество нелинейных моделей, которые, по сути, «внутренне линейны», не ограничивается рассмотренными выше случаями. Заинтересованный читатель без труда может предложить целый ряд других моделей.

### **5.8. Мера корреляционной связи. Выборочное корреляционное отношение**

Коэффициент корреляции  $r_{xy}$  позволяет определить силу линейной корреляционной связи. Если же зависимость между  $X$  и  $Y$  далека от линейной, то для определения тесноты корреляционной связи необходим другой более общий показатель – *коэффициент детерминации*.

Определение. Коэффициентом детерминации (причинности) называют число  $R^2$ :

$$R^2 = \frac{\sum_i n_{x_i} (\bar{y}_{x_i} - \bar{y})^2}{\sum_i n_{y_i} (y_i - \bar{y})^2}. \quad (5.20)$$

Коэффициент детерминации  $R^2$  указывает долю факторной дисперсии фактора  $X$  в общей дисперсии, т.е. численно выражает, какая часть вариации  $Y$  связана с воздействием фактора  $X$ .

Определение. Значение  $R = \sqrt{R^2}$  называют *выборочным корреляционным отношением*.

Сформулируем свойства корреляционного отношения  $R$ .

1.  $0 \leq R \leq 1$ .

2. При  $R = 0$  признаки  $X$  и  $Y$  не связаны корреляционной зависимостью.
3. При  $R = 1$  зависимость между  $X$  и  $Y$  функциональная.
4. Для любой выборки признаков  $X$  и  $Y$  справедливо соотношение:  $R \geq |r_B|$ , где  $r_B$  – выборочный коэффициент корреляции.

Если  $R = |r_B|$ , то корреляционная зависимость между  $X$  и  $Y$  – линейная, и все точки парных наблюдений  $(x_i, y_i)$  лежат на прямой линии регрессии.

Коэффициент детерминации, полученный в виде (5.20), исходя из выборочных данных, является оценкой истинного значения (теоретического) для коэффициента детерминации. Проверить значимость  $R^2$  можно, используя критерий Фишера. В этом случае нулевая гипотеза

$$H_0: R = 0$$

и соответственно конкурирующая гипотеза

$$H_1: R > 0.$$

В качестве критерия используется статистика

$$F = \frac{R^2}{1 - R^2} \cdot \frac{n - m}{m - 1}, \quad (5.21)$$

где  $m$  – число параметров уравнения регрессии. Статистика  $F$  имеет распределение Фишера с  $(m - 1; n - m)$  степенями свободы. Критическая точка находится по **табл. П.5**, а наблюдаемое значение критерия вычисляется непосредственно. Если наблюдаемое значение  $F$  превосходит критическое значение при заданном уровне значимости, то нулевая гипотеза отвергается, и величина  $R^2$  признается значимой. Отметим, что в компьютерных программах при нахождении уравнения регрессии вычисляется не только  $R^2$ , но и значение  $F$ , а также фактический уровень значимости этого значения.

При интерпретации  $R$  как показателя силы связи между  $X$  и  $Y$  можно применить качественную оценку зависимости. Для этих целей используется так называемая *шкала Чеддока* (табл. 5.7).

Т а б л и ц а 5.7. Шкала Чеддока

$R$	0,1 – 0,3	0,3 – 0,5	0,5 – 0,7	0,7 – 0,9	0,9 – 0,99
характеристика зависимости	слабая	умеренная	заметная	высокая	весьма высокая

При значениях  $R$ , меньших, чем 0,7 коэффициент детерминации  $R^2$  оказывается меньше 0,5, т. е. на долю вариации факторным признаком приходится менее 50% от влияния всех признаков, воздействующих на отклик. В этом случае полученная модель большой ценности не представляет, лучше попытаться найти другую модель (другой вид зависимости) с  $R^2$ , большим, чем 0,5, либо рассмотреть зависимость отклика от других факторов, которые могут оказаться более существенными. В реальных задачах при слабом влиянии факторов на отклик для получения общей картины зависимости желательно использовать модели множественной регрессии (см. раздел 5.10).

### 5.9. Простейшие случаи нелинейной регрессии

Выборочная функция регрессии подбирается в соответствии с имеющимися наблюдениями. Если график регрессии

$$\bar{y}_x = \varphi(x) \text{ или } \bar{x}_y = \psi(y)$$

не является прямой линией, т.е.  $\varphi(x)$  и  $\psi(y)$  – нелинейные функции, то и соответствующую корреляцию называют нелинейной. Наиболее распростра-

ненной нелинейной регрессией является параболическая модель (квадратическая\* регрессия).

$$\bar{y}_x = ax^2 + bx + c \text{ (или } y = ax^2 + bx + c). \quad (5.22)$$

Конкретные значения коэффициентов вычисляются, исходя из метода наименьших квадратов. Т.е. коэффициенты  $a$ ,  $b$ ,  $c$  определяются таким образом, чтобы сумма квадратов отклонений наблюдаемых значений от значений на линии регрессии была бы наименьшей:

$$\sum_i [y_i - (ax_i^2 + bx_i + c)]^2 = \min.$$

Вновь, как и в случае линейной регрессии, используя производные для нахождения экстремума функции, можно доказать, что коэффициенты  $a$ ,  $b$ ,  $c$  должны удовлетворять системе трех линейных алгебраических уравнений. Подставив найденные из системы параметры в (5.22), получим искомое уравнение регрессии.

Качество построенной нелинейной модели, также как и в линейном случае, характеризуют: коэффициент детерминации  $R^2$ , критерий  $F$ , уровень значимости  $F$ , стандартная ошибка оценки регрессии, стандартные ошибки коэффициентов регрессии.

Разумеется, можно попытаться построить и другие более пригодные для практического использования модели. При этом следует учесть, что на качество модели существенно влияет количество параметров модели. В частности, при двух наблюдениях в линейной модели линия регрессии (прямая) пройдет через обе точки (тогда все остатки равны нулю, а  $R^2 = 1$ ). При трех наблюдениях регрессионная кривая параболического типа пройдет через все три точки, и вновь при такой зависимости окажется  $R^2 = 1$ . Аналогично, как

---

\* Как легко заметить, квадратическая регрессия является частным случаем при  $n = 2$  полиномиальной регрессии  $\bar{y}_x = P_n(x)$ , где  $P_n(x)$  - полином степени  $n$ . Линейная регрессия - также частный случай полиномиальной при  $n = 1$ .

бы на плоскости  $XOY$  не располагались  $n$  точек, найдется кривая  $y = P_{n-1}(x)$ , описываемая полиномом порядка  $n-1$ , проходящая через все эти точки. Т.е. при увеличении количества параметров модели, сравнимым с количеством наблюдений, всегда можно получить модель зависимости с высоким значением  $R^2$ . Однако такие модели, не учитывающие случайный характер значений, не выявляющие возможные тенденции исследуемого процесса, и достаточно громоздкие и трудно интерпретируемые, по существу, лишь констатируют факт наличия конкретных значений. Для исследования и прогнозирования подобные модели особой ценности не представляют. Общепринято, что количество параметров модели должно быть меньше количества парных наблюдений как минимум в 4-5 раз. В прикладных задачах даже при большом количестве наблюдений, число параметров, обычно, полагают не выше четырех.

### **5.10. Методика построения модели множественной регрессии**

На начальной стадии исследования причинно-следственных связей между явлениями необходимо определиться, какой из рассматриваемых признаков будет фигурировать в качестве отклика (заболеваемость, продолжительность болезни, смертность и др.). Соответствующий показатель объективно должен аккумулировать в себе действие (изменение) всех остальных показателей, привлекаемых в модель зависимости в качестве факторов. Определившись с результативным признаком, откликом, далее, исходя из предшествующего опыта и здравого смысла, необходимо отметить возможные причины (факторы), влияющие на изменение результативного признака.

На следующем этапе исследования осуществляется переход от качественных признаков (причин и их следствий) к их количественным характеристикам. Зачастую, одно и то же явление количественно можно представить с разных сторон. Например, показатель заболеваемости населения может быть представлен как исчерпанная заболеваемость, как общая заболеваемость по обращаемости, как первичная заболеваемость, как накопленная заболеваемость.

мость, как частота заболеваний, выявленных при анализе причин смерти, как частота заболеваний, дополнительно выявленных при медицинских осмотрах и др. От правильного (информативного, объективно отражающего суть) выбора числовых характеристик явлений существенно зависит качество будущей модели, что наглядно было продемонстрировано в примере исследования зависимости показателя ЖЕЛ от стажа курильщика.

После сбора соответствующих числовых данных можно строить математическую модель зависимости переменной отклика ( $y$ ) от факторных переменных ( $x_1, x_2, \dots, x_m$ ). Выбираем подходящую функцию

$$y = f(x_1, x_2, \dots, x_m)$$

и строим модель. Заметим, что в случае линейной множественной регрессии, взяв вместо переменных  $y, x_1, x_2, \dots, x_m$  некоторые функции  $G(y), g_1(x_1), g_2(x_2), \dots, g_m(x_m)$ , теми же методами можно строить нелинейные модели относительно факторных переменных:

$$G(y) = a + b_1 g_1(x_1) + b_2 g_2(x_2) + \dots + b_m g_m(x_m).$$

Обратимся вновь к линейной модели. Прежде всего, из всех возможных факторных переменных необходимо выделить наиболее влиятельные по воздействию на отклик. Для этого находим множество коэффициентов корреляции показателей. На компьютере обычно выводится таблица, включающая все коэффициенты корреляции как между откликом и факторами, так и факторов между собой.

Далее: выделяем те факторы, которые наиболее сильно коррелируют с откликом  $y$ , их в первую очередь и включаем в модель. Начинаем с парной регрессии между  $y$  и той переменной  $x_k$ , корреляция которой с  $y$  наибольшая. Фиксируем соответствующие характеристики качества модели:  $R^2$ , стандартную ошибку регрессии и значимость коэффициентов модели.

Выделив указанные факторы, строим модель с двумя факторами: первый выбранный фактор остается в модели, а на место другого пробуем оставшиеся переменные. Среди двухфакторных моделей выбираем ту, которая дает ощутимый рост коэффициента детерминации  $R^2$ . Отметим, что  $R^2$  представляет собой долю вариации  $y$ , объясняемую изменением введенных в модель факторов. Следим за качеством модели по стандартной ошибке регрессии и значимости коэффициентов. При этом также учитываем коэффициент корреляции между факторами. Идеально, если факторы не коррелируют между собой, тогда каждый из коэффициентов модели характеризует вклад своего фактора в изменение показателя отклика. В противном случае, при сильной коррелированности факторов, происходит замещение факторов, что ведет к искажению истинной зависимости и неправильной интерпретации. Поэтому такие явно зависимые переменные вместе в модель включать нецелесообразно. Заметим, что в реальных задачах факторы, как правило, являются в той или иной степени коррелированными.

Выбрав подходящую модель с двумя факторами, оставляем их в модели и пробуем улучшить ее включением еще одного из оставшихся факторов. Как и ранее, следим за качеством модели. Может оказаться, что включение в модель новой переменной обесценивает вклад какой-то из переменных, ранее включенных в модель. В этом случае желательно пересмотреть целесообразность присутствия в модели этой старой переменной и выбрать модель лучшего качества, например, с большим значением  $R^2$ .

Поступая аналогично далее, продолжаем процесс до тех пор, пока включение новых переменных приводит к существенному росту коэффициента детерминации  $R^2$ . Причем важно обращать внимание не только на высокое значение  $R^2$ , но и на значение скорректированного коэффициента детерминации  $\bar{R}^2$ , особенно при большом количестве факторов в модели. Наиболее качественную из моделей и используем для дальнейших исследований.

**Пример 5.4.** Проиллюстрируем приведенные рекомендации на конкретном примере значений основных показателей фтизиатрической службы (данные табл. 5.8).

Т а б л и ц а 5.8. Значения основных показателей по фтизиатрической службе

Район	Смертность на 100 тыс. насе- ления (чел.)	Заболеваемость на 100 тыс. насе- ления (чел.)	Болезненность на 100 тыс. насе- ления (чел.)	Эффективность медицинских ос- мотров (%)
	$y$	$X_1$	$X_2$	$X_3$
1	13,3	76,4	281,9	47,7
2	11,8	53,3	278,3	71,4
3	12,4	62,2	273,1	50,1
4	9,2	73,3	214,2	57,3
5	17,2	72,7	294,5	40,3
6	10,7	53,7	254,1	61,3
7	12,3	55,6	253,3	47,8
8	24,4	95,4	301,8	39,2
9	16,3	73,7	285,3	39,4
10	10,0	55,3	236,2	63,6
11	11,4	75,8	231,8	72,7
12	23,1	89,1	279,4	28,5
13	12,0	52,6	275,4	58,6
14	9,8	50,6	248,8	58,6
	$\sum y = 193,9$	$\sum x_1 = 939,7$	$\sum x_2 = 3708,1$	$\sum x_3 = 736,5$

Рассматриваются показатели: смертность, заболеваемость, болезненность, эффективность медосмотров. При этом смертность – число больных, умерших в течение года от активных форм туберкулеза и состоявших на учете по данной территории на 100 тыс. среднегодовой численности населения; заболеваемость – число больных всеми формами активного туберкулеза, выявленных впервые в отчетном году на 100 тыс. среднегодовой численности населения; болезненность – число больных всеми формами активного туберкулеза, состоящих на учете на конец отчетного года на 100 тыс. населения, эффективность медосмотров - % впервые выявленных больных. Имеются соответствующие данные по 14 районам, т.е. всего 14

наблюдений. Вполне логично, что из имеющихся показателей в качестве отклика выбирается именно показатель смертности. В данном случае показатели уже заданы, поэтому подготовительный этап отбора факторов для исследования и их числовых характеристик считаем завершенным и приступаем к построению модели.

При отборе факторов в модель начинаем рассмотрение со значений коэффициентов корреляции, которые приведены в табл. 5.9.

Т а б л и ц а 5.9. Коэффициенты корреляции признаков

	y	$x_1$	$x_2$	$x_3$
y	1			
$x_1$	0,80089	1		
$x_2$	0,72478	0,35921	1	
$x_3$	-0,79492	-0,62914	-0,60731	1

Как видно из таблицы, все три факторных переменных  $x_1$ ,  $x_2$ ,  $x_3$  обладают высокими значениями коэффициентов корреляции с откликом  $y$ . Причем, факторы  $x_1$  и  $x_2$  положительно коррелированы с  $y$ , а фактор  $x_3$  имеет отрицательную корреляцию с  $y$ . Из этих коэффициентов выбираем наибольший по модулю:  $r_{x_1y} = 0,80089$ . Следовательно, первый фактор отбирается в будущую модель. Соответствующая модель представлена уравнением

$$y = - 3,90032 + 0,26445x_1, \quad (5.23)$$

$$(с.о.) \quad (3,91251) \quad (0,05708)$$

при этом  $R^2 = 0,64143$ , стандартная ошибка оценки регрессии (со) равна 2,97046, значимость  $F = 0,00058$ . Здесь и далее числа в скобках это стандартные ошибки соответствующих коэффициентов в модели, а значимость  $F$  - характеристика адекватности модели.

Далее попытаемся улучшить эту модель введением еще одного фактора (переменная  $x_1$  должна остаться в уравнении). Какую же из двух возможных переменных,  $x_2$  или  $x_3$ , предпочесть? Вновь обратимся к таблице 5.9, из которой следует, что  $x_2$  менее коррелирована с имеющейся факторной переменной  $x_1$  (0,35921 против -0,62914), следовательно, у нее больше шансов улучшить модель. Итак, выбирая факторные переменные  $X_1$ ,  $X_2$  приходим к модели

$$y = -24,66497 + 0,20493x_1 + 0,09349x_2, \quad (5.24)$$

$$(с.о.) \quad (5,59995) \quad (0,0398) \quad (0,02245)$$

где  $R^2 = 0,86079$ ,  $\bar{R}^2 = 0,8355$ , стандартная ошибка оценки регрессии = 1,93317, значимость  $F = 0,000019$ . Включение в модель новой переменной  $x_2$  заметно повы-

сило коэффициент детерминации  $R^2$ , уменьшило стандартную ошибку, значимость коэффициентов также достаточно высока.

Попробуем еще улучшить модель, добавив показатель третьего фактора  $x_3$ . Получим

$$y = -12,9506 + 0,16941x_1 + 0,07489x_2 - 0,08376x_3, \quad (5.25)$$

$$(с.о.) (9,96096) (0,04588) (0,02532) (0,05997)$$

где  $R^2 = 0,88351$ ,  $\bar{R}^2 = 0,8486$ , стандартная ошибка оценки регрессии = 1,85470, значимость  $F = 0,00006$ .

Как видим, в данной модели по сравнению с предыдущей произошло некоторое улучшение качества модели в смысле увеличения  $R^2$  (чуть более 2%) и уменьшения стандартной ошибки оценки регрессии (с 1,93317 до 1,85470). Однако вычисление показывает, что коэффициенты при  $x_3$  и свободный член в уравнении незначимы. Возможно, это вызвано наличием корреляции между  $x_3$  и переменными  $x_1$  и  $x_2$ . Следовательно, в найденной модели с тремя факторами вклад третьего фактора может оказаться нулевым. Фактически произошло замещение влияния  $x_3$  переменными  $x_1$  и  $x_2$ . Действительно, без учета влияния  $x_3$  в модели с двумя факторами коэффициенты при  $x_1$  и  $x_2$  больше соответствующих коэффициентов модели с тремя факторами. Уменьшение коэффициентов связано с обратной зависимостью между  $x_3$  и другими показателями.

Сравнение скорректированных коэффициентов детерминации в двух последних рассматриваемых моделях (5.24) - (5.25) также показывает лишь незначительное увеличение  $\bar{R}^2$  (с 0,8355 до 0,8486), т.е. с вводом третьего фактора прирост  $R^2$  примерно на 2% в значительной степени объясняется не влиянием, а просто наличием этого фактора. Исходя из вышесказанного, третий фактор в модель лучше не включать, а остановиться на модели (5.24). Заметим, что модель (5.24) можно попытаться улучшить изменением вида зависимости с линейной на нелинейную. Например, можно рассмотреть квадратичную зависимость от какого-либо фактора (или от обоих). В частности, модель вида

$$y = a + b_1x_1 + c_1x_1^2 + b_2x_2$$

приводит к результату

$$y = 11,01826 - 0,74188 x_1 + 0,00676 x_1^2 + 0,07875x_2 ,$$

$$(с.о.) (12,87236) (0,32337) (0,00230) (0,01796),$$

где  $R^2 = 0,9254$ ,  $\bar{R}^2 = 0,9030$ , стандартная ошибка оценки регрессии = 1,4847, значимость  $F = 0,000006$ .

Таким образом, использование указанной методики и творческий подход к исследованию приводят к вполне приемлемым моделям зависимостей.

### **5.11. Примеры построения регрессионных моделей**

Приведем примеры использования регрессионного анализа для построения характерных моделей зависимости в здравоохранении. В этих примерах отчетливо проявляются методика и специфика именно в творческой части исследования, которая обязательно имеет место при разработке качественных адекватных регрессионных моделей. Формальные компьютерные расчеты срабатывают только в простейших случаях, а при обработке реальных данных, как правило, требуется отдельный специфический подход, будь то методы группировки, выбор факторов, обоснованное исключение резко выделяющихся наблюдений и т.п.

#### **5.11.1. Соотношения параметров физического развития детей**

Физическое развитие является важнейшей составляющей в оценке уровня здоровья населения. В свою очередь, соотношения размеров тела - одна из составляющих при изучении физического развития индивидуума [7]. Для детского населения основными показателями физического развития считаются следующие: длина тела, масса тела, окружность грудной клетки (ОГК). Наряду с общепринятыми, можно рассмотреть и другие антропометрические показатели физического развития детей различных возрастов. Нахождение зависимостей между отдельными показателями способствует формированию стандартов и выявлению случаев, требующих коррекции или лечения. Понятно, что функциональной зависимости между значениями любой пары параметров нет. В данном случае зависимость может быть только статистической (корреляционной).

Опираясь на статистические данные порядка 200 наблюдений детей каждого из рассматриваемых возрастов, укажем приемлемые регрессионные

соотношения между параметрами физического развития детей\*, а также зависимости некоторых параметров от возраста [93]. Рассмотрим ряд моделей, как линейной, так и нелинейной структуры.

#### а) Дети 0 – 12 мес.

Все модели, приведенные ниже, адекватны и, как правило, имеют коэффициент детерминации не ниже 0,8, т.е. изменение признака-отклика более чем на 80% объясняется изменением признаков-факторов.

##### 1. Зависимость длины тела ( $l$ , см) от возраста ( $t$ , мес.)

$$\text{Мальчики: } l = 51,6387 + 3,1535 t - 0,0926 t^2.$$

$$\text{Девочки: } l = 51,4216 + 2,7438 t - 0,0680 t^2.$$

Вследствие квадратичной зависимости, прибавка в росте распределяется неравномерно: с увеличением возраста  $t$  рост замедляется. Если в первый месяц для мальчиков прирост составляет 3,0609см, то за 12-й месяц прирост всего 1,0237см.

Также приемлемыми являются линейные модели, рассчитываемые на основе постоянного среднемесячного прироста в течение года. Такие модели, в отличие от квадратичных, являются более грубыми, особенно вблизи границ интервала, зато они более просты в интерпретации.

$$\text{Мальчики: } l = 54,3113 + 1,9673 t.$$

$$\text{Девочки: } l = 53,6280 + 1,8396 t.$$

Согласно указанным моделям, прибавка в росте за 1 месяц в среднем в течение года составляет: у мальчиков – 1,97см, у девочек – 1,84см (значения коэффициентов при переменной  $t$ ). Данные соотношения показывают, что длина тела новорожденных мальчиков в среднем несколько больше, чем длина тела новорожденных девочек (на 0,7см.). Растут мальчики также быстрее, чем девочки, и уже к 1 году разрыв в среднем росте составляет около 2,2 см.

---

\* См. Медик В.А., Токмачев М.С. Соотношения параметров физического развития детей. - Сб. науч. тр. Новгородского научного центра Северо-Западного отделения РАМН - М.: Медицина, 2005.- Т.4.- С. 78 – 83.

Отметим еще один любопытный факт, связанный с длиной тела новорожденных. Группу из 200 наблюдений авторы разбили на две подгруппы, в зависимости от наличия заболеваний у матери в период беременности (117 наблюдений), и отсутствия таковых (83 наблюдения). Соответственно модели зависимости длины тела от возраста оказались следующими:

$$\text{при наличии заболеваний } l = 53,41162 + 1,97559 t;$$

$$\text{при отсутствии заболеваний } l = 55,02566 + 1,77328 t.$$

Следовательно, дети при заболеваемости матери в период беременности рождаются меньше ростом, чем дети, рожденные здоровой матерью. Однако в среднем растут такие дети быстрее и уже к 8 месяцам догоняют в росте своих сверстников.

## 2. Зависимость массы тела ( $m$ , кг) от длины тела ( $l$ , см)

$$\text{Мальчики: } m = -9,8536 + 0,2676 l.$$

$$\text{Девочки: } m = -10,6993 + 0,2822 l.$$

Согласно полученным моделям приращение массы тела на 1 см длины тела в среднем составляет: у мальчиков – 0,268 кг, у девочек – 0,282 кг. Так как в данной возрастной категории рост в значительной мере определяется возрастом (группа не однородна), то в приведенных моделях возраст косвенно учитывается через показатель роста: за показателем роста четко вырисовывается временной тренд.

## 3. Соотношение массы тела ( $m$ , кг), длины тела ( $l$ , см) и возраста ( $t$ , мес.)

$$\text{Мальчики: } m = -7,8919 + 0,0906 t + 0,2297 l$$

$$\text{Девочки: } m = -9,2813 + 0,0614 t + 0,2546 l$$

Из указанных моделей следует, что приращение массы тела у мальчиков на 1 месяц возраста составляет в среднем 0,091 кг, на 1 см длины тела – 0,230 кг. Однако точное соотношение этих приращений несколько условно, поскольку переменные возраста  $t$  и длины тела  $l$  имеют высокую корреляцию, которой и объясняется некоторое различие с предыдущими значениями приращений массы тела. Можно говорить лишь о тенденции: у мальчиков

первого года жизни приращение массы приближенно на 0,28 объясняется возрастом, и на 0,72 – ростом. У девочек тенденция аналогичная, но конкретные значения несколько иные: приращение массы приближенно на 0,19 объясняется возрастом, и на 0,81 – ростом ребенка.

#### 4. Зависимость окружности грудной клетки ( $l_{гр}$ , см) от массы тела ( $m$ , кг)

$$\text{Мальчики: } l_{гр} = 27,9292 + 1,9768 m.$$

$$\text{Девочки: } l_{гр} = 29,3863 + 1,7706 m.$$

Приращение окружности грудной клетки на 1кг приращенной массы в среднем составляет: у мальчиков – 1,98см, у девочек – 1,77см.

#### 5. Зависимость окружности грудной клетки ( $l_{гр}$ , см) от длины тела ( $l$ , см)

$$\text{Мальчики: } l_{гр} = 8,4978 + 0,5265 l.$$

$$\text{Девочки: } l_{гр} = 7,2227 + 0,5503 l.$$

Приращение длины окружности грудной клетки на 1 см прироста в среднем составляет у мальчиков – 0,53см, у девочек – 0,55см.

#### 6. Зависимость окружности головы ( $l_{гол}$ , см) от длины тела ( $l$ , см), массы тела ( $m$ , кг) и окружности грудной клетки ( $l_{гр}$ , см)

$$\text{Мальчики: } l_{гол} = 13,7158 + 0,4360 l, \quad R^2 = 0,7973;$$

$$l_{гол} = 30,3865 + 1,5504 m, \quad R^2 = 0,8216;$$

$$l_{гол} = 10,4846 + 0,7380 l_{гр}, \quad R^2 = 0,8303.$$

$$\text{Девочки: } l_{гол} = 15,2926 + 0,4105 l, \quad R^2 = 0,7627;$$

$$l_{гол} = 31,3076 + 1,4056 m, \quad R^2 = 0,7580;$$

$$l_{гол} = 11,1720 + 0,7181 l_{гр}, \quad R^2 = 0,7807.$$

Таким образом, для мальчиков приращение длины окружности головы в среднем составляет:

0,436см на 1см увеличения длины тела;

1,550см на 1кг приращения массы тела;

0,738см на 1см приращения длины окружности грудной клетки.

Разумеется, данные приращения не аддитивны, т.е. нет смысла рассматривать суммарное приращение.

Аналогично для девочек приращение длины окружности головы в среднем составляет:

0,411см на 1см увеличения длины тела;

1,406см на 1кг приращения массы тела;

0,718см на 1см приращения длины окружности грудной клетки.

### **7. Зависимость длины ног ( $l_{\text{ног}}$ , см) от длины тела ( $l$ , см)**

Мальчики:  $l_{\text{ног}} = -5,7904 + 0,5044 l$ ,  $R^2 = 0,5950$ .

Девочки:  $l_{\text{ног}} = -10,5256 + 0,5811 l$ ,  $R^2 = 0,6126$ .

Отметим, что в последних моделях значения коэффициента детерминации несколько ниже, чем в предыдущих. Т.е. значение показателя  $l_{\text{ног}}$  кроме длины тела существенно зависит и от других факторов, в данной модели не учтенных.

## **б) Дети 2-7 лет**

### **1. Зависимость роста ( $l$ , см) от возраста ( $t$ , мес.)**

Мальчики:  $l = 76,9272 + 0,5717 t$ ,  $R^2 = 0,8273$ .

Девочки:  $l = 71,4490 + 0,6244 t$ ,  $R^2 = 0,8384$ .

### **2. Зависимость массы тела ( $m$ , кг) от длины тела ( $l$ , см)**

Мальчики:  $m = -15,5546 + 0,3132 l$ ,  $R^2 = 0,8117$ .

Девочки:  $m = -11,9445 + 0,2737 l$ .  $R^2 = 0,8098$ .

### **3. Зависимость длины окружности грудной клетки ( $l_{\text{гр}}$ , см) от массы тела ( $m$ , кг)**

Мальчики:  $l_{\text{гр}} = 38,7840 + 0,9199 m$ ,  $R^2 = 0,7955$ .

Девочки:  $l_{\text{гр}} = 38,5148 + 0,9039 m$ ,  $R^2 = 0,8266$ .

### **4. Зависимость длины ног ( $l_{\text{ног}}$ , см) от длины тела ( $l$ , см)**

Мальчики:  $l_{\text{ног}} = -18,7293 + 0,6740 l$ ,  $R^2 = 0,8048$ .

Девочки:  $l_{\text{ног}} = -15,2292 + 0,6442 l, R^2 = 0,8849$ .

### **5. Зависимость ширины плеч ( $l_{\text{плеч}}$ , см) от длины тела ( $l$ , см)**

Мальчики:  $l_{\text{плеч}} = 3,8703 + 0,2334 l, R^2 = 0,6605$ .

Девочки:  $l_{\text{плеч}} = 5,0262 + 0,2194 l, R^2 = 0,7454$ .

Все приведенные выше модели адекватно отражают реальные данные и могут использоваться при изучении физического состояния детей.

В каждой регрессионной модели в качестве границ колеблемости признака-отклика общепринято использовать отклонения от теоретического среднего, кратные стандартному отклонению  $\sigma_R$ . На основе рассмотренных взаимосвязанных признаков могут быть разработаны нормативные оценочные таблицы. В настоящем сообщении приводятся оценочные таблицы некоторых показателей для детей в возрасте до 1 года (табл. 5.10 и 5.11). В качестве базового параметра используется показатель длины тела.

Общая оценка физического развития имеет свою градацию по совокупности признаков, распределенных по соответствующим сигмальным интервалам. Зависимости, указанные в работе, позволяют создать обновленные стандарты для оценки физического развития детей.

Площадь поверхности тела человека – это не только геометрическое понятие, примененное к конкретному объекту, а и важный показатель, используемый в медицине для стандартизации данных различных физиологических измерений. Показатель площади поверхности тела (BSA – body surface area) также задействован при расчете доз фармакологической нагрузки при ряде заболеваний, в частности, при ожоговых повреждениях. Имеется ряд показателей здоровья, использующих площадь поверхности тела, например, сердечный индекс, представляющий собой отношение сердечного выброса к площади поверхности тела человека.

Таблица 5.10. Оценочная таблица параметров физического развития для мальчиков первого года жизни

Рост	Рост (см)			Окружность грудной клетки (см)			Окружность головы (см)		
	$\bar{m} - \sigma_R$	$\bar{m}$	$\bar{m} + \sigma_R$	$\bar{l}_{гр} - \sigma_R$	$\bar{l}_{гр}$	$\bar{l}_{гр} + \sigma_R$	$\bar{l}_{гол} - \sigma_R$	$\bar{l}_{гол}$	$\bar{l}_{гол} + \sigma_R$
50	2,781	3,526	4,272	32,83	34,82	36,81	33,87	35,52	37,17
51	3,048	3,794	4,540	33,36	35,35	37,34	34,30	35,95	37,60
52	3,316	4,062	4,808	33,89	35,88	37,87	34,74	36,39	38,04
53	3,583	4,329	5,075	34,41	36,40	38,39	35,17	36,82	38,47
54	3,851	4,597	5,343	34,94	36,93	38,92	35,61	37,26	38,91
55	4,119	4,864	5,610	35,47	37,46	39,45	36,05	37,70	39,35
56	4,386	5,132	5,878	35,99	37,98	39,97	36,48	38,13	39,78
57	4,654	5,400	6,146	36,52	38,51	40,50	36,92	38,57	40,22
58	4,921	5,667	6,413	37,05	39,04	41,03	37,35	39,00	40,65
59	5,189	5,935	6,681	37,57	39,56	41,55	37,79	39,44	41,09
60	5,457	6,202	6,948	38,10	40,09	42,08	38,23	39,88	41,53
61	5,724	6,470	7,216	38,62	40,61	42,61	38,66	40,31	41,96
62	5,992	6,738	7,484	39,15	41,14	43,13	39,10	40,75	42,40
63	6,259	7,005	7,751	39,68	41,67	43,66	39,53	41,18	42,83
64	6,527	7,273	8,019	40,20	42,19	44,18	39,97	41,62	43,27
65	6,795	7,540	8,286	40,73	42,72	44,71	40,41	42,06	43,71
66	7,062	7,808	8,554	41,26	43,25	45,24	40,84	42,49	44,14
67	7,330	8,076	8,822	41,78	43,77	45,76	41,28	42,93	44,58
68	7,597	8,343	9,089	42,31	44,30	46,29	41,71	43,36	45,01
69	7,865	8,611	9,357	42,84	44,83	46,82	42,15	43,80	45,45
70	8,133	8,878	9,624	43,36	45,35	47,34	42,59	44,24	45,89
71	8,400	9,146	9,892	43,89	45,88	47,87	43,02	44,67	46,32
72	8,668	9,414	10,160	44,42	46,41	48,40	43,46	45,11	46,76
73	8,935	9,681	10,427	44,94	46,93	48,92	43,89	45,54	47,19
74	9,203	9,949	10,695	45,47	47,46	49,45	44,33	45,98	47,63
75	9,471	10,216	10,962	46,00	47,99	49,98	44,77	46,42	48,07
76	9,738	10,484	11,230	46,52	48,51	50,50	45,20	46,85	48,50
77	10,006	10,752	11,498	47,05	49,04	51,03	45,64	47,29	48,94
78	10,273	11,019	11,765	47,58	49,57	51,56	46,07	47,72	49,37
79	10,541	11,287	12,033	48,10	50,09	52,08	46,51	48,16	49,81
80	10,809	11,554	12,300	48,63	50,62	52,61	46,95	48,60	50,25

Таблица 5.11. Оценочная таблица параметров физического развития для девочек первого года жизни

Рост	Рост (см)			Окружность грудной клетки (см)			Окружность головы (см)		
	$\bar{m} - \sigma_R$	$\bar{m}$	$\bar{m} + \sigma_R$	$\bar{l}_{гр} - \sigma_R$	$\bar{l}_{гр}$	$\bar{l}_{гр} + \sigma_R$	$\bar{l}_{гол} - \sigma_R$	$\bar{l}_{гол}$	$\bar{l}_{гол} + \sigma_R$
50	2,605	3,411	4,216	33,11	34,74	36,36	34,30	35,82	37,33
51	2,888	3,693	4,498	33,66	35,29	36,91	34,71	36,23	37,74
52	3,170	3,975	4,780	34,21	35,84	37,46	35,12	36,64	38,15
53	3,452	4,257	5,063	34,76	36,39	38,01	35,53	37,05	38,56
54	3,734	4,540	5,345	35,31	36,94	38,56	35,94	37,46	38,97
55	4,016	4,822	5,627	35,86	37,49	39,11	36,35	37,87	39,39
56	4,299	5,104	5,909	36,41	38,04	39,66	36,77	38,28	39,80
57	4,581	5,386	6,191	36,96	38,59	40,22	37,18	38,69	40,21
58	4,863	5,668	6,474	37,51	39,14	40,77	37,59	39,10	40,62
59	5,145	5,951	6,756	38,06	39,69	41,32	38,00	39,51	41,03
60	5,427	6,233	7,038	38,61	40,24	41,87	38,41	39,92	41,44
61	5,710	6,515	7,320	39,16	40,79	42,42	38,82	40,33	41,85
62	5,992	6,797	7,602	39,71	41,34	42,97	39,23	40,74	42,26
63	6,274	7,079	7,885	40,26	41,89	43,52	39,64	41,15	42,67
64	6,556	7,362	8,167	40,81	42,44	44,07	40,05	41,56	43,08
65	6,838	7,644	8,449	41,37	42,99	44,62	40,46	41,98	43,49
66	7,121	7,926	8,731	41,92	43,54	45,17	40,87	42,39	43,90
67	7,403	8,208	9,013	42,47	44,09	45,72	41,28	42,80	44,31
68	7,685	8,490	9,296	43,02	44,64	46,27	41,69	43,21	44,72
69	7,967	8,773	9,578	43,57	45,19	46,82	42,10	43,62	45,13
70	8,249	9,055	9,860	44,12	45,74	47,37	42,51	44,03	45,54
71	8,532	9,337	10,142	44,67	46,29	47,92	42,92	44,44	45,95
72	8,814	9,619	10,424	45,22	46,84	48,47	43,33	44,85	46,36
73	9,096	9,901	10,707	45,77	47,39	49,02	43,74	45,26	46,77
74	9,378	10,184	10,989	46,32	47,94	49,57	44,15	45,67	47,18
75	9,660	10,466	11,271	46,87	48,49	50,12	44,56	46,08	47,60
76	9,943	10,748	11,553	47,42	49,04	50,67	44,98	46,49	48,01
77	10,225	11,030	11,835	47,97	49,59	51,22	45,39	46,90	48,42
78	10,507	11,312	12,118	48,52	50,14	51,77	45,80	47,31	48,83
79	10,789	11,595	12,400	49,07	50,70	52,32	46,21	47,72	49,24
80	11,071	11,877	12,682	49,62	51,25	52,87	46,62	48,13	49,65

Поскольку каждый измеряемый индивидуум имеет персонально нестандартную конфигурацию, то измерение площади поверхности человеческого тела – задача достаточно нетривиальная. В то же время, в реальной ситуации скрупулезная точность часто оказывается необязательной, и можно удовлетвориться приемлемой оценкой данной характеристики. Среди грубых методик отметим «правило ладони» и «правило девяток». Суть этих методик в следующем:

- практически вне зависимости от возраста площадь поверхности ладони приближенно равна 1% общей площади поверхности тела;

- площади поверхности основных частей тела взрослого человека приближенно оказываются кратными 9%, например, площади поверхности руки и головы составляют по 9%, спины – 18% от общей площади поверхности тела и т.д. (заметим, что для детей, особенно в младших возрастах, необходима модификация указанных соотношений).

Более точная оценка может быть получена статистически на основе других показателей, измеряемых достаточно просто. Известен ряд регрессионных формул [131 - 135], связывающих показатель площади поверхности тела ( $BSA$ ,  $m^2$ ) с показателями роста ( $L$ , см) и массы тела ( $M$ , кг):

$$BSA = 0,007184 \cdot L^{0,725} \cdot M^{0,425} \quad (\text{Дюбойс, 9 наблюдений}); \quad (5.26)$$

$$BSA = 0,017827 \cdot L^{0,5} \cdot M^{0,4838} \quad (\text{Бойд, 411 наблюдений}); \quad (5.27)$$

$$BSA = 0,008883 \cdot L^{0,663} \cdot M^{0,444} \quad (\text{Фуджимото, 201 наблюдение}); \quad (5.28)$$

$$BSA = 0,0235 \cdot L^{0,42246} \cdot M^{0,51456} \quad (\text{Гехан и Джордж, 401 наблюдение}); \quad (5.29)$$

$$BSA = 0,02465 \cdot L^{0,39646} \cdot M^{0,5378} \quad (\text{Хейкок, 81 наблюдение}). \quad (5.30)$$

Как легко заметить, все приведенные формулы имеют одну и ту же структуру:  $BSA = a \cdot L^b \cdot M^c$ . При этом, показатели степени,  $b$  и  $c$ , как правило, отличаются от 0,5 один в большую, другой в меньшую сторону. Полагая

$b = c = 0,5$  и регулируя изменившийся коэффициент  $a$ , можно получить [136] более простое соотношение:

$$BSA = 0,016667 \cdot \sqrt{L \cdot M} = \frac{\sqrt{L \cdot M}}{60} \quad (\text{Мостеллер, 401 наблюдение}). \quad (5.31)$$

Также отметим менее точную формулу Джексона, при целых  $L$  и  $M$  априори ограничивающую точность вычислений двумя знаками после запятой:

$$BSA = \frac{100 + M + (L - 160)}{100} = 0,01 \cdot (M + L - 60). \quad (5.32)$$

Можно предложить\* новые соотношения линейной и квадратичной структуры, выражающие площадь поверхности тела человека через его рост и массу тела:

$$BSA = -0,182637 + 0,006562 \cdot L + 0,012456 \cdot M \quad (\text{Токмачев 1}); \quad (5.33)$$

$$BSA = -0,185730 + 0,006607 \cdot L + 0,012388 \cdot M \quad (\text{Токмачев 2}); \quad (5.34)$$

$$BSA = -0,109464 + 0,005065 \cdot L + 0,017291 \cdot M - \\ - 0,000029 \cdot M^2 \quad (\text{Токмачев 3}); \quad (5.35)$$

$$BSA = -0,113447 + 0,005127 \cdot L + 0,017165 \cdot M - \\ - 0,000029 \cdot M^2 \quad (\text{Токмачев 4}); \quad (5.36)$$

$$BSA = -0,444759 + 0,009370(L + M) \quad (\text{Токмачев 5}). \quad (5.37)$$

Найденные зависимости являются уравнениями регрессии, построенными по данным 354 наблюдений возраста наблюдаемых от 3-х до 25 лет при росте индивидуумов от 100см до 200см и разнообразной массе. Некоторое различие в коэффициентах однотипных формул (5.33), (5.34), а также (5.35), (5.36), обусловлено различием в количестве знаков после запятой в наблюдаемых данных. Все приведенные соотношения характеризуются коэффициентом детерминации, близким к 1 (порядка 0,99...), и высокой значимостью коэффициентов регрессии. Это свидетельствует о качестве моделей и их адекватности реальным данным.

---

\* См. Медик В.А., Токмачев М.С. Соотношения параметров физического развития детей. - Сб. науч. тр. Новгородского научного центра Северо-Западного отделения РАМН - М.: Медицина, 2005.- Т.4.- С. 78 – 83.

Используя значение  $BSA_i$ , вычисляемое по конкретной формуле в каждом  $i$ -м наблюдении и соответствующее истинное значение площади поверхности тела  $s_i$ , сравним точность приведённых соотношений. Сравнение проводим по максимальным величинам абсолютного отклонения

$$\Delta_i = |BSA_i - s_i|, \quad (5.38)$$

относительного отклонения

$$\delta_i = \frac{\Delta_i}{s_i} = \frac{|BSA_i - s_i|}{s_i}, \quad (5.39)$$

а также соответствующих средних

$$\bar{\Delta} = \frac{1}{354} \sum_i \Delta_i, \quad (5.40)$$

$$\bar{\delta} = \frac{1}{354} \sum_i \delta_i. \quad (5.41)$$

Чем меньше значение рассматриваемой характеристики, тем выше точность применяемой формулы. Вычисленные по различным формулам значения оформим в виде таблицы с указанием ранга каждого из значений, для относительных величин дополнительно характеристики приведены в процентах (табл. 5.12).

Т а б л и ц а 5.12. Точность формул для BSA в зависимости от роста и массы

Формула	Дюбойс	Бойд	Фуджи-мото	Гехан и Джордж	Хейкок	Мостеллер
max $\Delta_i$ ранг	0,027499 (7)	0,024623 (4)	0,067022 (10)	0,034378 (8)	0,060662 (9)	0,021966 (3)
max $\delta_i$ (%) ранг	0,018840 1,88 (3)	0,027731 2,77 (5)	0,043765 4,38 (10)	0,038717 3,87 (6)	0,039563 3,96 (7)	0,025385 2,54 (4)
$\bar{\Delta}$ ранг	0,007130 (4)	0,007413 (7)	0,038254 (12)	0,012387 (8)	0,025456 (10)	0,005158 (3)
$\bar{\delta}$ $\bar{\delta}$ (%) ранг	0,005310 0,53 (4)	0,006063 0,61 (7)	0,026554 2,66 (11)	0,009622 0,96 (8)	0,016454 1,65 (10)	0,004181 0,42 (3)

Формула	Джексон	Токмачев 1	Токмачев 2	Токмачев 3	Токмачев 4	Токмачев 5
$\max \Delta_i$ ранг	0,113425 (12)	0,025903 (5)	0,026516 (6)	0,014566 (1)	0,016663 (2)	0,072129 (11)
$\max \delta_i$ (%) ранг	0,165966 16,60 (12)	0,041596 4,16 (8)	0,042581 4,26 (9)	0,007448 0,74 (1)	0,008091 0,81 (2)	0,054347 5,43 (11)
$\bar{\Delta}$ ранг	0,033422 (11)	0,00736 (6)	0,007249 (5)	0,003384 (2)	0,003346 (1)	0,017522 (9)
$\bar{\delta}$ $\bar{\delta}$ (%) ранг	0,030930 3,09 (12)	0,0058543 0,59 (6)	0,0058540 0,59 (5)	0,002209 0,22 (2)	0,002123 0,21 (1)	0,013219 1,32 (9)

Как видим из табл. 5.12, во всех случаях, за исключением формулы Джексона, точность вычислений оказывается вполне приемлемой. Однако следует отметить, что наилучшей точностью (наименьшими рангами) по всем параметрам обладают формулы квадратичной зависимости (5.35), (5.36), далее следует формула Мостеллера. Наихудшая точность с отклонением более чем на порядок у формулы Джексона. При этом заметим, что, несмотря на недопустимо большую относительную погрешность в отдельных наблюдениях, достигающую 16,6%, средняя относительная погрешность в формуле Джексона составляет всего 3,09%. А при значениях показателя роста более 130 см величина  $\delta_i$  не превышает 5%, т.е. формула Джексона пригодна лишь для ростового диапазона свыше 130см.

На практике площадь поверхности тела достаточно часто принято вычислять, исходя лишь из показателя массы тела. В частности, в программе STATISTICA, задействованы формулы:

$$BSA = -0,107 \cdot \sqrt[3]{M^2} \text{ (Дюбойс)}, \quad (5.42)$$

$$BSA = \frac{4M + 7}{M + 90} \text{ (Костефф)}, \quad (5.43)$$

где  $M$  – масса тела (кг),  $BSA$  - площадь поверхности тела ( $m^2$ ).

Согласно вышеупомянутым данным 354 наблюдений, получено регрессионное соотношение вида

$$BSA = 0,267011 + 0,027976 \cdot M - 0,000076 \cdot M^2. \quad (5.44)$$

Указанному уравнению регрессии соответствуют достаточно высокие значения коэффициента детерминации  $R^2 = 0,99532$  и стандартной ошибки 0,03204, что свидетельствует о хорошем качестве модели.

В табл. 5.13 приводятся результаты сравнения точность приведенных формул (5.42) – (5.44).

Таблица 5.13. Точность формул для BSA в зависимости от массы

Формула	Дюбойс	Костефф	Токмачев
$\max \Delta_i$ ранг	0,127911 (2)	0, 213011 (3)	0,101394 (1)
$\max \delta_i$ (%) ранг	0,078106 7,81 (2)	0,090446 9,04 (3)	0,071165 7,12 (1)
Продолжение табл. 5.13			
$\bar{\Delta}$ ранг	0,030264 (2)	0,050372 (3)	0,025687 (1)
$\bar{\delta}$ $\bar{\delta}_i$ (%) ранг	0,021155 2,12 (2)	0,031040 3,10 (3)	0,018464 1,85 (1)

Отметим, что значения, полученные по формуле Костефф, имеют существенные отклонения для роста свыше 170 см. Следовательно, применение формулы Костефф в этом случае может привести к ошибке. Средние значения  $\bar{\Delta}$  и  $\bar{\delta}$  во всех трех случаях вполне приемлемые. Однако точность формулы (5.44) наилучшая по всем параметрам.

Учитывая высокую корреляционную зависимость различных параметров физического развития, представим еще ряд регрессионных соотношений, связывающих площадь поверхности тела с другими показателями. Введем новые обозначения:

ОГК – длина окружности грудной клетки;

ОТ – длина окружности талии;

ВТ – высота талии от пола;

ст. – стандартная.

Во всех формулах показатели длины вычисляются в см, масса – в кг,  $BSA$  – в  $m^2$ .

$$BSA = 0,161459 + 0,005787 (\text{рост сидя}) + 0,013068 M, \quad (5.45)$$

где  $R^2 = 0,98600$ , ст. ошибка = 0,01851.

$$BSA = -1,507385 + 0,010131 L + 0,016964 (\text{ОГК}), \quad (5.46)$$

где  $R^2 = 0,868944$ , ст. ошибка = 0,056640.

$$BSA = -1,175086 + 0,011344 (\text{Рост сидя}) + 0,016308 (\text{ОГК}), \quad (5.47)$$

где  $R^2 = 0,823383$ , ст. ошибка = 0,065752.

$$BSA = -1,565112 + 0,012388 L + 0,015473 (\text{ОТ}), \quad (5.48)$$

где  $R^2 = 0,858623$ , ст. ошибка = 0,058828.

$$BSA = -0,788753 + 0,013120 (\text{ВТ}) + 0,014918 (\text{ОТ}), \quad (5.49)$$

где  $R^2 = 0,759884$ , ст. ошибка = 0,076666.

$$BSA = -1,621903 + 0,010995 L + 0,009756 (\text{ОГК}) + 0,007946 (\text{ОТ}), \quad (5.50)$$

где  $R^2 = 0,905495$ , ст. ошибка = 0,048220.

$$BSA = -1,264328 + 0,012328 (\text{рост сидя}) + \\ + 0,009541 (\text{ОГК}) + 0,007377 (\text{ОТ}), \quad (5.51)$$

где  $R^2 = 0,854964$ , ст. ошибка = 0,059736

Все представленные соотношения имеют достаточно высокие значения коэффициента детерминации и приемлемые величины стандартных ошибок, а сами модели зависимости адекватны реальным данным. Формулы (5.45) – (5.51) пригодны для практического применения. Указанные соотношения также могут оказаться весьма полезными при изучении пропорций тела человека.

### 5.11.3. Модели зависимости заболеваемости детей от степени загрязнения атмосферного воздуха

Северо-Западный регион России является одним из неблагоприятных по климатическим условиям (долгая зима, высокая влажность, перепады атмосферного давления, недостаток  $Ca^{++}$ ,  $Mg^{++}$ ,  $K^+$ , йода, фосфора, фтора, кобальта) и по экологической обстановке. Результаты эколого-гигиенических

исследований показали, что в атмосфере городов региона содержатся значительные примеси таких загрязнителей как пыль, оксиды азота, формальдегид, диоксид серы, оксид углерода, фенол [19]. Все это – факторы риска заболеваемости населения.

В последнее время появились исследования, свидетельствующие о наличии зависимости между состоянием здоровья людей и кратковременными и незначительными (в пределах ПДК) колебаниями концентраций аэрополлютантов в атмосфере. Публикуются данные об увеличении смертности, частоты госпитализации по поводу респираторных и сердечно-сосудистых заболеваний в зависимости от содержания отдельных загрязнителей в воздухе (пыли, аммиака, угарного и сернистого газа), увеличение уровня обращаемости за скорой медицинской помощью детей в возрасте до 14 лет [308]. Ближайшими эффектами негативного воздействия факторов окружающей среды являются аллергические заболевания, являющиеся мультифакторной патологией с ярко выраженной средовой компонентой [309].

Поставим задачу\* выявить связи между числом обращений детей различного возраста с обострениями бронхиальной астмы (БА) к специалисту в поликлинику, а также количеством вызовов педиатра на дом по поводу острой респираторной вирусной инфекции (ОРВИ) и уровнем содержания в воздухе максимально-разовых концентраций аэрополлютантов.

Обращаемость детского населения г. Великий Новгород в поликлинику изучалась ежедневно на основании данных статистических талонов в течение 2002 - 2003 гг. Содержание пыли, фенола, формальдегида, диоксида азота, аммиака, оксида углерода в атмосферном воздухе определяли в этот период подразделения Новгородского Областного Центра по Гидрометеорологии (НЦГМС). При химическом анализе проб воздуха использованы методики, изложенные в «Руководстве по контролю загрязнения атмосферы» РД 52.04.186-89. Отметим, что Великий Новгород не входит в число территорий,

---

\* См. Оконенко Т.И., Токмачев М.С., Вебер В.Р. Экологические подходы к оценке влияния загрязнения атмосферного воздуха на детей с заболеваниями дыхательной системы. - Экология человека. № 4. 2006. С. 6-9.

в которых средние показатели загрязнения атмосферного воздуха превышают предельно допустимые концентрации в несколько раз, хотя максимально-разовые показатели постоянно регистрируются на высоких значениях.

При подсчете коэффициентов корреляции в 2002 г. выявлены достаточно устойчивые связи между уровнем содержания в воздухе формальдегида, диоксида азота и обращаемостью детей с бронхиальной астмой в возрасте 10–14 лет, а также количеством вызовов педиатра на дом к заболевшим ОРВИ пациентам в возрасте 1–3 лет. Причем зависимость отмечалась между концентрацией ксенобиотика и уровнем обращаемости детей на следующий день после выброса.

Коэффициенты корреляции между концентрацией формальдегида и количеством обратившихся к аллергологу в возрасте 11–14 лет практически в течение всего года были более 0,5; лишь в зимние месяцы (январь, декабрь) отмечается менее тесная связь.

Прослеживается корреляционная связь и между содержанием в воздухе диоксида азота и числом обратившихся в поликлинику с бронхиальной астмой. Причем в этом случае зависимость регистрируется в более теплый период времени (апрель – октябрь). Известно, что в теплый период времени вследствие того, что воздух менее подвижен, концентрации токсических веществ в атмосфере увеличиваются, а токсичность формальдегида в присутствии диоксида азота возрастает, к тому же количество автотранспорта в этот период времени возрастает.

Здоровье человека определяется триадой, включающей факторы наследственности; факторы качества жизни (социально-экономическое и психологическое благополучие, доступность и качество медицинского обслуживания, образ жизни и наличие вредных привычек и др.); факторы состояния окружающей среды. Доля влияния отдельных факторов на состояние здоровья населения в среднем составляет: образ жизни (курение, употребление алкоголя и наркотиков, характер питания, материально-бытовые условия и др.)

– до 50 %, окружающая среда – 17–20%, состояние здравоохранения – 8–10 %, генетическая составляющая – около 20% [77].

Таким образом, влияние экологической составляющей (окружающая среда), оцениваемой 17–20%, весьма существенно и вдвое превышает соответствующую долю влияния фактора состояния здравоохранения. В то же время для регрессионной модели зависимость в 17–20% интерпретируется как слабая и установить такую зависимость гораздо сложнее, нежели зависимость порядка 90–100%. Здесь требуется четкая продуманная организация самого исследования.

Поскольку между временем выброса в атмосферу вредных веществ, обострением заболевания и обращением за поликлинической медицинской помощью проходит некоторое время, переменная-отклик  $y$  по отношению к соответствующим факторным переменным  $x_i$  во всех моделях берётся со сдвигом на 1 день.

Специфика построения соответствующих регрессионных моделей в том, что имеющиеся связи изучаются на фоне других объективно существующих факторов, не включаемых в модель. В данных условиях кроме временного сдвига значений факторов важную роль играет группировка наблюдений по степени однородности. Этот прием позволяет в некоторой степени нейтрализовать фоновые факторы и сосредоточиться именно на изучаемой связи типа «доза - эффект».

### ***Обращаемость по поводу ОРВИ (вызов врача на дом)***

Отметим, что зависимость обращаемости для детей всех возрастов всех возрастов (0–14 лет) от уровня загрязнения атмосферного воздуха явно не прослеживается, поскольку эта группа совершенно неоднородна, и имеются другие более мощные факторы, влияющие на уровень количества вызовов педиатра на дом.

Выделяя более локальные возрастные группы, получаем наличие зависимости лишь в группе 1–3 лет (376 наблюдений 2002–2003г.г.) и то доста-

точно слабой. Адекватной является линейная регрессионная модель с двумя факторами ( $y$ -уровень обращений, чел. в день;  $x_1$ -уровень концентрации  $\text{NO}_2$ ,  $\text{мг/м}^3$ ;  $x_2$ -уровень концентрации формальдегида,  $\text{мг/м}^3$ ):

$$y = 3,0980 + 11,2323x_1 + 42,8778x_2; \quad R^2 = 0,1126;$$

$$(c.o.) \quad (0,2252) \quad (2,3071) \quad (10,4342) \quad F = 23,6606.$$

Как видим, данная модель объясняет изменение числа обращений вследствие указанных факторов лишь на 11,26%.

Введение в модель фактора времени  $t$  существенно повышает значение коэффициента детерминации  $R^2$  и наблюдаемого значения  $F$ -критерия, а, следовательно, и качество модели:

$$y = 0,9184 + 0,3410t + 6,2585x_1 + 75,0481x_2; \quad R^2 = 0,4083;$$

$$(c.o.) \quad (0,2438) \quad (0,0250) \quad (1,9213) \quad (8,8516) \quad F = 85,5722.$$

В приведенной модели  $t$  - переменная, характеризующая месячную сезонность, в соответствии с изменением средней температуры воздуха:

Порядковый номер месяца	1	2	3	4	5	6	7	8	9	10	11	12
Значение $t$	12	10	9	7	5	3	1	2	4	6	8	11

Соответствующая линейная модель зависимости уровня обращений детей в возрасте 1–3 лет лишь от переменной сезонности  $t$  имеет вид:

$$y = 2,6938 + 0,3000t; \quad R^2 = 0,2536;$$

$$(c.o.) \quad (0,1809) \quad (0,0265) \quad F = 128,4059.$$

Та же модель, но для всех детей 0–14 лет:

$$y = 4,9160 + 1,6983t; \quad R^2 = 0,3644;$$

$$(c.o.) \quad (0,7878) \quad (0,1152) \quad F = 217,2990.$$

Как видим, влияние сезонности на формирование уровня обращений детей ОРЗ весьма существенно.

Далее попытаемся нейтрализовать переменную сезонности, используя только данные за время года со стабильно положительными температурами (97 наблюдений, май – август).

Вновь прослеживается влияние загрязнения атмосферного воздуха на число вызовов педиатра на дом лишь в группе детей в возрасте 1–3 лет. Значимыми оказываются факторы-переменные  $x_1$ -уровень концентрации  $\text{NO}_2$ ,  $\text{мг/м}^3$ ;  $x_2$ -уровень концентрации формальдегида,  $\text{мг/м}^3$ . Соответствующие регрессионные модели имеют вид:

а) линейная модель

$$y = 0,8462 + 23,6053x_1 + 61,6379x_2; \quad R^2 = 0,4588; \quad F = 39,8446;$$

(с.о.) (0,2895) (3,5408) (11,4028)

б) нелинейная модель

$$y = -1 + 2,1312e^{6,1896x_1} e^{15,9876x_2}; \quad R^2 = 0,4325; \quad F = 35,8127.$$

Сравнивая полученную линейную модель с соответствующей вышеприведённой моделью (376 наблюдений в течение всего года), отметим не только существенное повышение коэффициента детерминации  $R^2$  с 0,1126 до 0,4588, но и заметный рост влияния факторов на уровень обращений к педиатру. Увеличение значения  $x_1$  на 0,1  $\text{мг/м}^3$  при фиксированном  $x_2$  увеличивает уровень вызовов по поводу ОРВИ на 2,361 чел.; аналогично, увеличение значения  $x_2$  на 0,1  $\text{мг/м}^3$  при фиксированном  $x_1$  увеличивает его на 6,164 чел.

Персональные модели для факторов,  $x_1$ - уровень концентрации  $\text{NO}_2$ ,  $x_2$ -уровень концентрации формальдегида, следующие:

$$y = 1,7866 + 25,0798x_1; \quad R^2 = 0,2906;$$

(с.о.) (0,2636) (4,0206)  $F = 38,9114;$

$$y = 2,1226 + 67,5155x_2; \quad R^2 = 0,2031;$$

(с.о.) (0,2599) (13,6474)  $F = 24,4741.$

Последнюю модель можно улучшить введением квадратичного слагаемого:

$$y = 1,3451 + 159,7937x_2 - 2085,0438x_2^2; \quad R^2 = 0,2456;$$

(с.о.) (0,4216) (42,0930) (902,0064)  $F = 15,4623.$

Приведём модель зависимости уровня обращений в группе детей 1–3 лет от тех же факторов за летний период времени (38 наблюдений, с 26 мая по 30 июля):

$$y = 1,0543 + 18,8643x_1 + 2767,7493x_2^2; \quad R^2 = 0,5072;$$

(с.о.) (0,4154) (5,1087) (677,5421)  $F = 18,0105$ .

В этой модели использованы данные за более однородный по температуре промежуток времени. Таким образом, уменьшен вклад сезонности в уровень заболеваемости, что и привело к некоторому росту коэффициента детерминации  $R^2$ .

Также представим адекватную модель зависимости уровня обращений ОРВИ детей в возрасте 0-1 лет от уровня концентрации  $\text{NH}_3$  (28 наблюдений со сдвигом на 1 день с 26 мая по 31 июля в дни, когда был приём):

$$y = -1 + 6,8934 x^{0,3517}; \quad R^2 = 0,1648; \quad F = 5,1310.$$

### ***Обращаемость детей в связи с обострением бронхиальной астмы***

Зависимость уровня рассматриваемых обращений от уровня загрязнения атмосферного воздуха проявляется лишь в возрастном интервале 11-14 лет (переменная-отклик  $y$ ). Приведём ряд адекватных регрессионных моделей, характеризующих указанную зависимость. Полагаем  $x_1$  - уровень концентрации  $\text{NO}_2$ ,  $\text{мг/м}^3$ ;  $x_2$  - уровень концентрации формальдегида,  $\text{мг/м}^3$ ,  $x_3$  - уровень концентрации пыли,  $\text{мг/м}^3$ .

Тогда

$$y = 1,0578 + 4,1350x_1 + 85,3244x_2 + 0,5365x_3; \quad R^2 = 0,3618;$$

(с.о.) (0,2295) (1,9432) (10,7665) (0,2765)  $F = 28,9183$ ;

(157 наблюдений с января по декабрь);

$$y = 2,2958 + 7,6546 x_1; \quad R^2 = 0,0684;$$

(с.о.) (0,2088) (2,2689)  $F = 11,3819$ ;

(157 наблюдений с января по декабрь);

$$y = 1,4047 + 94,8425 x_2; \quad R^2 = 0,3284;$$

(с.о.) (0,1750) (9,9185)  $F = 91,4353$ ;

(189 наблюдений с января по декабрь);

$$y = 2,4891 + 0,9025 x_3; \quad R^2 = 0,0461;$$

(с.о.) (0,1551) (0,3001)  $F = 9,0408$ ;

(189 наблюдений с января по декабрь).

Как легко заметить, на рассматриваемом временном промежутке влияние факторов  $NO_2$  и «пыли», хотя и имеет место, весьма незначительно. Например, как следует из полученной модели зависимости  $y$  от  $x_3$ , повышение уровня обращений по поводу обострения бронхиальной астмы у детей пылевой фактор объясняет лишь на 4,61%. Указанный результат закономерно согласуется с данными [128] о крайне незначительном влиянии этого фактора на заболеваемость верхних дыхательных путей и лёгких у взрослых. В то же время для группы 11–14 лет следует отметить существенную зависимость  $y$  от фактора  $x_2$  ( $R^2=0,3284$ ): при росте концентрации формальдегида на  $0,1 \text{ мг/м}^3$  при прочих равных условиях уровень заболеваемости бронхиальной астмой возрастает на 9,48 чел.

Рассматривая заболеваемость бронхиальной астмой в летний период, находим в смысле интерпретации более качественные модели. В частности, при 43 наблюдениях, с 5 мая по 29 августа, для заболеваемости в возрастном интервале 11-14 лет (зависимая переменная  $y$ ) получаем

$$y = 0,7644 + 17,8978x_1 + 2480,1906x_2^2; \quad R^2 = 0,4519;$$

(с.о.) (0,3789) (5,3440) (553,0942)  $F = 16,4872$ ;

$$y = 1,3631 + 19,1429 x_1; \quad R^2 = 0,1763;$$

(с.о.) (0,4293) (6,4618)  $F = 8,7762$ ;

$$y = 1,1181 + 94,2187 x_2; \quad R^2 = 0,2767;$$

(с.о.) (0,3920) (23,7926)  $F = 15,6817$ ;

В представленных моделях, как и ранее,  $x_1$  - уровень концентрации  $NO_2$ ,  $\text{мг/м}^3$ ;  $x_2$  - уровень концентрации формальдегида,  $\text{мг/м}^3$ .

Таким образом, можно сделать следующие выводы:

- Выявлена достоверная корреляционная связь между содержанием в воздухе г. Великий Новгород формальдегида, диоксида азота и количеством обращений на следующий день к аллергологу с бронхиальной астмой и числом вызовов педиатра на дом к больным ОРВИ.

- Среди детей, страдающих бронхиальной астмой, наиболее метеочувствительными оказались больные в возрасте 11–14 лет. Среди заболевших ОРВИ – дети 1–3 лет.
- Построенные регрессионные модели, описывающие зависимость числа обращений от концентрации аэрополлютантов, подтверждают установленные связи и позволяют решать ряд проблем гигиенического характера и прогнозировать нагрузку на врачей детских поликлиник.

Итак, представленная в данном разделе методика указывает не столько на формально-математическую часть построения адекватных регрессионных моделей, сколько на способы формирования однородных данных с целью с целью отсекаемых внешних, не рассматриваемых в модели факторов.

#### **5.11.4. Регрессионные модели заболеваемости и смертности населения**

На основе базы данных заболеваемости и смертности населения Новгородского научного центра СЗО РАМН практически всё население Новгородской области можно классифицировать по половозрастным уровням и состояниям здоровья  $E_0, E_1, E_2, \dots, E_{23}$ . При этом состояние  $E_0$  вводится для лиц, не обращавшихся по поводу тяжёлых форм заболеваний, а состояния  $E_{20}, \dots, E_{23}$  соответствуют смертности по различным причинам за рассматриваемый период времени. Все остальные состояния характеризуют заболеваемость, классифицированы в соответствии с МКБ-10 и сформированы с учётом тяжести заболевания. При наличии нескольких заболеваний состояние  $E_j$  соответствует только одному из них, но наиболее тяжёлому. При имеющихся нескольких тяжёлых заболеваниях введена система приоритетов. Таким образом, например, в состоянии  $E_1$  - болезни системы кровообращения - на находятся не только индивидуумы с одними заболеваниями (системы кровообращения), и больные с целыми «букетами заболеваний», среди которых обязательно зафиксировано заболевание системы кровообращения в тяжёлой форме.

Оставляя в стороне общую картину заболеваемости и смертности, поставим более локальную задачу: нахождение по статистическим данным количественных зависимостей между смертностью, состояниями здоровья  $\{E_j\}$  и возрастом. Построим регрессионные модели для ряда возможных зависимостей. Нахождение уравнения регрессии по выборочным данным, по существу, является классической задачей с чётко разработанной методикой. В частности, в медицине сложилась практика нахождения регрессионной зависимости между заболеваемостью, как результирующим признаком, и конкретными факторными признаками, влияющими на формирование уровня заболеваемости, например, зависимость заболеваемости верхних дыхательных путей и лёгких от уровней содержания в воздухе концентраций аэрополлютантов. Суть каждого подобного исследования – включение в модель факторов, действительно связанных причинно-следственной зависимостью, сбор статистически корректных данных, характеризующих используемые факторы, и подбор типа модели (линейная, экспоненциальная и т. д.), адекватно отражающей существующую зависимость.

Исходя из имеющихся статистических данных, сформируем соответствующие последовательности наблюдений, классифицированные по половозрастным признакам, заболеваемости и смертности. Для факторов смертности, состояний здоровья (типа заболеваемости) и возраста приведём регрессионные модели наиболее характерных зависимостей. При этом набор наблюдаемых значений смертности и заболеваемости, связанных с возрастом  $t$ , по существу, будет временным рядом, а зависимость признака от времени  $t$  – его трендом.

Не вдаваясь в математические и технические подробности формирования состояний, адекватно отражающих реальные наблюдения (кроме первого случая), приведём некоторые конкретные модели зависимостей по данным мужской части населения. Используем временной интервал от 0 до 84 лет с периодом наблюдений в 4 года.

1. Зависимость смертности состоящих на учёте с болезнями системы кровообращения ( $E_1$ ) из-за этих болезней ( $E_{20}$ ) от возраста.

Соответствующая линейная модель\*, построенная для всех возрастов 0 – 84 лет, имеет большую стандартную ошибку коэффициентов. Построим модель с квадратичной зависимостью.

$$y = 0,01802 - 0,00082 t + 0,00002 t^2,$$

(с.о.) (0,00740) (0,00041) (0,00000)

стандартная ошибка  $y$  равна 0,01125;

значимость  $F$  равна  $2,27 \cdot 10^{-9}$ ;

$R^2 = 0,89045$ ;  $R = 0,94364$ .

Прокомментируем полученную зависимость.

Представленная модель в целом адекватна реальным данным и объясняет смертность по причине  $E_{20}$  на 89% фактором возраста. При этом рост смертности с увеличением возраста начинается примерно с 20-летнего возраста (точка минимума на параболе) с квадратичной зависимостью. До 20-летнего возраста, согласно модели, зависимость обратная: с возрастом смертность уменьшается. Этот артефакт объясняется весьма просто. У детей (точнее, до 20-ти лет) смертность для данного состояния практически наблюдается лишь в самых ранних возрастах, причём, согласно структуре наших данных, рассматриваются лишь те случаи, когда состояние  $E_1$  вначале зафиксировано, а смерть наступает в течение последующих 4-х лет (период наблюдений). Таким образом, при изучении данной зависимости вполне логично отбросить первые два наблюдения, которые почти полностью зависят от младенческой смертности. Для возраста 8–84 лет оказывается адекватной не квадратичная, а линейная модель зависимости:

$$y = -0,02616 + 0,00142t,$$

(с.о.) (0,00537) (0,00011)

---

\* Токмачев М. С. Регрессионные модели заболеваемости и смертности населения. - Сб. науч. тр. Новгородского научного центра Северо-Западного отделения РАМН; М.: Медицина, 2004. - Т.3.- С. 151 – 155.

стандартная ошибка  $y$  равна 0,01007;

значимость  $F$  равна  $1,75 \cdot 10^{-10}$ ;

$$R^2 = 0,91398 \quad R = 0,95602.$$

Как легко заметить, во-первых, новая модель по всем характеристикам превосходит предыдущую, во-вторых, линейная модель проще квадратичной, в-третьих, удалось избавиться от противоречивой зависимости.

Также отметим, что для возраста 20–84 лет в целом имеет место более качественная линейная модель

$$y = -0,03751 + 0,00161t,$$

(с.о.)    (0,00686)    (0,00012)

стандартная ошибка  $y$  равна 0,00917;

$$R^2 = 0,92245 \quad R = 0,96044.$$

2. *Зависимость общей смертности среди состоящих на учёте с болезнями системы кровообращения (т.е.  $E_{20} - E_{23}$  для группы  $E_1$ ) от возраста.*

$$y = 0,02501 - 0,00156t + 0,00006t^2,$$

(с.о.)    (0,00736)    (0,00041)    (0,00000)

стандартная ошибка  $y$  равна 0,01120;

значимость  $F$  равна  $3,33 \cdot 10^{-18}$ ;

$$R^2 = 0,98857; \quad R = 0,99427.$$

Модель полностью адекватна реальным данным и среди больных в состоянии  $E_1$  на 98,86% объясняет смертность фактором возраста.

3. *Общая зависимость смертности среди состоящих на учёте с новообразованиями (т.е.  $E_{20} - E_{23}$  для группы  $E_3$ ) от возраста.*

$$y = -0,03114 + 0,00440t,$$

(с.о.)    (0,01232)    (0,00023)

стандартная ошибка  $y$  равна 0,02821;

значимость  $F$  равна  $4,39 \cdot 10^{-13}$ ;

$$R^2 = 0,94028; R = 0,96968.$$

4. Зависимость заболеваемости системы кровообращения ( $E_1$ ) от возраста.

$$y = -0,06557 + 0,00384t,$$

(с.о.)      (0,01737)      (0,00036)

стандартная ошибка  $y$  равна 0,03977;

значимость  $F$  равна  $1,72 \cdot 10^{-9}$ ;

$$R^2 = 0,85792; R = 0,92623.$$

5. Зависимость смертности по причине болезней системы кровообращения ( $E_{20}$ ) и возраста  $t$ , где  $t \in [8, 83)$ .

$$y = -0,0062 + 0,21035 E_1 + 0,00050t,$$

(с.о.)      (0,00598)      (0,04896)      (0,00023)

стандартная ошибка  $y$  равна 0,00707;

значимость  $F$  равна  $6,48 \cdot 10^{-12}$ ;

$$R^2 = 0,95006; R = 0,97983.$$

6. Общая зависимость смертности больных с болезнями органов дыхания  $E_{20} - E_{23}$  для группы ( $E_4$ ) от возраста.

$$y = -0,04580 + 0,00255t,$$

(с.о.)      (0,01227)      (0,00025)

стандартная ошибка  $y$  равна 0,02809;

значимость  $F$  равна  $4,8 \cdot 10^{-9}$ ;

$$R^2 = 0,88418; R = 0,91753.$$

7. Зависимость общей смертности ( $E_{20} - E_{23}$ ) от заболеваемости системы кровообращения ( $E_1$ ), новообразований и возраста.

$$y = -235,79878 + 0,13424 E_1 + 1,16284 E_3 + 0,59852 t^2,$$

(с.о.)      (55,36412)      (0,03511)      (0,19903)      (0,05270)

стандартная ошибка  $y$  равна 138,5265;

значимость  $F$  равна  $8,26 \cdot 10^{-23}$ ;

$R^2 = 0,99782$ ;  $R = 0,99891$ .

Модели других зависимостей по имеющимся данным строятся аналогично.

В заключение отметим, что фактор, определяемый, как возраст, является весьма общим и ёмким. Он аккумулирует в себе всё множество возможных заболеваний, статистически не учтённых, но которые суммарно и определяют результативный признак: заболеваемость или смертность. Этим фактом всеохватывающего свойства возраста (времени)  $t$ , а также относительно малой доли других случайных причин, формирующих результативные признаки, и объясняются высокие значения коэффициента детерминации  $R^2$  моделях.

## **ГЛАВА 6. ПРОГНОЗИРОВАНИЕ ПОКАЗАТЕЛЕЙ ЗДОРОВЬЯ НАСЕЛЕНИЯ НА ОСНОВЕ “СТАНДАРТНЫХ” ПАРАМЕТРОВ**

## 6.1 Классические методы прогнозирования временных рядов

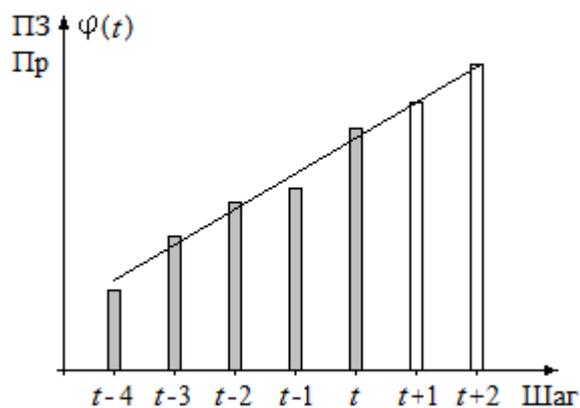
Последовательности ПЗ, ежегодно публикуемые ГОСКОМСТАТОм РФ, а также последовательности значений ИП, представляют собой одномерные временные ряды. Поэтому для обоснованного выбора и усовершенствования методов прогнозирования, адаптированных к закономерностям изменения соответствующих показателей здоровья, рассмотрим и проанализируем известные методы прогнозирования временных рядов.

Большинство классических методов прогнозирования [15, 46, 79, 120 и др.] основано на сглаживании известного участка временного ряда некоторой аппроксимирующей, сглаживающей его функцией  $\varphi(t)$ . Если рассматриваемый участок временного ряда насчитывает несколько его последовательных значений  $X(t)$ , то в большинстве случаев допустимо считать, что функция  $\varphi(t)$  является участком тренда данного ряда. Продолжая эту функцию от текущего дискретного момента (шага)  $t$ , на котором осуществляется прогнозирование, на последующие дискретные моменты времени, максимальное значение которых в иностранной литературе обычно называется временем упреждения ряда, получают значения прогнозов на предстоящих шагах. Известны методы прогнозирования и не использующие сглаживания временного ряда некоторой аппроксимирующей функцией.

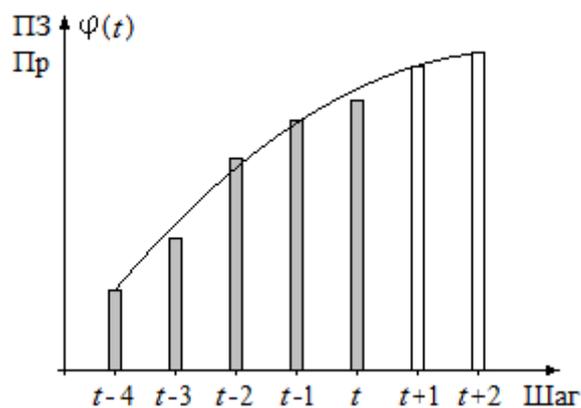
Выбор вида аппроксимирующей функции производится с учётом того, чтобы на рассматриваемом участке временного ряда обеспечить наименьшее значение ошибки аппроксимации, в качестве которой могут быть использованы среднее значение модуля отклонения этой функции на каждом шаге  $t$  интерполируемого участка от фактического значения элемента ряда  $X(t)$ , максимальное значение указанного модуля и др. Наиболее популярным методом сглаживания полученных значений  $x_t$  элементов временного ряда является [10, 15, 21, 25] метод наименьших квадратов, обеспечивающий наименьшую сумму квадратов отклонений  $\varphi(t) - x_t$  на рассматриваемом участке

временного ряда. Для получения значений прогнозов данный метод обычно используется совместно с методом скользящего окна, определяющим на каждом шаге участок ряда, по значениям которого определяется сглаживающая функция. Как правило этот участок заканчивается значением ряда, полученным последним. С получением каждого следующего значения временного ряда указанный участок сдвигается на один шаг вправо, т.е. как бы скользит по оси времени, и поэтому называется скользящим окном. Число элементов ряда, попадающих в скользящее окно, называется шириной этого окна. На каждом шаге  $t$  функция задаёт в текущем положении скользящего окна изменение усреднённого значения  $x_t$  обеспечивающего минимальное значение погрешности аппроксимации выбранного вида.

Проиллюстрируем на двух примерах решение задачи сглаживания временного ряда некоторого ПЗ и использования аппроксимирующей функции для получения прогнозов. Пусть на шаге  $t$  на основе попадающих в скользящее окно с пятью значениями  $ПЗ_{i-4}$ ,  $ПЗ_{i-3}$ ,  $ПЗ_{i-2}$ ,  $ПЗ_{i-1}$  и  $ПЗ_i$  ряда значений некоторого показателя здоровья (ширина окна равна пяти) строится сглаживающая функция  $\varphi(t)$  для этого участка. Вид этой функции выбирается в зависимости от того, каковы значения ПЗ, попавших в скользящее окно. Так, для участка ряда ПЗ, представленного на рис. 6.1, в качестве указанной функции логично взять линейную функцию  $\varphi(t) = at + b$ , а для участка ПЗ, приведённого на рис. 6.2, – параболу  $\varphi(t) = at^2 + bt + c$ , где  $a$ ,  $b$  и  $c$  – постоянные коэффициенты, которые нужно определить. Значения  $\varphi_{t+1}$  и  $\varphi_{t+2}$  являются прогнозами  $Пр_{i+1}$  и  $Пр_{i+1}$  на один и на два шага вперёд (на этих и последующих рисунках в отличие от ПЗ, значениями которых известны,  $Пр$  обозначаются прямоугольниками с белым внутренним фоном).



**Рис. 6.1.** Линейное сглаживание участка ряда и определение прогнозов



**Рис. 6.2.** Параболическое сглаживание участка ряда и определение прогнозов

На следующем шаге в рассмотренных примерах аппроксимирующие функции будут строиться по значениям  $\text{ПЗ}_{t-3}, \dots, \text{ПЗ}_{t-1}$ , попадающим в сдвинутое на один шаг вправо скользящее окно. При этом значение прогноза  $\text{ПР}_{t+2}$  может быть уточнено, так как в общем случае с увеличением интервала прогнозирования точность прогнозирования уменьшается. Последовательность прогнозов на каждом шаге  $t$  определяется, таким образом, на с помощью рекуррентных выражений, аргументами которых являются значения элементов ряда, попавших в скользящее окно.

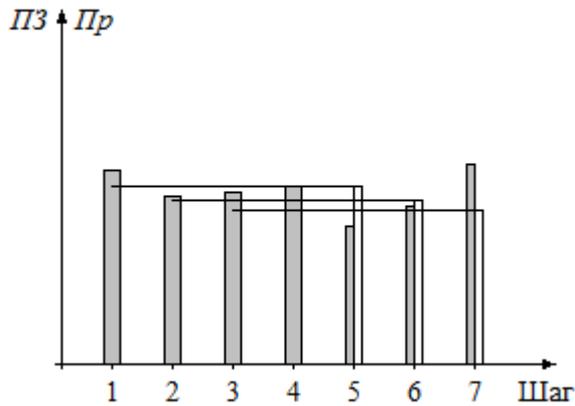
Получаемые на основе рассматриваемых методов прогнозы можно рассматривать как условные математические ожидания значений ряда, определённые при условии, что на каждом шаге  $t$  значения элементов ряда в скользящем окне принимали значения  $\varphi(t)$ . Поэтому в соответствии с используемой в математической статистике для таких функций терминологией (гл. 2) их уравнения часто называют регрессионными (регрессия – зависимость среднего значения какой либо величины от некоторой другой величины или от нескольких величин). С увеличением ширины скользящего окна аппроксимирующая функция приближается к участку тренда ряда на шагах, соответствующих этому окну. Если ширина окна равна, например, двум шагам, то ни о каком тренде говорить не приходится.

В последние годы в литературе выражения для расчета значений прогнозов чаще называют не методами прогнозирования, а прогнозирующими моделями или алгоритмами прогнозирования. Авторы считают логичным использовать в дальнейшем термин «модель прогнозирования» для задающих модель выражений в обобщённом виде, а термин «алгоритм прогнозирования» – для указанных выражений с конкретными значениями всех коэффициентов.

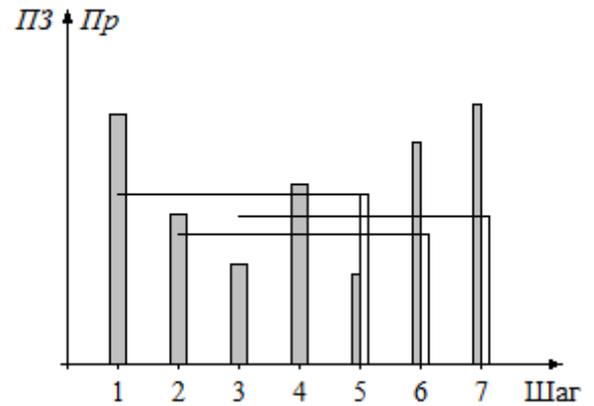
В зависимости от вида аппроксимирующей функции различают большое число моделей определения значений прогноза  $X_T$  на  $T$  шагов для временных рядов (полагаем  $T = t + 1$ ). Коэффициенты  $a, b, c, d, \dots$  являются величинами, зависящими от вида модели, от числа и конкретных значений элементов в скользящем окне, т.е. они изменяются с каждым смещением окна. В различных приложениях наиболее широко используются следующие модели прогнозирования:

- Скользящего среднего (скользящих констант, наивная модель):  $X_{t+T} = a$ .
- Линейная:  $X_{t+T} = aT + b$ .
- Параболическая (квадратичная):  $X_{t+T} = aT^2 + bT + c$ .
- Кубическая:  $X_{t+T} = aT^3 + bT^2 + cT + d$ .
- Экспоненциальная:  $X_{t+T} = ae^{bT} + c$ .
- Синусоидальная и косинусоидальная:  $X_{t+T} = a \sin \frac{2\pi(t+T)}{c} + b$ ,  $X_{t+T} = a \cos \frac{2\pi(t+T)}{c} + b$ .

Рис. 6.3 и 6.4 иллюстрируют динамику получения прогнозов некоторого ПЗ на примере модели скользящего среднего [129] с шириной окна  $m=4$  при различном значении отношения  $\sigma_X/M(X)$ , т.е. при различной величине относительного разброса значений  $X_t$ . С помощью этой модели выполняется одношаговое прогнозирование, т.е. прогнозирование на один шаг. По приведённым примерам не сложно заключить, что в общем случае с увеличением разброса значений элементов ряда точность прогнозирования уменьшается. Такой вывод справедлив для любых моделей прогнозирования временных рядов.



**Рис. 6.3.** Динамика прогнозирования значений временного ряда на основе модели скользящего среднего при  $\sigma_X/M(X) = 0,065$



**Рис. 6.4.** Динамика прогнозирования значений временного ряда на основе модели скользящего среднего при  $\sigma_X/M(X) = 0,115$

Поясним методику определения параметров  $a, b, c, \dots$  аппроксимирующей функции  $\varphi(t)$  (алгоритма выбранной модели) на основе метода наименьших квадратов. При этом будем предполагать, что прогнозирование на  $T$  шагов ( $T \geq 1$ ) осуществляется на шаге  $t_0$ , на котором получено значение  $\text{ПЗ}_0$ . Если выбранная функция  $\varphi(t)$  имеет  $n$  неизвестных параметров, значения которых на каждом шаге прогнозирования обычно изменяются, то для их определения на каждом указанном шаге необходимо использовать не менее  $n$  значений ПЗ, полученных последними. Поэтому для определения прогноза используется  $m \geq n$  последних ПЗ временного ряда, полученных на отрезке  $[t_0 - m + 1, t_0]$ , т.е. ширина скользящего окна выбирается равной  $m$ . Разность  $m - n$  выбирается исходя из особенностей динамики изменения ПЗ на аппроксимируемом интервале (можно ли достаточно точно рассматриваемый участок ряда аппроксимировать функцией с числом параметров меньше  $m$ ).

Согласно методу наименьших квадратов неизвестные параметры  $a, b, c, \dots$  аппроксимирующей функции, находятся исходя из условия [10, 25]:

$$\sum_{t=1-m}^0 [\text{ПЗ}_i - \varphi(t, a, b, c, \dots)]^2 = \min. \quad (6.1)$$

Для обеспечения минимального значения выражения (6.1) необходимо, чтобы все частные производные этого выражения по искомым параметрам были равны нулю. Поэтому дифференцируя выражение (6,1) по параметрам  $a, b, c, \dots$  и приравнивая получаемые выражения для производных нулю, получают следующую систему  $m$  уравнений с  $m$  неизвестными параметрами:

$$\left\{ \begin{array}{l} \sum_{t=1-m}^0 [ПЗ_i - \varphi(t, a, b, c, \dots)](\partial\varphi/\partial a)_t = 0, \\ \sum_{t=1-m}^0 [ПЗ_i - \varphi(t, a, b, c, \dots)](\partial\varphi/\partial b)_t = 0, \\ \dots\dots\dots \\ \dots\dots\dots \end{array} \right. \quad (6.2)$$

Если, например, используется линейная аппроксимация  $\varphi(t, a, b) = at + b$  при  $n = 2$ , то  $\partial\varphi/\partial a = t, \partial\varphi/\partial b = 1$ . В этом случае при  $m = 3$  система уравнений (6.2) после подстановки в неё значений  $t$  и частных производных принимает следующий промежуточный вид:

$$\left\{ \begin{array}{l} (ПЗ_{-2} + 2a - b) \times (-2) + (ПЗ_{-1} + a - b) \times (-1) + (ПЗ_0 - b) \cdot 0 = 0, \\ (ПЗ_{-2} + 2a - b) \times 1 + (ПЗ_{-1} + a - b) \times 1 + (ПЗ_0 - b) \cdot 1 = 0. \end{array} \right.$$

После приведения подобных членов получаем:

$$\left\{ \begin{array}{l} 5a - 3b = -2ПЗ_{-2} - ПЗ_{-1}, \\ -3a + 3b = ПЗ_{-2} + ПЗ_{-1} + ПЗ_0. \end{array} \right.$$

Следовательно,  $a = 0,5(ПЗ_0 - ПЗ_{-2}), b = (5ПЗ_0 + 2ПЗ_{-1} - ПЗ_{-2})/6$ . Согласно этим выражениям значения  $a$  и  $b$  следует находить на каждом шаге прогнозирования значений ПЗ.

Аналогично находятся параметры и других аппроксимирующих функций, если при применении метода наименьших квадратов их определения связано с решением системы линейных алгебраических уравнений.

## **6.2 Анализ точности прогнозирования показателей здоровья на основе полиномиальных моделей**

Модели скользящего среднего, линейная, параболическая и кубическая относятся к классу полиномиальных моделей, так как каждая из них представляет собой полином степени  $n$ , где  $n \geq 0$ . Для того, чтобы эти модели выполняли сглаживание ряда в скользящем окне, необходимо, чтобы значение  $n$  было бы меньше ширины  $m$  используемого скользящего окна хотя бы на 2, то есть выбирать модель с  $n \leq m - 2$ . Если же использовать полиномиальную модель на основе полинома степени  $n = m - 1$ , то в этом случае все значения ПЗ, входящие в окно, будут принадлежать аппроксимирующей функции, график которой пройдёт через соответствующие точки  $(t, ПЗ_t)$  окна. В этом случае функция  $\varphi(t)$  не будет сглаживающей. Но прогнозирование с помощью такой функции тоже возможно. Принимая, например, на шаге  $t$   $Пр_{t+1} = ПЗ_t$ , приходим к так называемой наивной модели [120], являющейся частным случаем модели скользящего среднего при  $n = 0$ ,  $m = 1$ .

Полиномиальные модели отличаются простотой определения их параметров и являются наиболее удобными для прогнозирования широкого класса временных рядов, в том числе и для временных рядов показателей здоровья. Поэтому читателю предлагаются результаты статистического исследования точности прогнозирования показателей здоровья с помощью полиномиальных моделей.

Алгоритмы прогнозирования для трёх полиномиальных моделей приведены в табл. 6.1 (алгоритмы для моделей скользящего среднего, линейной и параболической). При этом для каждой модели алгоритмы с наименьшим значением  $m$  не являются сглаживающими. Они реализуют функции, значения которых на принадлежащих скользящему окну шагах равны значениям ПЗ на

этих шагах. Остальные алгоритмы реализуют сглаживание значений элементов скользящего окна на основе метода наименьших квадратов.

Т а б л и ц а 6.1. Полиномиальные алгоритмы прогнозирования временных рядов на  $T$  шагов

Модель	$m$	Алгоритм
Скользящего среднего ( $\text{Пр}_T = a$ )	1	$\text{Пр}_T = \text{ПЗ}_0.$
	2	$\text{Пр}_T = 0,5(\text{ПЗ}_0 + \text{ПЗ}_{-1}).$
	3	$\text{Пр}_T = (\text{ПЗ}_0 + \text{ПЗ}_{-1} + \text{ПЗ}_{-2})/3.$
	4	$\text{Пр}_T = 0,25(\text{ПЗ}_0 + \text{ПЗ}_{-1} + \text{ПЗ}_{-2} + \text{ПЗ}_{-3}).$
	5	$\text{Пр}_T = 0,2(\text{ПЗ}_0 + \text{ПЗ}_{-1} + \text{ПЗ}_{-2} + \text{ПЗ}_{-3} + \text{ПЗ}_{-4}).$
Линейная ( $\text{Пр}_T = aT + b$ )	2	$\text{Пр}_T = (\text{ПЗ}_0 - \text{ПЗ}_{-1})T + \text{ПЗ}_0.$
	3	$\text{Пр}_T = 0,5(\text{ПЗ}_0 - \text{ПЗ}_{-2})T + (5\text{ПЗ}_0 + 2\text{ПЗ}_{-1} - \text{ПЗ}_{-2})/6.$
	4	$\text{Пр}_T = 0,1(3\text{ПЗ}_0 + \text{ПЗ}_{-1} - \text{ПЗ}_{-2} - 3\text{ПЗ}_{-3})T + 0,1(7\text{ПЗ}_0 + 4\text{ПЗ}_{-1} + \text{ПЗ}_{-2} - 2\text{ПЗ}_{-3}).$
	5	$\text{Пр}_T = 0,1(2\text{ПЗ}_0 + \text{ПЗ}_{-1} - \text{ПЗ}_{-3} - 2\text{ПЗ}_{-4})T + 0,2(3\text{ПЗ}_0 + 2\text{ПЗ}_{-1} + \text{ПЗ}_{-2} - \text{ПЗ}_{-4}).$
Параболическая ( $\text{Пр}_T = aT^2 + bT + c$ )	3	$\text{Пр}_T = 0,5(\text{ПЗ}_0 - 2\text{ПЗ}_{-1} + \text{ПЗ}_{-2})T^2 + 0,5(3\text{ПЗ}_0 - 4\text{ПЗ}_{-1} + \text{ПЗ}_{-2})T + \text{ПЗ}_0.$
	4	$\text{Пр}_T = 0,25(\text{ПЗ}_0 - \text{ПЗ}_{-1} - \text{ПЗ}_{-2} + \text{ПЗ}_{-3})T^2 + 0,05(21\text{ПЗ}_0 - 13\text{ПЗ}_{-1} - 17\text{ПЗ}_{-2} + 9\text{ПЗ}_{-3})T + 0,05(19\text{ПЗ}_0 + 3\text{ПЗ}_{-1} - 3\text{ПЗ}_{-2} + \text{ПЗ}_{-3}).$
	5	$\text{Пр}_T = (2\text{ПЗ}_0 - \text{ПЗ}_{-1} - 2\text{ПЗ}_{-2} - \text{ПЗ}_{-3} + 2\text{ПЗ}_{-4})T^2/14 + (54\text{ПЗ}_0 - 13\text{ПЗ}_{-1} - 40\text{ПЗ}_{-2} - 27\text{ПЗ}_{-3} + 26\text{ПЗ}_{-4})T/70 + (31\text{ПЗ}_0 + 9\text{ПЗ}_{-1} - 3\text{ПЗ}_{-2} - 5\text{ПЗ}_{-3} + 3\text{ПЗ}_{-4})/35.$

Приведённые модели и алгоритмы реализуют соответственно выборочные постоянные среднеквадратической регрессии  $\varphi$  на  $X$ , прямые среднеквадратической регрессии  $\varphi$  на  $X$ , параболы указанной регрессии и т. д.

Как уже указывалось, предпочтение той или иной модели отдают в зависимости от точности прогнозирования, которую может обеспечить модель. Однако на различных статистиках одна и та же модель или один и тот же алгоритм могут обеспечивать различную точность прогнозирования. Поэтому авторами было проведено статистическое исследование точности прогнозирования различных ПЗ и ИП на реальных статистических данных по соответствующим ПЗ и рассчитанным на их основе ИП здоровья населения. По результатам этого исследования даются рекомендации по использованию соответствующих моделей и алгоритмов прогнозирования рассматриваемых показателей. Исследована также точность прогнозирования соотношения среднедушевого дохода и прожиточного минимума (показатель социально-экономических условий жизни населения, влияющих на его здоровье). В процессе исследования в качестве погрешности прогнозов использовалось среднее относительное значение модуля разности  $Pr$  и фактического значения рассматриваемого показателя для одного и того же шага  $t$ .

При статистическом анализе точности прогнозирования временных рядов,  $n$  последовательных значений элементов которых известны, указанные разности не сложно определить, если прогнозирование осуществлять для шагов ряда с известными значениями элементов. В этом случае при ширине скользящего окна в  $m$  шагов можно получить  $n - m$  указанных разностей, что и делалось при проведении исследования. Кроме того, с целью установления влияния характеристик временных рядов на результаты прогнозирования были найдены значения оценок коэффициентов корреляции  $R(\tau)$  для временных рядов соответствующих ПЗ, приводимых в государственной статистике, и для рядов интегральных показателей здоровья, получаемых на основе многопараметрических линейных моделей (табл. 4.4). Ограниченность количества элементов отдельных исследованных рядов ПЗ обусловлена малым объемом статистических данных, приведённых в соответствующей литературе. Для обеспечения приемлемой достоверности выводы делались только для рядов с  $n > 20$ .

В табл. 6.2 ÷ 6.4 приведены указанные характеристики временных рядов различных ПЗ и ИП для РФ (регион с большой численностью населения), Северо-Западного федерального округа и Новгородской области (малая численность населения). При этом  $\overline{\delta\text{Пр}}(T) = |\text{Пр}(T) - \text{ПЗ}(T)| / \text{ПЗ}(T)$  – оценки относительных погрешностей прогнозов математических ожиданий соответствующих показателей на время от 1 до 4 лет с шагом 1 год (в таблицах приводятся средние значения этих оценок). Минимальные значения оценок этих погрешностей, соответствующие предпочтительным алгоритмам прогнозирования, выделены жирным шрифтом.

Таблица 6.2. Результаты статистического анализа точности прогнозирования показателей здоровья населения РФ

Общий коэффициент рождаемости (1980 ÷ 2007 гг, $R(1) = -0,018$ , $R(2) = -0,114$ )													
Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$		
Скользящего среднего	1	1	0,0513	Линейная	2	1	<b>0,0401</b>	Параболическая	3	1	0,1887		
		2	0,0920			2	<b>0,0852</b>			2	0,2772		
		3	0,1333			3	<b>0,1286</b>			3	0,4842		
		4	<b>0,1819</b>			4	0,1877			4	0,7722		
	2	1	0,0687		3	1	0,0451		4	1	0,0664		
		2	0,1117			2	0,0900			2	0,1686		
		3	0,1564			3	0,1357			3	0,2942		
		4	0,2095			4	0,2105			4	0,4680		
	3	1	0,0896		4	1	0,0542		5	1	0,0660		
		2	0,1340			2	0,0984			2	0,1446		
		3	0,1847			3	0,1591			3	0,2593		
		4	0,2391			4	0,2481			4	0,3955		
	4	1	0,1113		5	1	0,0621		Предпочтительные модели: линейная при $m = 2$ , $T < 4$ и скользящего среднего при $m = 1$ , $T = 4$				
		2	0,1615			2	0,1153						
		3	0,2148			3	0,1913						
		4	0,2679			4	0,2833						

Общий коэффициент смертности (1980 ÷ 2007 гг, $R(1) = -0,055$ , $R(2) = -0,019$ )														
Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$			
Скользящего среднего	1	1	<b>0,0388</b>	Линейная	2	1	0,0475	Параболическая	3	1	0,2304			
		2	<b>0,0680</b>			2	0,0935			2	0,3641			
		3	<b>0,0891</b>			3	0,1378			3	0,5358			
		4	0,1023			4	0,1919			4	0,6823			
	2	1	0,0500		3	1	0,0487		4	1	0,0716			
		2	0,0765			2	0,0971			2	0,1617			
		3	0,0919			3	0,1508			3	0,2821			
		4	<b>0,1022</b>			4	0,1960			4	0,4429			
	3	1	0,0618		4	1	0,0534		5	1	0,0659			
		2	0,0807			2	0,1058			2	0,1521			
		3	0,0927			3	0,1585			3	0,2580			
		4	0,0997			4	0,1988			4	0,3966			
	4	1	0,0670		5	1	0,0655		Предпочтительная модель: скользящего среднего при $m = 1$ для $T < 4$ и $m = 2$ для $T = 4$					
		2	0,0833			2	0,1161							
		3	0,0925			3	0,1594							
		4	0,1049			4	0,1874							
	Общий коэффициент младенческой смертности (1980 ÷ 2007 гг, $R(1) = -0,200$ , $R(2) = 0,114$ )													
	Модель	$m$	$T$		$\overline{\delta\text{Пр}}(T)$	Модель	$m$		$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$
	Скользящего среднего	1	1		<b>0,0417</b>	Линейная	2		1	0,0545	Параболическая	3	1	0,2504
			2		<b>0,0685</b>				2	0,0881			2	0,3556
3			<b>0,0939</b>	3	0,1230			3	0,5459					
4			<b>0,1198</b>	4	0,1408			4	1,0646					
2		1	0,0525	3	1		0,0448	4	1	0,0676				
		2	0,0811		2		0,0768		2	0,1477				
		3	0,1067		3		0,1030		3	0,2736				
		4	0,1323		4		0,1408		4	0,4568				
3		1	0,0639	4	1		0,0494	5	1	0,0573				
		2	0,0918		2		0,0746		2	0,1219				
		3	0,1207		3		0,0986		3	0,2251				
		4	0,1432		4		0,1215		4	0,3429				
4		1	0,0765	5	1		0,0495	Предпочтительная модель: скользящего среднего при $m = 1$						
		2	0,1060		2		0,0796							
		3	0,1332		3		0,0992							
		4	0,1504		4		0,1235							

Общая заболеваемость (1985 ÷ 2007 гг, $R(1)=-0,273$ , $R(2)=0,072$ )													
Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$		
Скользящего среднего	1	1	<b>0,0238</b>	Линейная	2	1	0,0273	Параболическая	3	1	0,5900		
		2	<b>0,0413</b>			2	0,0515			2	0,5980		
		3	<b>0,0585</b>			3	0,0734			3	0,8692		
		4	0,0789			4	0,0808			4	1,2511		
	2	1	0,0318		3	1	0,0253		4	1	0,0374		
		2	0,0496			2	0,0473			2	0,0766		
		3	0,0695			3	0,0700			3	0,1294		
		4	0,0904			4	0,0776			4	0,2044		
	3	1	0,0404		4	1	0,0261		5	1	0,0292		
		2	0,0597			2	0,0473			2	0,0601		
		3	0,0814			3	0,0652			3	0,1105		
		4	0,1010			4	<b>0,0749</b>			4	0,1517		
	4	1	0,0504		5	1	0,0292		Предпочтительные модели: скользящего среднего при $m = 1$ , $T < 4$ и линейная при $m = 4$ , $T = 4$				
		2	0,0719			2	0,0469						
		3	0,0925			3	0,0645						
		4	0,1116			4	0,0796						
Первичная инвалидность (1985 ÷ 2007 гг, $R(1) = -0,137$ , $R(2) = -0,216$ )													
Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$		
Скользящего среднего	1	1	<b>0,1135</b>	Линейная	2	1	0,1664	Параболическая	3	1	0,2730		
		2	<b>0,1466</b>			2	0,2477			2	0,6854		
		3	<b>0,1859</b>			3	0,3806			3	1,2702		
		4	<b>0,2130</b>			4	0,5456			4	2,0690		
	2	1	0,1180		3	1	0,1712		4	1	0,2395		
		2	0,1657			2	0,2334			2	0,4868		
		3	0,2007			3	0,3769			3	0,8231		
		4	0,2220			4	0,4387			4	1,3174		
	3	1	0,1420		4	1	0,1321		5	1	0,1765		
		2	0,1776			2	0,2444			2	0,3409		
		3	0,2130			3	0,3360			3	0,5837		
		4	0,2192			4	0,4416			4	0,8996		
	4	1	0,1543		5	1	0,1578		Предпочтительная модель: скользящего среднего при $m = 1$				
		2	0,1934			2	0,2411						
		3	0,2099			3	0,3330						
		4	<b>0,2310</b>			4	0,4162						

Средняя продолжительность предстоящей жизни (1985 ÷ 2007 гг, $R(1) = -0,046, R(2) = 0,072$ )														
Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$			
Скользящего среднего	1	1	<b>0,0105</b>	Линейная	2	1	0,0128	Параболическая	3	1	0,6522			
		2	<b>0,0188</b>			2	0,0270			2	0,7163			
		3	<b>0,0252</b>			3	0,0441			3	0,8102			
		4	0,0295			4	0,0658			4	1,0308			
	2	1	0,0147		3	1	0,0138		4	1	0,0189			
		2	0,0221			2	0,0303			2	0,0478			
		3	0,0267			3	0,0507			3	0,0931			
		4	0,0291			4	0,0680			4	0,1590			
	3	1	0,0180		4	1	0,0170		5	1	0,0204			
		2	0,0233			2	0,0365			2	0,0509			
		3	0,0274			3	0,0559			3	0,0994			
		4	0,0289			4	0,0722			4	0,1603			
	4	1	0,0198		5	1	0,0225		Предпочтительные модели: скользящего среднего при $m = 1$ для $T < 4$ и $m = 4$ для $T = 4$					
		2	0,0240			2	0,0411							
		3	0,0263			3	0,0588							
		4	<b>0,0281</b>			4	0,0704							
	Интегральный показатель общественного здоровья населения (модель № 1, 1985 ÷ 2007 гг, $R(1) = -0,137, R(2) = 0,0072$ )													
	Модель	$m$	$T$		$\overline{\delta\text{Пр}}(T)$	Модель	$m$		$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$
	Скользящего среднего	1	1		<b>0,0214</b>	Линейная	2		1	0,0293	Параболическая	3	1	0,4929
			2		<b>0,0329</b>				2	0,0493			2	0,6221
3			<b>0,0463</b>	3	0,0772			3	0,7629					
4			<b>0,0618</b>	4	0,1015			4	0,8092					
2		1	0,0244	3	1		0,0267	4	1	0,0350				
		2	0,0383		2		0,0479		2	0,0863				
		3	0,0537		3		0,0752		3	0,1593				
		4	0,0680		4		0,1103		4	0,2531				
3		1	0,0306	4	1		0,0293	5	1	0,0360				
		2	0,0462		2		0,0539		2	0,0722				
		3	0,0602		3		0,0806		3	0,1275				
		4	0,0780		4		0,1115		4	0,2077				
4		1	0,0387	5	1		0,0333	Предпочтительная модель: скользящего среднего при $m = 1$						
		2	0,0532		2		0,0569							
		3	0,0703		3		0,0855							
		4	0,0885		4		0,1028							

Интегральный показатель общественного здоровья населения (модель № 2, 1985 ÷ 2007 гг, $R(1) = -0,137$ , $R(2) = -0,024$ )													
Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$		
Скользящего среднего	1	1	<b>0,0203</b>	Линейная	2	1	0,0256	Параболическая	3	1	0,5030		
		2	<b>0,0327</b>			2	0,0437			2	0,6272		
		3	<b>0,0479</b>			3	0,0734			3	0,7369		
		4	<b>0,0627</b>			4	0,0976			4	0,7747		
	3	1	0,0255		3	1	0,0256		4	1	0,0322		
		2	0,0394			2	0,0486			2	0,0788		
		3	0,0557			3	0,0768			3	0,1422		
		4	0,0708			4	0,1088			4	0,2319		
	4	1	0,0321		4	1	0,0284		5	1	0,0333		
		2	0,0470			2	0,0525			2	0,0694		
		3	0,0639			3	0,0812			3	0,1242		
		4	0,0822			4	0,1096			4	0,2038		
	4	1	0,0393		5	1	0,0329		Предпочтительная модель: скользящего среднего при $m = 1$				
		2	0,0559			2	0,0576						
		3	0,0740			3	0,0841						
		4	0,0922			4	0,1042						

Таблица 6.3. Результаты статистического анализа точности прогнозирования показателей здоровья населения Северо-Западного федерального округа

Общий коэффициент рождаемости (1985 ÷ 2007 гг, $R(1) = -0,046$ , $R(2) = -0,048$ )													
Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$		
Скользящего среднего	1	1	0,0711	Линейная	2	1	<b>0,0620</b>	Параболическая	3	1	0,2483		
		2	<b>0,1331</b>			2	0,1353			2	0,4156		
		3	<b>0,1827</b>			3	0,2198			3	0,7800		
		4	<b>0,2309</b>			4	0,3103			4	1,2687		
	2	1	0,0961		3	1	0,0811		4	1	0,0968		
		2	0,1543			2	0,1589			2	0,2604		
		3	0,2060			3	0,2475			3	0,4838		
		4	0,2470			4	0,3528			4	0,6758		
	3	1	0,1241		4	1	0,0933		5	1	0,1160		
		2	0,1754			2	0,1768			2	0,2550		
		3	0,2190			3	0,2636			3	0,4268		
		4	0,2502			4	0,3429			4	0,6430		
	4	1	0,1427		5	1	0,1065		Предпочтительные модели: скользящего среднего при $m = 1$ , $T > 1$ и линейная при $m = 2$ , $T = 1$				
		2	0,1904			2	0,1954						
		3	0,2229			3	0,2829						
		4	0,2490			4	0,3315						

Средняя продолжительность предстоящей жизни (1985 ÷ 2007 гг, $R(1) = -0,068, R(2) = -0,120$ )													
Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$		
Скользящего среднего	1	1	0,0125	Линейная	2	1	<b>0,0115</b>	Параболическая	3	1	0,0586		
		2	<b>0,0217</b>			2	0,0263			2	0,6313		
		3	<b>0,0305</b>			3	0,0465			3	0,7529		
		4	<b>0,0359</b>			4	0,0690			4	0,8538		
	2	1	0,0170		3	1	0,0162		4	1	0,0193		
		2	0,0264			2	0,0356			2	0,0551		
		3	0,0324			3	0,0533			3	0,1057		
		4	0,0372			4	0,0746			4	0,1702		
	3	1	0,0207		4	1	0,0210		5	1	0,0237		
		2	0,0289			2	0,0404			2	0,0585		
		3	0,0343			3	0,0597			3	0,1143		
		4	0,0376			4	0,0799			4	0,1849		
	4	1	0,0241		5	1	0,0256		Предпочтительные модели: скользящего среднего при $m = 1, T > 1$ и линейная при $m = 2, T = 1$				
		2	0,0301			2	0,0442						
		3	0,0351			3	0,0633						
		4	0,0368			4	0,0771						
Интегральный показатель общественного здоровья населения (модель № 2, 1985 ÷ 2007 гг, $R(1) = -0,091, R(2) = -0,072$ )													
Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$		
Скользящего среднего	1	1	0,0399	Линейная	2	1	<b>0,0282</b>	Параболическая	3	1	0,3052		
		2	<b>0,0580</b>			2	0,0643			2	0,4159		
		3	<b>0,0799</b>			3	0,1097			3	0,5461		
		4	<b>0,1021</b>			4	0,1498			4	0,5916		
	2	1	0,0460		3	1	0,0459		4	1	0,0555		
		2	0,0686			2	0,0861			2	0,1376		
		3	0,0915			3	0,1258			3	0,2382		
		4	0,1156			4	0,1582			4	0,3331		
	3	1	0,0577		4	1	0,0548		5	1	0,0633		
		2	0,0805			2	0,0943			2	0,1329		
		3	0,1052			3	0,1241			3	0,2140		
		4	0,1262			4	0,1536			4	0,3219		
	4	1	0,0688		5	1	0,0632		Предпочтительные модели: линейная при $m = 1$ и скользящего среднего при $m > 1$				
		2	0,0935			2	0,0998						
		3	0,1165			3	0,1234						
		4	0,1292			4	0,1416						

Таблица 6.4. Результаты статистического анализа точности прогнозирования показателей здоровья населения Новгородской области

Общий коэффициент рождаемости (1975 ÷ 2008 гг, $R(1) = -0,121$ , $R(2) = -0,079$ )														
Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$			
Скользящего среднего	1	1	<b>0,0544</b>	Линейная	2	1	0,0584	Параболическая	3	1	0,1979			
		2	<b>0,0953</b>			2	0,1175			2	0,3694			
		3	<b>0,1335</b>			3	0,1590			3	0,7017			
		4	<b>0,1787</b>			4	0,2112			4	1,1608			
	2	1	0,0732		3	1	0,0563		4	1	0,0860			
		2	0,1144			2	0,1009			2	0,2173			
		3	0,1561			3	0,1383			3	0,3845			
		4	0,2056			4	0,1928			4	0,5964			
	3	1	0,0924		4	1	0,0631		5	1	0,0824			
		2	0,1362			2	0,1051			2	0,1724			
		3	0,1816			3	0,1472			3	0,2928			
		4	0,2336			4	0,2271			4	0,4319			
	4	1	0,1144		5	1	0,0685		Предпочтительная модель: скользящего среднего при $m = 1$					
		2	0,1599			2	0,1122							
		3	0,2092			3	0,1756							
		4	0,2607			4	0,2550							
	Интегральный показатель общественного здоровья населения (модель № 1, 1985 ÷ 2008 гг, $R(1) = -0,182$ , $R(2) = -0,120$ )													
	Модель	$m$	$T$		$\overline{\delta\text{Пр}}(T)$	Модель	$m$		$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$
	Скользящего среднего	1	1		<b>0,0567</b>	Линейная	2		1	0,0779	Параболическая	3	1	0,2403
			2		<b>0,0985</b>				2	0,1320			2	0,5070
3			<b>0,1385</b>	3	0,1990			3	0,8681					
4			<b>0,1937</b>	4	0,3076			4	1,3988					
2		1	0,0725	3	1		0,0758	4	1	0,1257				
		2	0,1171		2		0,1081		2	0,2354				
		3	0,1667		3		0,1631		3	0,4363				
		4	0,2274		4		0,2245		4	0,6897				
3		1	0,0945	4	1		0,0653	5	1	0,0910				
		2	0,1433		2		0,1073		2	0,1786				
		3	0,2009		3		0,1491		3	0,2669				
		4	0,2597		4		0,1901		4	0,4291				
4		1	0,1198	5	1		0,0804	Предпочтительная модель: скользящего среднего при $m = 1$						
		2	0,1756		2		0,1033							
		3	0,2308		3		0,1454							
		4	0,2914		4		0,1929							

Интегральный показатель общественного здоровья населения (модель № 2, 1985 ÷ 2007 гг, $R(1) = -0,182$ , $R(2) = -0,024$ )													
Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	Модель	$m$	$T$	$\overline{\delta\text{Пр}}(T)$		
Скользящего среднего	1	1	<b>0,0520</b>	Линейная	2	1	0,0510	Параболическая	3	1	0,2123		
		2	<b>0,0821</b>			2	0,0897			2	0,3815		
		3	<b>0,1104</b>			3	0,1399			3	0,5724		
		4	<b>0,1517</b>			4	0,2042			4	0,8596		
	2	1	0,0617		3	1	0,0589		4	1	0,0785		
		2	0,0927			2	0,0935			2	0,1567		
		3	0,1310			3	0,1205			3	0,2566		
		4	0,1715			4	0,1626			4	0,4145		
	3	1	0,0757		4	1	0,0583		5	1	0,0615		
		2	0,1129			2	0,0944			2	0,1181		
		3	0,1533			3	0,1229			3	0,1561		
		4	0,1998			4	0,1654			4	0,2676		
	4	1	0,0958		5	1	0,0685		Предпочтительная модель: скользящего среднего при $m = 1$				
		2	0,1355			2	0,0989						
		3	0,1799			3	0,1365						
		4	0,2271			4	0,1680						

Из результатов исследования точности прогнозирования значений ПЗ и ИП, приведённым в табл. 6.2 ÷ 6.4, следует, что среди предпочтительных по точности прогнозирования моделей отсутствует параболическая модель. Такой вывод можно было ожидать и по полученным значениям коэффициентов корреляции  $R(1)$  и  $R(2)$ . Малые значения их модулей и большой процент отрицательных значений указанных коэффициентов свидетельствуют о значительном количестве малых знакопеременных колебаний значений показателей рядов относительно текущих значений их математических ожиданий. Для такого случая для прогнозирования обычно наиболее подходят модель скользящего среднего и линейная модель, причём обычно при небольшой ширине скользящего окна.

Кроме относительной погрешности получения прогнозов рассматриваемых показателей, нередко для оценивания точности прогнозирования используют среднеквадратические оценки погрешностей. Согласно исследова-

ниям в этом случае выводы по сравнению методов и алгоритмов прогнозирования показателей здоровья практически не отличаются от приведённых выше.

Все проанализированные алгоритмы прогнозирования, были получены на основе метода наименьших квадратов при одинаковой точности определения значений элементов временного ряда, т.е. при одинаковой “степени доверия” этим значениям. В некоторых случаях, когда, например, погрешности получения значений отдельных элементов ряда больше, чем других, при выводе выражений для аппроксимирующих функций с использованием метода наименьших квадратов можно ввести весовые коэффициенты для элементов ряда [10, 120]. Это означает, что элементы ряда, значениям можно доверять больше, должны иметь и большие весовые коэффициенты, чем элементы ряда, значения которых вызывают определённые сомнения.

Рассмотрим пример. Пусть с помощью алгоритма скользящего среднего прогнозируется значение показателя исчерпанной заболеваемости в  $t + 1$ -м году по его значениям в  $t - 2 \div t$  годы, равным соответственно 1400, 2000 и 2000 заболеваний на 1000 человек населения. С помощью указанного алгоритма получаем:  $\text{Pr}(t + 1) = 1800$ . Однако исследование в  $t - 2$ -м году было недостаточно полным, и значение ПЗ, равное 1400, представляется сомнительным. Поэтому вводим весовые коэффициенты для значений показателей заболеваемости в скользящем окне с тремя элементами ряда. Примем сумму весовых коэффициентов равной ширине скользящего окна, т.е. трём, и выберем эти коэффициенты равными 0,6, 1,2 и 1,2. Тогда получаем:  $\text{Pr}(t + 1) = (0,6 \cdot 1400 + 1,2 \cdot 2000 + 1,2 \cdot 2000) / 3 = 1880$ .

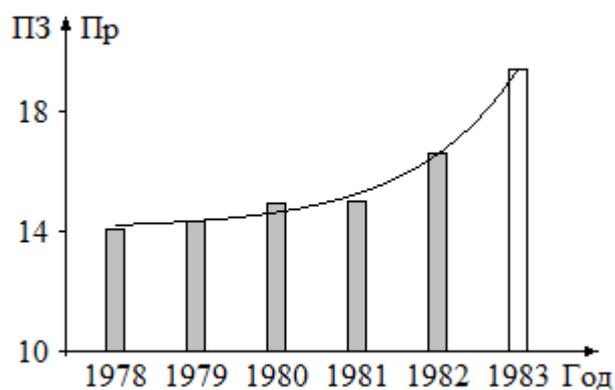
Следует однако заметить, что вопрос выбора значений весовых коэффициентов остаётся открытым: методов его решения не известно. Поэтому при сомнении в значениях отдельных элементов полученного временного ряда показателей здоровья целесообразнее воспользоваться методом, предлагаемым в § 6.4.

### 6.3 Прогнозирование показателей здоровья на основе “неполиномиальных” моделей

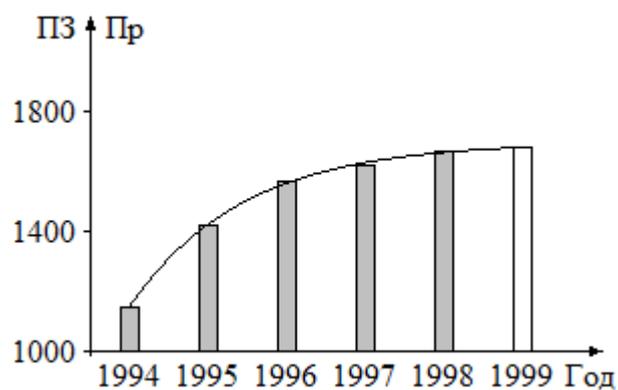
Все полиномиальные модели являются линейными относительно их параметров. Поэтому определение параметров таких моделей с помощью метода наименьших квадратов сводится к решению простой системы линейных алгебраических уравнений. Применение же метода наименьших квадратов для определения параметров экспоненциальной, синусоидальной и других моделей, имеющих один или более параметров, относительно которых значение выходной величины модели изменяется нелинейно, приводит к необходимости решения системы нелинейных алгебраических уравнений. В связи с указанным подбор параметров таких моделей проще выполнять путём моделирования или с использованием итерационных методов (методов последовательных приближений). При этом в качестве критерия качества сглаживания рассматриваемого участка ряда аппроксимирующей функцией  $\varphi(x)$  по-прежнему может служить значение суммы квадратов отклонений этой функции от элементов ряда.

Рассмотрим в качестве примера вариант решения данной задачи для экспоненциальной модели  $\text{ПЗ}_{t+T} = ae^{bt} + c$  с параметрами  $a$ ,  $b$  и  $c$ . Модель используется для прогнозирования значений некоторого ПЗ на  $T$  шагов ( $T > 0$ ). Она линейна относительно параметров  $a$  и  $c$  и нелинейна относительно параметра  $b$ . Параметры модели нужно выбрать так, чтобы минимизировать сумму квадратов (6.1) отклонений сглаживающей функции от значений ПЗ в скользящем окне. С высокой точностью это можно сделать с помощью ЭВМ, использовав один из методов поиска экстремальных значений функций [23]. На рис. 6.6 приведены примеры графиков, полученных разработанной программой на основе метода перебора значений параметров сглаживающей соответствующие участки временных рядов ПЗ населения РФ экспоненциальной функции  $\varphi(t) = ae^{bt} + c$ . Эти участки подобраны с учётом удобства сглаживания их именно экспоненциальной функцией, причём четырёх воз-

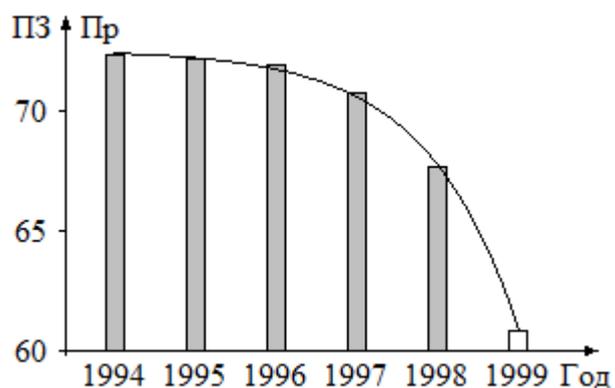
можных типов: возрастающей с увеличением темпа роста, возрастающей с сокращением темпа роста, уменьшающейся с возрастанием темпа уменьшения и уменьшающейся с сокращением темпа уменьшения. Полученные значения параметров приведены под рисунками. Относительная погрешность прогнозирования в рассмотренных случаях для  $T = 1$  оказалась равной 0,0351 (рис. 6.5 а), 0,0072 (рис. 6.5 б), 0,1069 (рис. 6.5 в) и 0,0509 (рис. 6.5 г).



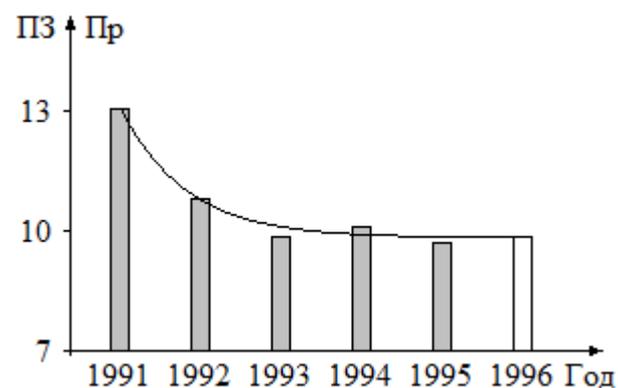
а) Общий коэффициент младенческой смертности  
( $a = 0,118$ ,  $b = 0,762$ ,  $c = 14,067$ )



б) Заболеваемость детского населения по обращениям  
( $a = -549,32$ ,  $b = -0,697$ ,  $c = 1692,54$ )



в) Средняя продолжительность предстоящей жизни  
( $a = -0,128$ ,  $b = 0,887$ ,  $c = 99,652$ )



г) Общий коэффициент рождаемости  
( $a = 2,758$ ,  $b = -16248$ ,  $c = 9,348$ )

**Рис. 6.5.** Прогнозирование показателей здоровья на основе экспоненциальной модели

Аналогичная методика сглаживания участков временных рядов ПЗ может быть применена и для нахождения оптимальных значений параметров и других нелинейных моделей: синусоидальной, логарифмической и т. д. При этом кроме метода перебора значений параметров можно использовать и

другие методы поиска оптимальных решений (градиентный метод, методы Монте-Карло и пр.). Однако в большинстве случаев для прогнозирования показателей здоровья предпочтительными являются полиномиальные модели нулевого и первого порядка (скользящего среднего и линейная). Их алгоритмы чрезвычайно просты (они не требуют оптимизации параметров модели на каждом шаге прогнозирования), а точность прогнозирования на  $1 \div 2$  шага получается удовлетворительной.

В целом можно заключить, что для улучшения точности прогнозирования необходимо на каждом шаге анализировать участок ряда, по которому будет определяться прогноз. Это позволит обоснованно выбрать тип модели прогнозирования для данного шага и ширину текущего скользящего окна, т.е. число последних элементов ряда, по которым находятся значения параметров выбранной модели и определяется прогноз.

#### **6.4. Прогнозирование при “нетипичных” выбросах значений показателей здоровья**

Временные ряды показателей здоровья населения могут иметь большие выбросы значений ПЗ (случайной составляющей ряда), которые явно не соответствуют их “типичным колебаниям”. Такие выбросы чаще всего являются следствием эпидемий, крупных катастроф различного характера и т.д. Однако они могут иметь место и в рядах ПЗ населения административных единиц с малой численностью населения.

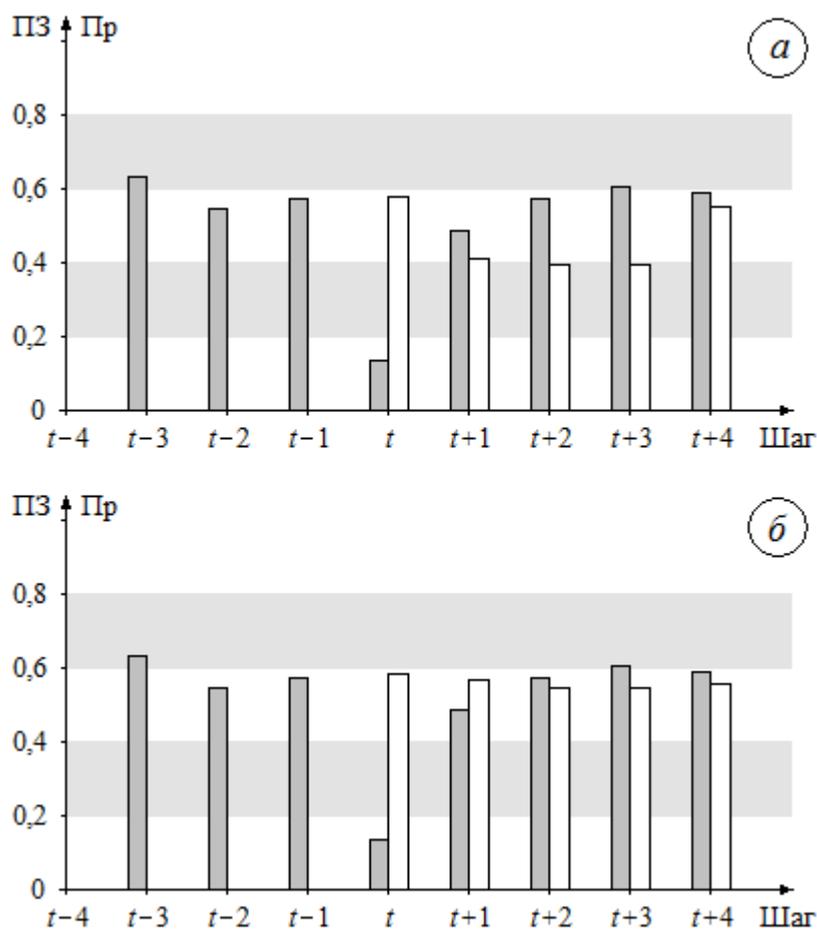
Не трудно убедиться в том, что наличие во временном ряде нетипичных колебаний значений ПЗ может привести к значительным погрешностям в прогнозировании рассматриваемых ПЗ на те шаги, значение прогноза на которые определяется с использованием нетипичных ПЗ. Такой случай иллюстрируется на рис. 6а, где значение ПЗ на шаге  $t$  явно “выпадает” из типичных значений ряда данного показателя здоровья. Поэтому при прогнози-

ровании, например, с помощью алгоритма скользящего среднего по трём последним значениям ПЗ получаем значения прогнозов  $Pr$  для шагов  $t+1$ ,  $t+2$  и  $t+3$ , значительно отличающиеся от полученных на этих шагах фактических значений ПЗ.

Для уменьшения влияния нетипичных значений ПЗ на точность прогнозирования с помощью алгоритмов, основанных на методе наименьших квадратов, можно вводить в этот метод весовые коэффициенты для значений

ряда ПЗ в скользящем окне [10]. Однако при этом подбирать весовые коэффициенты приходится индивидуально для каждого участка временного ряда с нетипичным ПЗ, что является недостатком такого метода.

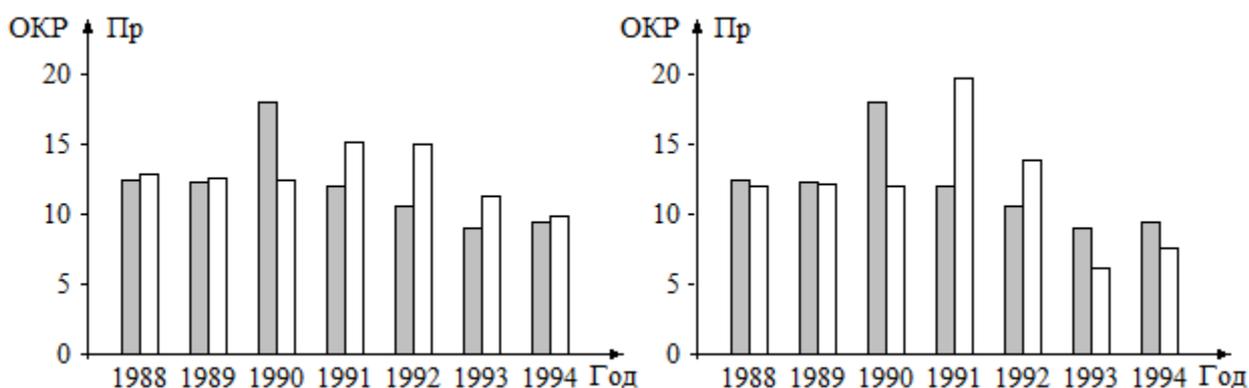
В работах [65, 84 и 86] был предложен простой метод повышения точности прогнозирования временных рядов с нетипичными значениями ПЗ, не требующий сравнения и выбора весов ПЗ в каждой реализации скользящего окна. Сущность этого метода состоит в следующем: *если в скользящее окно попадает ПЗ, значение которого является нетипичным, то при прогнозировании на основе этого скользящего окна вместо значения указанного ПЗ сле-*



**Рис. 6.6.** Прогнозирование при попадании в скользящее окно ПЗ с нетипичным значением с применением: а) обычного алгоритма; б) предлагаемого алгоритма

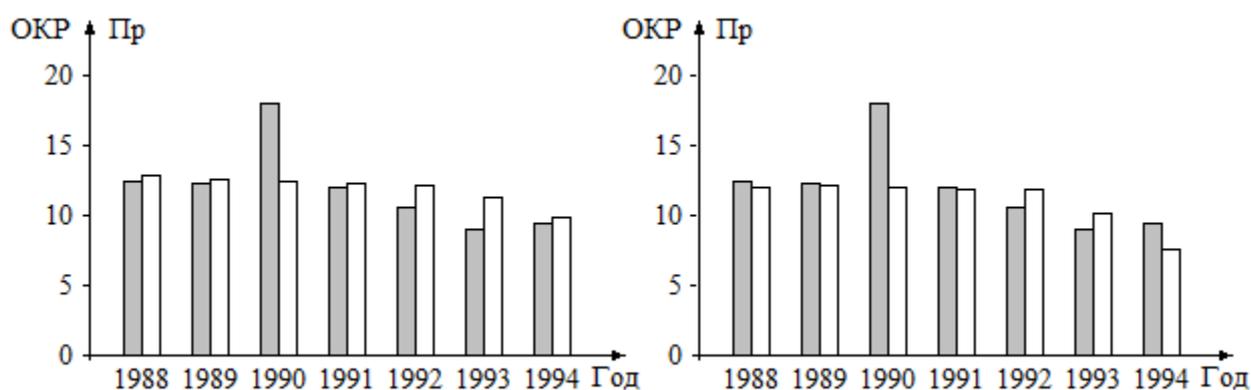
дует использовать его прогноз. Рис. 6.6б поясняет результативность предложенного метода.

На практике значительные отклонения значений ПЗ от типовых встречаются не очень часто. Однако они встречаются. Так на рис. 9.7 приводятся диаграммы участка временного ряда показателя рождаемости и прогнозов этого показателя для населения Шимского района Новгородской области за 1988 ÷ 1994 г. Выделяющееся на диаграммах значение показателя рождаемости за 1990 год (17,90) превысило значения прогнозов на этот год более чем на 45%. Оно существенно повлияло на значения прогнозов на последующие  $m$  лет ( $m$  – ширина скользящего окна в используемой модели). Поэтому при прогнозировании следует учитывать появление нетипичных выбросов в рассматриваемых временных рядах.



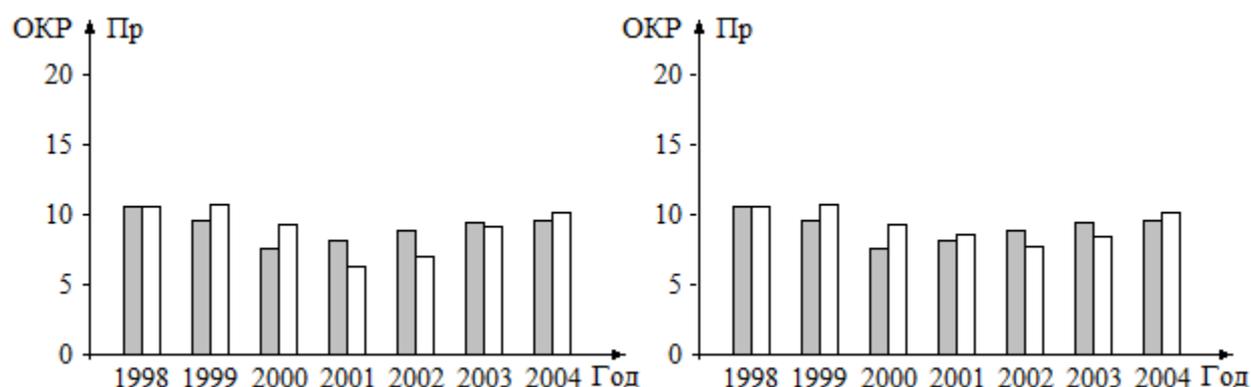
**Рис. 6.7.** Диаграммы участка временного ряда показателя общей рождаемости населения Шимского района Новгородской области и прогнозов этого показателя на один год на основе модели скользящего среднего при  $m = 2$  (слева) и линейной модели при  $m = 3$  без устранения влияния нетипичного значения ОКР в 1990 г.

На рис. 6.8 представлены те же диаграммы, что и на рис. 6.7, но при прогнозировании с использованием предложенного метода повышения точности прогнозирования. Сравнение этих диаграмм показывает, что значения прогнозов на 1991 и 1992 годы на основе модели скользящего среднего и на 1991 ÷ 1993 годы на основе линейной модели стали существенно точнее.



**Рис. 6.8.** Диаграммы участка временного ряда показателя общей рождаемости населения Шимского района Новгородской области и прогнозов этого показателя на один год на основе модели скользящего среднего при  $m = 3$  с помощью обычного метода (слева) и линейной модели при  $m = 3$  с использованием предложенного метода устранения влияния нетипичности значений элементов временного ряда на результаты прогнозирования

Аналогичные выводы можно сделать и по результатам прогнозирования рождаемости на 2001 ÷ 2003 годы для Северо-Западного федерального округа, в котором в 2000 году имело место падение рождаемости на 18,3% по сравнению с прогнозом. Соответствующие диаграммы приведены на рис. 6.9.



**Рис. 6.9.** Диаграммы участка временного ряда показателя общей рождаемости населения Северо-Западного федерального округа на основе линейной модели при  $m = 3$  без использования метода устранения влияния нетипичного значения ОКР в 2000 г. (слева) и с использованием этого метода

Вопрос о том какие значения ПЗ следует считать нетипичными и вводить соответствующие изменения в значения прогнозов зависит от дисперсии значений элементов ряда и, по-видимому, в зависимости от назначения вре-

менного ряда ПЗ должен решаться индивидуально. В дальнейшем во всех случаях принимается, что значения ПЗ считаются нетипичными, если они отличаются от прогнозов на данный шаг более чем на 20% для регионов с населением не менее 0,5 млн человек и более чем на 30% для регионов с населением меньше 0,5 млн человек. Для ИП здоровья характерна меньшая дисперсия, чем для отдельных ПЗ. Поэтому для временных рядов ИП указанный критерий можно соответственно снизить, например, до 20% и до 15%.

Для сравнительного анализа точности прогнозирования на основе разных моделей при наличии нетипичных значений во временных рядах показателей здоровья необходим достаточный объём статистических данных с указанными значениями. Поскольку нетипичные значения во временных рядах показателей здоровья встречаются не часто, то имеющийся статистический материал не позволяет дать заключение по рассматриваемой задаче. Поэтому для её решения использовались модели временных рядов, характеристики которых приведены в главе 3.

Что касается трендов рассматриваемых временных рядов с шагом один год, то согласно анализу для исследования точности алгоритмов прогнозирования при наличии в рядах нетипичных элементов тренды этих рядов можно представить как сумму медленно изменяющегося колебательного процесса и постоянной составляющей. Таким образом, при исследовании использовалась следующая модель временного ряда некоторого показателя  $X$ :

$$X(t) = X_{\text{тр}}(t) + X_{\text{сл}}(t) = A \sin \frac{2\pi t}{T} + B + X_{\text{сл}}(t), \quad (6.3)$$

где  $X_{\text{тр}}(t)$  – трендовая составляющая,  $X_{\text{сл}}(t)$  – случайная составляющая ряда, значения которой имеют колоколообразное распределение в промежутке  $[a, b]$ ,  $A$  – амплитуда синусоиды,  $t$  – текущий год,  $T$  – период синусоиды в годах,  $B$  – постоянная составляющая тренда ( $B > 0$ ). При этом некоторые значения случайной составляющей моделировались как нетипичные. За погрешность прогнозирования принималось отношение модуля разности  $X(t) - \text{Пр}(t)$  к

модулю  $X(t)$  в год прогноза, т.е. использовалась относительная погрешность. Поэтому значения  $A$ ,  $a$  и  $b$  выбирались пропорционально значению  $B$ . В соответствии со статистическими данными было принято:  $A = 0,1B$ ,  $a = -0,1B$ ,  $b = 0,1B$ .

Моделирование случайной составляющей с колоколообразной функцией плотности производилось путём сложения  $n$  взвешенных случайных величин  $C_i$ , равномерно распределённых в  $[a, b]$ , согласно выражению

$$X_{сл} = \frac{1}{n} \sum_{i=1}^n C_i. \quad (6.4)$$

При этом значение  $n$  выбиралось исходя из требования получения необходимого среднего квадратического отклонения  $\sigma(X_{сл})$ :  $n = \text{round}[(b - a)^2 / 12 / \sigma^2(X_{сл})]$ . В соответствии со статистическими данными значение  $\sigma(X_{сл})$  было принято равным  $0,02B$ . Поэтому получили:  $n = 8$ .

Нетипичные значения  $X_{сл}$  воспроизводились со случайным периодом, равномерно распределённым в  $[1, 40]$ . Они моделировались также согласно алгоритму (6.4), но при этом на каждом последнем шаге указанного периода значения  $X_{сл}$  изменялись более чем на  $0,2B$  (коррекция  $X_{сл}$ ) согласно алгоритму:

$$X_{слн} = X_{сл} \pm 0,2B \pm \text{abs}(X_{сл}). \quad (6.5)$$

Алгоритмы (6.3) ÷ (6.5) реализуют обобщённую модель временного ряда показателей здоровья с нетипичными значениями этих показателей, позволяющую получать ряды сколь угодно большой длительности. С помощью данной модели была исследована точность прогнозирования временных рядов показателей здоровья на основе полиномиальных алгоритмов (табл. 9.1) с использованием предложенного метода устранения влияния нетипичности значений элементов ряда на результаты прогнозирования. При этом число появлений нетипичных значений ПЗ для каждого интервала прогнозирования

и при каждом из используемых алгоритмов было выбрано равным 20000. Полученные результаты приведены в табл. 6.5, в которой как и ранее  $m$  – ширина скользящего окна, а  $\overline{\delta\text{Пр}}(T)$  – оценка относительной погрешности прогноза на  $T$  шагов.

Таблица 6.5. Результаты прогнозирования временного ряда показателей здоровья, генерируемого обобщённой моделью

Прогнозирование ряда без нетипичных значений ПЗ на основе алгоритмов скользящего среднего											
$m$	$T$	$\overline{\delta\text{Пр}}(T)$	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	$m$	$T$	$\overline{\delta\text{Пр}}(T)$
1	1	0,0477	2	1	0,0424	3	1	<b>0,0417</b>	4	1	0,0425
	2	0,0497		2	<b>0,0457</b>		2	0,0459		2	0,0471
	3	0,0532		3	<b>0,0504</b>		3	0,0507		3	0,0552
	4	0,0579		4	<b>0,0552</b>		4	0,0559		4	0,0571
Прогнозирование ряда без нетипичных значений ПЗ на основе линейных алгоритмов											
$m$	$T$	$\overline{\delta\text{Пр}}(T)$	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	$m$	$T$	$\overline{\delta\text{Пр}}(T)$
2	1	0,0806	3	1	0,0601	4	1	0,0523	5	1	<b>0,0486</b>
	2	0,1233		2	0,0800		2	0,0651		2	<b>0,0582</b>
	3	0,1682		3	0,1026		3	0,0799		3	<b>0,0700</b>
	4	0,2156		4	0,1268		4	0,0964		4	<b>0,0836</b>
Прогнозирование ряда с нетипичными значениями ПЗ на шагах определения Прогнозов на основе алгоритмов скользящего среднего											
Без использование метода устранения нетипичности значений элементов ряда						С использованием метода устранения нетипичности значений элементов ряда					
$m$	$T$	$\overline{\delta\text{Пр}}(T)$	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	$m$	$T$	$\overline{\delta\text{Пр}}(T)$	$m$	$T$	$\overline{\delta\text{Пр}}(T)$
1	1	0,2328	2	1	0,1177	1	1	0,0494	2	1	0,0460
	2	0,2333		2	0,1185		2	0,0527		2	0,0498
	3	0,2337		3	0,1194		3	0,0575		3	0,0552
	4	0,2336		4	0,1200		4	0,0615		4	0,0594
3	1	<b>0,0817</b>	4	1	0,0913	3	1	<b>0,0451</b>	4	1	0,0459
	2	<b>0,0836</b>		2	0,0918		2	<b>0,0494</b>		2	0,0505

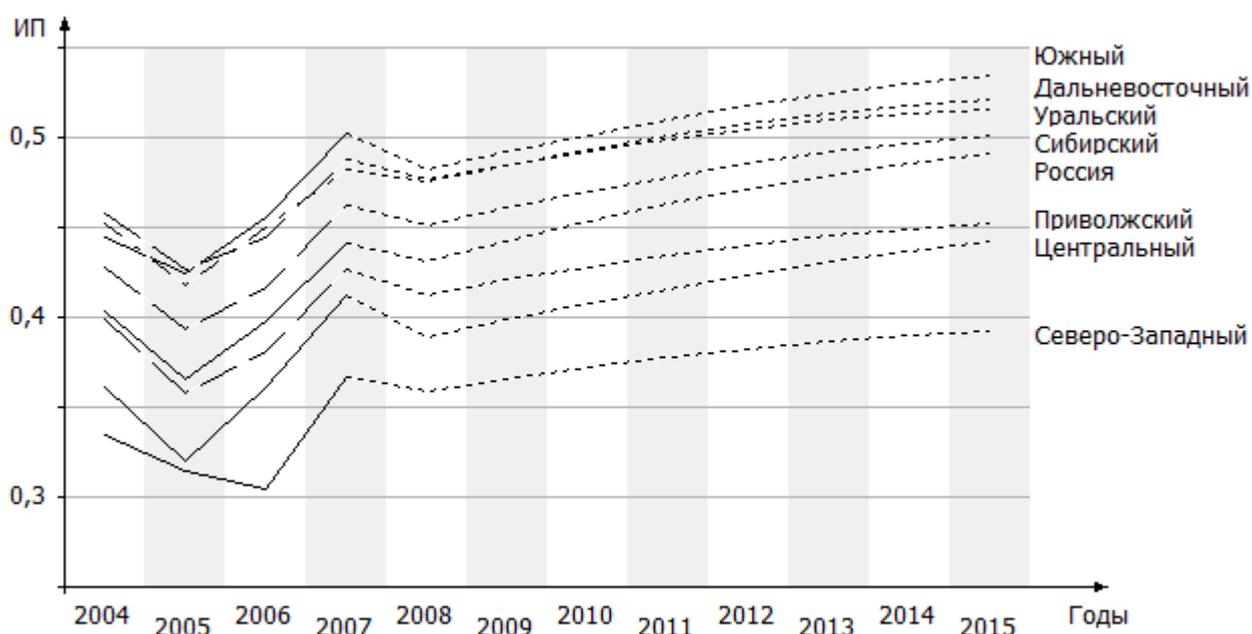
	3	<b>0,0856</b>		3	0,0931		3	<b>0,0553</b>		3	0,0563
	4	<b>0,0877</b>		4	0,0950		4	<b>0,0594</b>		4	0,0603
Прогнозирование ряда с нетипичными значениями ПЗ на шагах определения прогнозов на основе линейных алгоритмов											
Без использование метода устранения нетипичности значений элементов ряда						С использованием метода устранения нетипичности значений элементов ряда					
<i>m</i>	<i>T</i>	$\overline{\delta\Pi p}(T)$	<i>m</i>	<i>T</i>	$\overline{\delta\Pi p}(T)$	<i>m</i>	<i>T</i>	$\overline{\delta\Pi p}(T)$	<i>m</i>	<i>T</i>	$\overline{\delta\Pi p}(T)$
2	1	0,6997	3	1	0,4672	2	1	0,0630	3	1	0,0515
	2	1,0488		2	0,6420		2	0,0908		2	0,0627
	3	1,3981		3	0,8168		3	0,1206		3	0,0753
	4	1,7481		4	0,9925		4	0,1490		4	0,0877
4	1	0,3500	5	1	<b>0,2787</b>	4	1	0,0479	5	1	<b>0,0463</b>
	2	0,4544		2	<b>0,3489</b>		2	0,0551		2	<b>0,0525</b>
	3	0,5589		3	<b>0,4183</b>		3	0,0642		3	<b>0,0610</b>
	4	0,6643		4	<b>0,4884</b>		4	0,0730		4	<b>0,0694</b>

Статистические данные, полученные с использованием модели для рядов без нетипичных значений их элементов (верхняя часть табл. 6.5), близки к данным по относительным погрешностям прогнозирования, полученным на реальной статистике показателей здоровья (табл. 6.2 и 6.3). Это свидетельствует о достоверности данных по эффективности предложенного метода повышения точности прогнозирования временных рядов с нетипичными значениями ПЗ, согласно которым относительная погрешность прогнозирования на шаге с указанным значением ПЗ улучшается примерно на порядок (нижняя часть табл. 6.5). Жирным шрифтом по-прежнему выделены минимальные значения погрешностей прогнозирования, т.е. для прогнозирования с помощью алгоритма скользящего среднего предпочтительнее выбирать  $m = 3$ , а при использовании для этой цели линейного алгоритма –  $m = 5$ .

### 6.5 Прогнозирование здоровья населения регионов России

Рассмотренные модели и алгоритмы прогнозирования позволяют оценивать здоровье населения на ближайшие годы на основе имеющейся стати-

стики показателей здоровья. На рис. 6.10 представлены результаты прогнозирования значений ИП здоровья населения на 2008 ÷ 2015 годы (точечные линии) – [51]. Сплошные линии соответствуют фактическим значениям ИП, для определения которых использовалась 2-я модель (табл. 4.1). Учитывая выводы по точности получения средних значений прогнозов, приведённые в табл. 6.2 ÷ 6.4, прогнозирование проводилось с помощью алгоритмов скользящего среднего при ширине скользящего окна равной трём (до 2010 г.) и линейного при ширине скользящего окна равной пяти (начиная с 2010 г.).



**Рис. 6.10.** Прогнозы динамики интегрального показателя здоровья населения России в целом и федеральных округов

В табл. 6.6 приводятся значения прогнозов как по основным показателям здоровья, так и по интегральному показателю здоровья населения Новгородской области, вычисляемому на основе 2-й модели на 2008 ÷ 2015 годы. При этом во всех алгоритмах использовался предложенный метод повышения точности прогнозирования временных рядов с нетипичными значениями ПЗ, т.е. при отсутствии таких значений ПЗ все алгоритмы работали

обычным образом. Замену ПЗ на его прогноз для расчёта последующих прогнозов программа выполнила только для показателя первичной инвалидности за 2005 г., в котором при прогнозе 14,7 инвалида на 1000 человек населения фактическое значение этого показателя оказалось равным 19,9. Для сравнения в таблице приведены и фактические значения ПЗ за 2007 г.

Для прогнозирования показателей общей заболеваемости по обращениям и первичной инвалидности использовался только алгоритм скользящего средневзвешенного, согласно которому сложение показателей в окне выполняется с весами, а для остальных показателей прогнозирование проводилось так же как и рис. 6.10.

Т а б л и ц а 6.6. Значения прогнозов показателей здоровья населения Новгородской области

Год	2007	2008	2009	2010	2011	2012	2013	2014	2015
ОКР	10,70	11,42	12,00	12,46	12,83	13,12	13,36	13,55	13,70
СППЖ	62,70	62,72	62,73	62,75	62,76	62,77	62,78	62,79	62,81
ОЗО	1999,2	1986,7	1978,1	1970,3	1962,7	1955,0	1947,0	1938,6	1929,7
ПИНВ	15,20	15,04	14,91	14,81	14,73	14,66	14,61	14,57	14,53
ОКС	20,10	19,11	18,46	18,03	17,74	17,55	17,43	17,35	17,29
ИП	0,275	0,289	0,300	0,308	0,315	0,320	0,324	0,328	0,331

Согласно результатам прогнозирования, после многолетнего уменьшения интегрального и других показателей здоровья населения, обусловленного сложными социально-экономическими условиями жизни в перестроечные и послеперестроечные годы, можно ожидать медленного увеличения этих показателей.

## ГЛАВА 7. МОДЕЛИРОВАНИЕ НА ОСНОВЕ ЦЕПЕЙ МАРКОВА

### 7.1. Основные понятия цепей Маркова

В разделе 2.1 рассмотрена последовательность испытаний по схеме Бернулли, т.е. однотипных и независимых относительно некоторого случайного события  $A$ . В указанной схеме всего два исхода испытания: происходит событие  $A$  или происходит противоположное событие  $\bar{A}$ .

Во многих реальных задачах схема испытаний оказывается существенно сложнее: количество возможных исходов, вообще говоря, больше двух (может быть даже бесконечное число исходов), и вероятности появления некоторого исхода в каждом испытании зависят от исходов предыдущих испытаний. Более того, вероятности одного и того же исхода в разных испытаниях могут различаться.

Полагаем, что множество возможных исходов испытаний конечно или счётно. Обозначим это множество возможных несовместных событий (исходов испытания) как множество  $\{E_j\}$ , т.е.  $E_1, E_2, \dots$ . (Заметим, что нумерация событий может быть произвольной, например,  $E_0, E_1, E_2, \dots$ .) В череде испытаний указанные события располагаются случайным образом и образуют некоторую последовательность исходов, рассматриваемую как произведение случайных событий. Для независимых испытаний происходящие события независимы, и вероятность последовательности исходов равна произведению безусловных вероятностей этих исходов, например

$$P(E_j, E_k, E_s) = P(E_j) \cdot P(E_k) \cdot P(E_s) \quad (7.1)$$

При произвольных испытаниях все события - исходы испытаний - оказываются зависимыми, и тогда вероятность произведения событий, согласно теореме умножения, равна произведению, содержащему условные вероятности, в частности (см. разд.2.1):

$$P(E_j, E_k) = P(E_j) P(E_k|E_j);$$

$$P(E_j, E_k, E_s) = P(E_j) P(E_k|E_j) P(E_s|E_j, E_k);$$

$$P(E_j, E_k, E_s, E_m) = P(E_j) P(E_k|E_j) P(E_s|E_j, E_k) P(E_m|E_j, E_k, E_s) \quad (7.2)$$

и т.д.

Как видим, с ростом числа испытаний вычисление условных вероятностей существенно усложняется. В то же время, для большого класса задач зависимость вероятности осуществления события от далекой предыстории пренебрежимо мала. Исходя из этого, модель последовательности зависимых испытаний можно упростить: предположим, что исход каждого последующего испытания зависит лишь от исхода предыдущего испытания и не зависит от всех остальных. Тогда соотношения (7.2) принимают вид:

$$\begin{aligned} P(E_j, E_k) &= P(E_j) P(E_k|E_j); \\ P(E_j, E_k, E_s) &= P(E_j) P(E_k|E_j) P(E_s|E_k); \\ P(E_j, E_k, E_s, E_m) &= P(E_j) P(E_k|E_j) P(E_s|E_k) P(E_m|E_s) \end{aligned} \quad (7.3)$$

и т.д.

Для удобства восприятия обозначим условные вероятности проще:  $P(E_k|E_j) = p_{jk}$ . Таким образом,  $p_{jk}$  – вероятность появления события  $E_k$  в некотором испытании при условии, что предыдущее испытание закончилось исходом  $E_j$ . Так как начальное испытание не имеет предыдущего, то начальному испытанию соответствует безусловная вероятность  $a_j = P(E_j)$  – вероятность осуществления события  $E_j$  в начальном испытании. Следовательно, в общем виде для последовательности  $n$  испытаний с исходами  $E_{j_0}, E_{j_1}, E_{j_2}, \dots, E_{j_n}$  получаем соотношение вероятностей

$$P(E_{j_0}, E_{j_1}, E_{j_2}, \dots, E_{j_n}) = a_{j_0} \cdot p_{j_0 j_1} \cdot p_{j_1 j_2} \cdot \dots \cdot p_{j_{n-1} j_n} \quad (7.4)$$

*Определение. Цепью Маркова* называют последовательность испытаний с возможными исходами  $E_1, E_2, \dots$ , если вероятности любой последовательности исходов  $E_{j_0}, E_{j_1}, E_{j_2}, \dots, E_{j_n}$  определяются по формуле (7.4) через распределение вероятностей  $\{a_j\}$  для исходов  $\{E_j\}$  в начальном испытании и через фиксированные условные вероятности  $\{p_{jk}\}$  в последующих испытаниях.

Терминологически слово «цепь» подразумевает, что каждое последующее звено соединено лишь с предыдущим, т.е. зависит только от него. Свойство независимости от предыстории называют марковским свойством.

Отметим, что поскольку начальное испытание обязательно должно завершиться одним из возможных исходов и все события – слагаемые  $E_j, E_k$  попарно несовместны, то справедливо условие нормировки:

$$\sum_j a_j = 1. \quad (7.5)$$

Также отметим, что после появления некоторого фиксированного исхода  $E_j$  в конкретном испытании следующее испытание непременно должно закончиться осуществлением какого-то исхода  $E_k$ , т.е. вновь сумма всех возможных исходов – событие достоверное, и для соответствующих вероятностей  $p_{jk}$  при фиксированном  $j$  и изменяющемся  $k$  справедливо условие нормировки:

$$\sum_k p_{jk} = 1, \quad j = 1, 2, \dots \quad (7.6)$$

Для цепей Маркова исторически сложилась поясняемая ниже терминология.

Объект, к которому относится исследование, принято называть *системой*. Обозначим ее  $S$ . Под системой  $S$  может пониматься любой объект, подверженный возможным контролируемым изменениям: техническое устройство, предприятие, поликлиника, отдельный индивидуум, группа людей, объединенных по какому-либо признаку, например, проживающие в данном городе, родившиеся в определенном году, имеющие определенное заболевание и т.д.

Возможные исходы испытаний  $E_j$  называют *состояниями системы*. Все возможные состояния системы, случайные события, должны быть несовместны, т.е. в каждый момент времени система может находиться только в одном из рассмотренных состояний. Полагаем, что испытания происходят через равные промежутки времени. Номер испытания называют *номером ша-*

га, следовательно, номер шага служит временным параметром. Таким образом, вместо «в  $(n - 1)$ -м испытании был исход  $E_j$ , а  $n$ -е испытание закончилось появлением  $E_k$ » говорят «на  $n$ -м шаге система из состояния  $E_j$  перешла в состояние  $E_k$ ». Условные вероятности  $p_{jk}$  называют *переходными вероятностями* из состояния  $E_j$  в состояние  $E_k$ .

Множество переходных вероятностей  $\{p_{jk}\}$  можно упорядочить в виде матрицы  $\mathbf{P}$ , называемой *стохастической матрицей*:

$$\mathbf{P} = \begin{pmatrix} p_{11} & p_{12} & p_{13} & \cdots \\ p_{21} & p_{22} & p_{23} & \cdots \\ p_{31} & p_{32} & p_{33} & \cdots \\ \dots\dots\dots\dots\dots\dots \end{pmatrix} \quad (7.7)$$

Согласно (7.6) в стохастической матрице сумма элементов в каждой строке равна 1.  $\mathbf{P}$  – квадратная матрица, конечная или бесконечная в зависимости от количества возможных состояний системы.

Цепь Маркова, в которой переходные вероятности  $p_{jm}$  не зависят от номера шага,  $p_{jm}(k) = p_{jm}(l) = p_{jm}$ , называют *однородной*. В противном случае, когда переходные вероятности меняются в зависимости от номера шага, т.е.  $p_{jm}(k) \neq p_{jm}(l)$ , цепь Маркова называют *неоднородной*.

**Пример 7.1.** Пусть однородная цепь Маркова включает лишь два возможных состояния:  $E_1$  и  $E_2$ . Стохастическая матрица  $\mathbf{P}$  имеет вид

$$\mathbf{P} = \begin{pmatrix} 1-\alpha & \alpha \\ \beta & 1-\beta \end{pmatrix}.$$

Заметим, что для наглядности переходов из состояния в состояние матрицу  $\mathbf{P}$  можно записать совместно с состояниями:

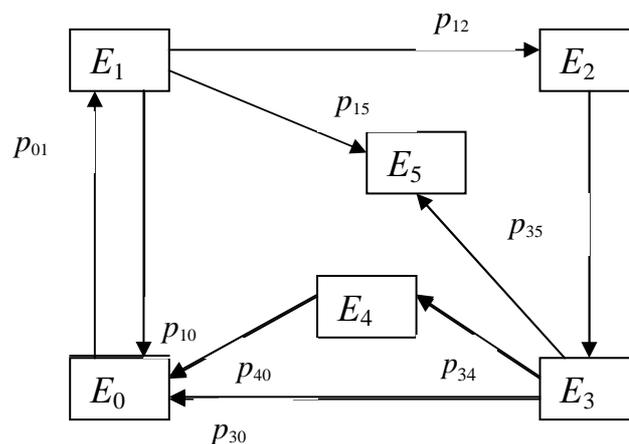
$$\mathbf{P} = \begin{matrix} & E_1 & E_2 \\ \begin{matrix} E_1 \\ E_2 \end{matrix} & \begin{pmatrix} 1-\alpha & \alpha \\ \beta & 1-\beta \end{pmatrix} \end{matrix}.$$

Данную цепь Маркова интерпретируем следующим образом. В период реабилитации ежедневно в определенное время фиксируется состояние больного: улучшение (состояние  $E_1$ ) и ухудшение (состояние  $E_2$ ). При этом, если в предыдущий день у больного наблюдалось улучшение, то в текущий день с вероятностью  $\alpha$

будет ухудшение (т.е.  $\alpha = p_{12}$ ), а с вероятностью  $1-\alpha$  вновь будет улучшение (т.е.  $1-\alpha = p_{11}$ ). Аналогично, если в предыдущий день у больного наблюдалось ухудшение, то улучшение наступит с вероятностью  $\beta$  (т.е.  $p_{21}=\beta$ ), а продолжится ухудшение с вероятностью  $1-\beta$  (т.е.  $1-\beta = p_{22}$ ).

Для анализа цепи Маркова удобно использовать геометрическую иллюстрацию процесса, называемую *графом состояний*. Граф состояний  $G(S)$  наглядно представляет состояния системы  $S$ , возможные переходы из состояния в состояние и вероятности, соответствующие этим переходам. На графе состояния принято изображать в виде прямоугольника, возможные переходы из одного состояния в другое соответствующими стрелками с указанием переходных вероятностей. Вероятность **того**, что система останется в том же состоянии, обычно на графе не указывается. Значение этой вероятности легко рассчитывается, исходя из нормирующего условия (7.6), справедливого для каждого состояния.

**Пример 7.2.** Рассмотрим элементарный граф состояний человека относительно его здоровья (рис. 7.1).



**Рис. 7.1.** Элементарный граф состояний человека.

На графе представлены шесть возможных состояний человека:

$E_0$  – здоров (чувствует себя здоровым);

$E_1$  – чувствует себя заболевшим, пытается вылечиться самостоятельно;

$E_2$  – обратился за медицинской помощью, проводится обследование;

$E_3$  – диагноз установлен, проводится курс лечения;

$E_4$  – реабилитация;

$E_5$  – летальный исход.

Стрелки на графе указаны лишь для ненулевых переходных вероятностей, меняющих состояние. Вероятности остаться в том же состоянии вычисляются. Например, переходная вероятность  $p_{11}$ , означающая что человек, находясь в состоянии  $E_1$ , на следующем шаге также будет в состоянии  $E_1$ , на графе не указана. Но согласно (7.6)  $p_{10} + p_{11} + p_{12} + p_{13} + p_{14} + p_{15} = 1$ , т.е. в рассматриваемом случае  $p_{10} + p_{11} + p_{12} + 0 + 0 + p_{15} = 1$ . Поэтому  $p_{11} = 1 - (p_{10} + p_{12} + p_{15})$ .

Аналогично:  $p_{22} = 1 - p_{23}$ ;  $p_{33} = 1 - (p_{30} + p_{34} + p_{35})$ .

Отметим, что из состояния  $E_5$  на графе не выходит ни одной стрелки. Следовательно, попав в него однажды, система (человек) остается в этом состоянии навсегда. Такие состояния принято называть *поглощающими*.

По значениям вероятностей на графе однозначно формируется стохастическая матрица. Представленная на рис. 7.1 схема является некой упрощенной, элементарной моделью реальной ситуации.

Пусть система имеет  $n$  возможных состояний  $E_1, E_2, \dots, E_n$ , образующих однородную цепь Маркова. Таким образом, на каждом шаге матрица переходных вероятностей (7.7) одна и та же: стохастическая матрица  $\mathbf{P}$ , элементы которой известны. Вектор начальных вероятностей  $\{a_j\}$  также полагаем заданным. Рассмотрим алгоритм нахождения вероятностей состояний системы на 1-м, 2-м, ...  $k$ -м шаге. Обозначим безусловную вероятность нахождения системы в состоянии  $E_j$  на  $k$ -м шаге, как  $p_j(k)$ . Тогда для всех возможных состояний  $E_1, E_2, \dots, E_n$  на каждом шаге  $k$  получаем свой набор (вектор) соответствующих вероятностей

$$p_1(k), p_2(k), \dots, p_n(k), \quad (7.8)$$

где  $k = 1, 2, \dots$



$$\{p_j(1)\} = \{p_j(0)\} \cdot \mathbf{P}. \quad (7.11)$$

Далее аналогично, полагая вероятности  $p_1(1), p_2(1), \dots, p_n(1)$  в качестве полной системы гипотез, можно найти соответствующие вероятности на втором шаге  $p_1(2), p_2(2), \dots, p_n(2)$  по формуле:

$$p_j(2) = \sum_{i=1}^n p_i(1) p_{ij}.$$

Для вычисления всей совокупности безусловных вероятностей на втором шаге вновь запишем в матричном виде

$$\{p_j(2)\} = \{p_j(1)\} \cdot \mathbf{P} = \{p_j(0)\} \cdot \mathbf{P}^2, \quad (7.12)$$

где  $\mathbf{P}^2$  - квадрат стохастической матрицы.

Продолжая далее, в общем случае на  $k$ -м шаге получаем

$$\{p_j(k)\} = \{p_j(0)\} \cdot \mathbf{P}^k. \quad (7.13)$$

В случае неоднородной марковской цепи, когда стохастическая матрица изменяется на каждом шаге, соотношение вида (7.13) для вектора безусловных вероятностей на  $k$ -м шаге имеет вид

$$\{p_j(k)\} = \{p_j(0)\} \mathbf{P}(1) \cdot \mathbf{P}(2) \cdot \dots \cdot \mathbf{P}(k), \quad (7.14)$$

где  $\mathbf{P}(m)$  – стохастическая матрица системы для  $m$ -го шага.

**Пример 7.3.** Исследуется успеваемость в группе студентов. Исходя из среднего балла оценок в сессию, сформированы четыре возможных состояния: «двоечник» (студент отчислен за неуспеваемость), «троечник», «хорошист», «отличник». В соответствии с вкладываемым в указанное понятие смыслом удобно данные состояния обозначить, как  $E_2, E_3, E_4, E_5$ . Согласно статистическим данным конкретной специальности стохастическая матрица  $\mathbf{P}$  имеет вид (запишем матрицу вместе с состояниями):

$$\mathbf{P} = \begin{matrix} & E_2 & E_3 & E_4 & E_5 \\ \begin{matrix} E_2 \\ E_3 \\ E_4 \\ E_5 \end{matrix} & \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0,13 & 0,70 & 0,16 & 0,01 \\ 0,08 & 0,26 & 0,56 & 0,10 \\ 0,02 & 0,05 & 0,53 & 0,40 \end{pmatrix} \end{matrix}$$

Шаг процесса – 1 семестр. Все время исследования – 9 семестров. Пусть в начальном состоянии, после сессии 1 семестра, студент находится, например, в состоянии  $E_4$ . Найдем вектора безусловных вероятностей после 1-го и 2-го шагов, т.е. осуществим вероятностный прогноз на два временных интервала: после 2-го и 3-го семестров. После первого семестра (начальный вектор):  $\{p_j(0)\} = (0; 0; 1; 0)$ . После второго семестра (см. формулу (7.11)):

$$\begin{aligned} \{p_j(1)\} &= \{p_j(0)\} \cdot \mathbf{P} = (0; 0; 1; 0) \cdot \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0,13 & 0,70 & 0,16 & 0,01 \\ 0,08 & 0,26 & 0,56 & 0,10 \\ 0,02 & 0,05 & 0,53 & 0,40 \end{pmatrix} = \\ &= (0,08; 0,26; 0,56; 0,10). \end{aligned}$$

После третьего семестра (см. формулу (7.12)):

$$\begin{aligned} \{p_j(2)\} &= \{p_j(1)\} \cdot \mathbf{P} = (0,08; 0,26; 0,56; 0,10) \cdot \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0,13 & 0,70 & 0,16 & 0,01 \\ 0,08 & 0,26 & 0,56 & 0,10 \\ 0,02 & 0,05 & 0,53 & 0,40 \end{pmatrix} = \\ &= (0,1606; 0,3326; 0,4082; 0,0986). \end{aligned}$$

Таким образом, «хорошист» в первом семестре после сессии третьего семестра с вероятностью 0,1606 окажется отчисленным, с вероятностью 0,3326 будет «троечником», с вероятностью 0,4082 – «хорошистом» и с вероятностью 0,0986 – «отличником».

Аналогично вычисление безусловных вероятностей можно продолжить и на последующие шаги.

Далее проведем соответствующий расчет не для индивидуума с заданным начальным состоянием, а для целой группы. Пусть известно, что после первого семестра среди студентов в учебной группе 30% «троечников», 50% «хорошистов» и 20% «отличников». Вновь осуществим вероятностный прогноз на два семестра.

После первого семестра (начальный вектор)  $\{p_j(0)\} = (0; 0,3; 0,5; 0,2)$ .

После второго семестра:

$$\begin{aligned} \{p_j(1)\} &= \{p_j(0)\} \cdot \mathbf{P} = (0; 0,3; 0,5; 0,2) \cdot \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0,13 & 0,70 & 0,16 & 0,01 \\ 0,08 & 0,26 & 0,56 & 0,10 \\ 0,02 & 0,05 & 0,53 & 0,40 \end{pmatrix} = \\ &= (0,083; 0,350; 0,434; 0,133). \end{aligned}$$

После третьего семестра:

$$\begin{aligned} \{p_j(2)\} &= \{p_j(1)\} \cdot \mathbf{P} = (0,083; 0,350; 0,434; 0,133) \cdot \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0,13 & 0,70 & 0,16 & 0,01 \\ 0,08 & 0,26 & 0,56 & 0,10 \\ 0,02 & 0,05 & 0,53 & 0,40 \end{pmatrix} = \\ &= (0,16588; 0,36449; 0,36953; 0,1001). \end{aligned}$$

Округляя полученные вероятности до третьего знака после запятой, запишем прогнозируемое соотношение студентов после третьего семестра: отчисленных – 16,6%, «троечников» – 36,4%, «хорошистов» – 37%, «отличников» – 10%.

Заметим, что в рассмотренном примере 7.3 фигурирует однородная цепь Маркова. В реальной ситуации представляется более разумным обратиться к неоднородной цепи, так как ввиду различной сложности семестров и «взросления» студентов с течением времени переходные вероятности, вообще говоря, оказываются зависящими от номера шага (семестра).

Неоднородная цепь Маркова в задаче исследования общественного и индивидуального здоровья рассмотрена далее в разделе 7.3.

## 7.2. Некоторые модели марковских процессов

Во многих реальных задачах система переходит из состояния в состояние не в фиксированные, а в случайные моменты времени  $t$ . В этом случае время оказывается не дискретной величиной, как в цепи Маркова, а непрерывной переменной. Такую цепь с непрерывным временем принято называть *марковским случайным процессом*. Характерным свойством марковского процесса, в отличие от случайных процессов других типов, является зависимость переходов из состояния в состояние только от ближайшей предыстории: «будущее состояние системы зависит лишь от настоящего и не зависит от прошлого».

Для марковских процессов вместо переходных вероятностей из состояния  $E_j$  в состояние  $E_k$  за фиксированное время  $\Delta t$ , которые в разделе 7.1 обозначены как  $p_{jk}$ , введем понятие *интенсивности переходов*  $\lambda_{jk}$ :

$$\lambda_{jk}(t) = \lim_{\Delta t \rightarrow 0} \frac{p_{jk}}{\Delta t} = \frac{dp_{jk}}{dt}.$$

В данном случае вероятность  $p_{jk}$  зависит от  $\Delta t$ , т.е. является функцией времени, а интенсивность  $\lambda_{jk}(t)$  оказывается не чем иным как плотностью этой вероятности перехода из  $E_j$  в  $E_k$ , определяемой для каждого момента времени  $t$ .

В конкретных моделях (в частности, при исследовании здоровья населения) для малых промежутков времени  $\Delta t$  принято функцию  $\lambda_{jk}(t)$  считать постоянной в течение этого времени, что, как правило, соответствует реальности. В качестве единицы измерения времени используются 1 час, 1 сутки, 1 месяц, 1 год и т.д. Случайный процесс с постоянными интенсивностями переходов  $\{\lambda_{jk}\}$  называют *однородным*, в противном случае, когда  $\lambda_{jk}$  является функцией  $t$  – *неоднородным*.

Для вычисления интенсивности  $\lambda_{jk}$  по статистическим данным можно использовать приближенные соотношения

$$\lambda_{jk} = \frac{p_{jk}}{\Delta t} \quad (\text{т.е. } p_{jk} = \lambda_{jk} \cdot \Delta t). \quad (7.15)$$

Следовательно, согласно статистическому определению вероятности

$$\lambda_{jk} = \frac{n_k}{n_j \Delta t}, \quad (7.16)$$

где  $n_j$  - количество элементов системы, находящихся в начальный момент времени в состоянии  $E_j$ , а  $n_k$  - количество элементов (из них) перешедших в состояние  $E_k$  за время  $\Delta t$ . Если переход из  $E_j$  в  $E_k$  невозможен, то интенсивность  $\lambda_{jk}$  также как и вероятность  $p_{jk}$ , оказывается равной нулю.

Для марковских процессов при построении графа состояний вместо переходных вероятностей  $p_{jk}$  принято указывать интенсивность переходов  $\lambda_{jk}$ . Интенсивности тривиальных переходов из  $E_j$  в  $E_j$  на графе не указываются.

Важнейшей составляющей исследования марковских процессов является задача нахождения безусловных вероятностей для любого момента времени  $t$ . Соответствующий вектор безусловных вероятностей в момент времени  $t$  для состояний  $E_0, E_1, \dots, E_n$  обозначим как

$$p_0(t), p_1(t), \dots, p_n(t), \quad (7.17)$$

причем для каждого  $t$  имеет место условие нормировки

$$\sum_j p_j(t) = 1 \quad (7.18)$$

Тогда начальный вектор, определяемый в момент времени  $t = 0$ , оказывается

$$p_0(0), p_1(0), \dots, p_n(0), \quad \text{где } \sum_j p_j(0) = 1 \quad (7.19)$$

В отличие от цепи Маркова, (т.е. марковского процесса с дискретным временем) вычисление вектора безусловных вероятностей (7.17) при началь-

ном векторе (7.19) для марковского случайного процесса осуществляется не матричным методом, а путем решения соответствующей системы линейных дифференциальных уравнений, называемых *системой Колмогорова*.

Решение систем дифференциальных уравнений – это задача достаточно не тривиальная, трудоемкая и при большом количестве уравнений и произвольных коэффициентах  $\lambda_{jk}(t)$  практически невыполнимая вручную. Обычно решение системы уравнений находится с помощью соответствующего программного обеспечения.

Задачей специалистов-нематематиков при решении конкретной проблемы является разработка графа состояний и, как следствие, формирование системы дифференциальных уравнений Колмогорова, исходя из реальных данных и структуры выбранной модели. Приведем конкретную методику.

1. Каждое из уравнений характеризует одно из состояний  $E_j$  и является линейным дифференциальным уравнением вида

$$\frac{dp_j}{dt} = -\sum_k \lambda_{jk} p_j + \sum_k \lambda_{kj} p_k, \quad (7.20)$$

где  $\lambda_{jk}$  - интенсивности переходов из состояния  $E_j$  (берутся со знаком «минус»),  $\lambda_{kj}$  - интенсивности переходов в состояние  $E_j$  (берутся со знаком «плюс»), индекс в записи вероятностей  $p_k$  соответствует состоянию, из которого осуществляется переход в состояние  $E_j$ . Заметим, что переходы из  $E_j$  в  $E_j$  в уравнениях игнорируются (т.е. слагаемых  $\lambda_{jj} p_j$  нет).

2. Количество уравнений системы равно количеству возможных состояний (в рассматриваемом случае -  $n+1$ ).

3. Начальными значениями (условиями) при решении системы служат конкретные числа (7.19). Например, если в начальный момент времени система находится в состоянии  $E_1$ , то

$$p_0(0) = 0, \quad p_1(0) = 1, \quad p_2(0) = 0, \quad \dots, \quad p_n(0) = 0.$$

4. Решениями системы (искомыми функциями) являются безусловные вероятности (7.17). Соответствующие значения вероятностей нахождения системы в каждом из состояний можно найти для любого значения времени  $t$ .

Заметим, что условие нормировки (7.18), связывающее искомые функции, может заменить любое из уравнений системы. В этом случае из условия нормировки  $p_0 + p_1 + \dots + p_n = 1$  выражаем выбранную вероятность  $p_i$  через оставшиеся переменные. Подставляем полученное выражение в каждое из дифференциальных уравнений, кроме уравнения, с производной  $\frac{dp_i}{dt}$ , которое отбрасываем. Таким образом, порядок системы уменьшается на единицу, что упрощает ее решение. Зависимая переменная  $p_i$  в этой системе отсутствует. Однако после решения полученной системы  $p_i$  также оказывается определенной.

Задача составления системы дифференциальных уравнений значительно упрощается, если предварительно начертить граф состояний с указанием интенсивностей переходов  $\lambda_{jk}$ . Входящие в вершину графа или выходящие из вершины графа стрелки указывают на знак («плюс» или «минус»), который нужно приписать соответствующей интенсивности в дифференциальном уравнении. Наглядность графа позволяет избежать досадных ошибок, возможных при формальной записи системы уравнений.

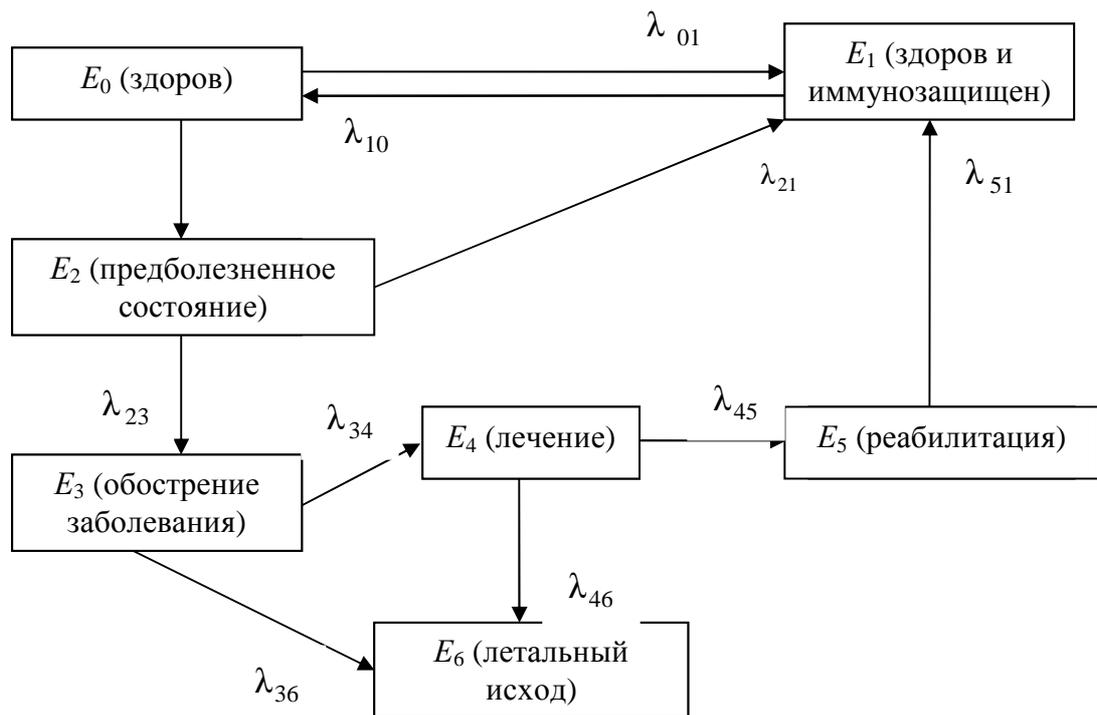
**Пример 7.4.** *Марковский процесс инфекционного заболевания.*

Появление и развитие инфекционного заболевания можно структурировать введением характерных состояний и указанием возможных переходов из состояния в состояние и соответствующих интенсивностей этих переходов. При этом полагаем, что вероятности переходов из состояния  $E_j$  в состояние  $E_k$  не зависят от предыстории, т.е. состояний, предшествовавших  $E_j$ . Это позволяет считать исследуемый процесс марковским и построить соответствующую теоретическую модель.

Возможные состояния и интенсивности переходов, характерные для различных инфекционных заболеваний, представлены на рис 7.2. В данной модели представлены основные состояния при возникновении и протекании инфекционного за-

болевания. Переход из состояния  $E_0$  (здоров) в состояние  $E_1$  (здоров и иммунозащищен) происходит непосредственно через состояние «вакцинация», которое ввиду его кратковременности в модели отсутствует. Указанная обратная связь из  $E_1$  в  $E_0$  (т.е. потеря иммунозащищенности) для некоторых инфекционных заболеваний на небольшом временном интервале может отсутствовать вовсе (тогда  $\lambda_{10}=0$ ).

Попасть из состояния  $E_0$  в состояние  $E_1$  также можно, пройдя последовательно через состояния  $E_2, E_3, E_4, E_5$  (т.е. переболев данным заболеванием). Отметим, что переход из  $E_2$  в  $E_1$  также может быть связан с вакцинацией или легким протеканием заболевания без стадии обострения, что фактически соответствует состоянию  $E_2$ .



**Рис. 7.2.** Граф состояний при инфекционном заболевании.

Попадание в состояние  $E_6$  (летальный исход) в данной модели возможно лишь из двух состояний:  $E_3$  (обострение заболевания) и  $E_4$  (лечение). В этом случае летальный исход происходит именно вследствие инфекционного заболевания. Ра-

зумеется, попадание в  $E_6$  теоретически возможно из любого состояния. Но в этом случае причиной летального исхода, скорее всего, будет какая-то внешняя причина, а не исследуемое заболевание, и вероятности таких переходов за период протекания инфекционного заболевания крайне малы. Поэтому интенсивности таких переходов практически равны нулю и в модели проигнорированы.

Отметим также, что на графе из состояния  $E_6$  нет выходящих стрелок, т.е. это состояние поглощающее.

В случае  $\lambda_{10} = 0$  (т.е. за рассматриваемый промежуток времени иммунозащищенность устойчива) состояние  $E_0$  оказывается источником: из него можно выйти, но за период исследования нельзя вернуться.

Переходы из состояния в состояние могут происходить в произвольные моменты времени, что связано с индивидуальной защищенностью каждого живого организма и типом вируса.

Исходя из возможных переходов и их интенсивностей, указанных на графе, можно найти безусловные вероятности нахождения системы в любом из состояний в любой момент времени. Для этого нужно составить и решить систему дифференциальных уравнений Колмогорова при заданном начальном векторе (7.19). Например, для индивидуума, находящегося в состоянии  $E_0$  (здоров), начальный вектор имеет вид  $(1; 0; 0; 0; 0; 0; 0)$ .

Согласно общей теории в нашем случае требуется составить семь уравнений по числу состояний в системе. Искомыми функциями будут функции аргумента  $t$   $p_0(t), p_1(t), \dots, p_6(t)$ . далее, ради удобства записи, аргумент  $t$  указывать не будем.

Состояние  $E_0$  на графе (рис. 7.2) имеет одну входящую стрелку (из  $E_1$ ) и две выходящих (в  $E_1$  и  $E_2$ ). Следовательно, дифференциальное уравнение, соответствующее  $E_0$ , оказывается таким (см. (7.20)):

$$\frac{dp_0}{dt} = -\lambda_{01} p_0 - \lambda_{02} p_0 + \lambda_{10} p_1.$$

Еще раз подчеркнем, что интенсивности при выходящих стрелках на графе в дифференциальном уравнении берутся со знаком «минус». Разумеется, в правой части полученного уравнения можно привести подобные члены.

Составляя по указанной выше методике уравнение для каждого состояния, приходим к системе уравнений Колмогорова:

$$\left\{ \begin{array}{l} \frac{dp_0}{dt} = -(\lambda_{01} + \lambda_{02})p_0 + \lambda_{10}p_1, \\ \frac{dp_1}{dt} = \lambda_{01}p_0 - \lambda_{10}p_1 + \lambda_{21}p_2 + \lambda_{51}p_5, \\ \frac{dp_2}{dt} = \lambda_{02}p_0 - (\lambda_{21} + \lambda_{23})p_2, \\ \frac{dp_3}{dt} = \lambda_{23}p_2 - (\lambda_{34} + \lambda_{36})p_3, \\ \frac{dp_4}{dt} = \lambda_{34}p_3 - (\lambda_{45} + \lambda_{46})p_4, \\ \frac{dp_5}{dt} = \lambda_{45}p_4 - \lambda_{51}p_5, \\ \frac{dp_6}{dt} = \lambda_{36}p_3 + \lambda_{46}p_4. \end{array} \right. \quad (7.21)$$

Решая найденную систему, получаем её общее решение, которое в данном случае состоит из семи длинных громоздких равенств с произвольными постоянными. Причем при больших  $n$  и особенно переменных коэффициентах  $\lambda_{ij}$  могут возникнуть проблемы с нахождением общего решения, т.е. эта задача достаточно сложная даже при компьютерной реализации. Поэтому, как отмечено выше, можно понизить порядок системы на 1, поскольку одна из искомым функций (любая) линейно выражается через остальные. Удалим из системы, например, второе уравнение (уравнение с производной  $dp_1/dt$ ). Тогда в оставшихся уравнениях системы необходимо заменить соответствующую переменную  $p_1$  ее выражением из условия нормировки:  $1 - p_0 - p_2 - p_3 - p_4 - p_5 - p_6$ . В нашем примере  $p_1$  имеется лишь в первом уравнении. Поэтому получаем систему дифференциальных уравнений, порядка на 1 меньше, чем порядок системы (7.21) и одно алгебраическое соотношение:

$$\left\{ \begin{array}{l} \frac{dp_0}{dt} = -(\lambda_{01} + \lambda_{02})p_0 + \lambda_{10}(1 - p_0 - p_2 - p_3 - p_4 - p_5 - p_6), \\ \frac{dp_2}{dt} = \lambda_{02}p_0 - (\lambda_{21} + \lambda_{23})p_2, \\ \frac{dp_3}{dt} = \lambda_{23}p_2 - (\lambda_{34} + \lambda_{36})p_3, \\ \frac{dp_4}{dt} = \lambda_{34}p_3 - (\lambda_{45} + \lambda_{46})p_4, \\ \frac{dp_5}{dt} = \lambda_{45}p_4 - \lambda_{51}p_5, \\ \frac{dp_6}{dt} = \lambda_{36}p_3 + \lambda_{46}p_4, \\ p_1 = 1 - p_0 - p_2 - p_3 - p_4 - p_5 - p_6. \end{array} \right. \quad (7.22)$$

Разумеется, правую часть первого из уравнений этой системы можно преобразовать (выполнить умножение и перегруппировать слагаемые). Аналогичным образом, можно было изъять из системы (7.21) другую переменную. В частности, удобно избавиться в системе от переменной  $p_6$ , поскольку она явно отсутствует во всех уравнениях, кроме отбрасываемого.

Практически системы рассматриваемого вида решают при заданных начальных условиях (задача Коши), что позволяет определить в общем решении конкретные значения произвольных постоянных. Полученное таким образом частное решение системы (7.21) при заданных начальных условиях представляет собой систему соотношений для координат вектора безусловных вероятностей  $\{p_j(t)\}$  при любых  $t > 0$ . Полагая в частном решении конкретное значение времени  $t$ , получим набор значений – вероятностей состояний рассматриваемого процесса развития и лечения инфекционного заболевания именно в данный момент времени  $t$ . Однако к тому же результату можно придти вычислительно проще: не находя ни общего, ни частного решения, а решая одним из приближенных методов непосредственно систему (7.21) при заданных начальных условиях и фиксированном  $t$ .

Найдем, например, координаты вектора  $\{p_j(4)\}$  для некоторого индивидуума (или совокупности индивидуумов), классифицированного в состоянии  $E_0$  (здоров), т.е. решаем систему при начальном условии  $p_0(0) = 1, p_1(0) = 0, \dots, p_6(0) = 0$ . **Предположим**, что интенсивности переходов за 1 неделю  $\lambda_{ij}$  имеют значения:

$$\lambda_{01} = 0,05; \lambda_{02} = 0,03; \lambda_{10} = 0,005; \lambda_{21} = 0,005; \lambda_{23} = 0,05;$$

$$\lambda_{34} = 0,09; \lambda_{36} = 0,001; \lambda_{45} = 0,08; \lambda_{46} = 0,001; \lambda_{51} = 0,08.$$

Тогда, используя систему (7.21) или систему (7.22), находим:

$$p_0(4) = 0,727765; \quad p_1(4) = 0,170483; \quad p_2(4) = 0,091708; \quad p_3(4) = 0,008885;$$

$$p_4(4) = 0,001060; \quad p_5(4) = 0,000085; \quad p_6(4) = 0,000014.$$

Для совокупности индивидуумов, находящихся в начальный момент времени в состоянии  $E_0$  (здоров), найденные вероятности представляют четырехнедельный прогноз количественного разбиения на группы по всем состояниям здоровья. Полученные соотношения легко перевести в проценты (или в промилле): 72,78% будут в состоянии  $E_0$ , 17,05% будут в состоянии  $E_1$  и т.д.

*Замечание.* В практической реализации примера 7.4 схема развития инфекционного заболевания несколько упрощена. В реальности интенсивности переходов, особенно в стадии заражения, не являются постоянными, а зависят от ряда причин и, в частности, от времени  $t$ . Полагая в качестве интенсивностей  $\lambda_{jk}$  определенного рода функции, (в зависимости от модели, процесса заражения, активности вируса, учета контактов заболевших и т.д.), мы усложняем лишь решение системы дифференциальных уравнений. Сама структура системы уравнений Колмогорова и методика их построения остаются неизменными.

Если в качестве системы в марковском случайном процессе рассматривается группа  $N$  индивидуумов, то наряду с безусловными вероятностями можно непосредственно найти их количественное распределение по состояниям в произвольный момент времени  $t$ :  $(m_0, m_1, \dots, m_n)$ , где  $m_i$  - количество индивидуумов в состоянии  $E_i$ ,  $\sum m_i = N$ . Поскольку вероятности представляют собой доли количества объектов, обладающих определенным признаком,  $p_i = \frac{m_i}{N}$ , то распределение  $N$  индивидуумов по состояниям получается в результате умножения вероятностей на  $N$ :  $m_i = p_i N$ ,  $i = \overline{0, n}$ .

Линейный характер зависимости  $m_i$  от  $p_i$  позволяет по тем же правилам составить систему дифференциальных уравнений Колмогорова с искомыми функциями  $m_0, m_1, \dots, m_n$ :



Согласно разработанной методике всё население классифицируется по возрастам  $t_i$  и состояниям здоровья  $E_j$ , при этом  $E_0$ - «относительно здоров» и  $E_n$ - «смерть» (состояние  $E_n$  обычно подразделяется на несколько состояний согласно смертности по причинам).  $E_1, E_2, \dots, E_{n-1}$  - состояния, соответствующие наиболее тяжелому типу заболеваемости. Такое определение состояний (по наиболее тяжелому, доминирующему заболеванию) призвано привести в соответствие реальность, когда каждый индивидуум обладает целым «букетом» различных заболеваний, и математическую модель, согласно которой в любой момент времени индивидуум может находиться лишь в одном из состояний. Отметим, что состояние можно определить и как совокупность нескольких распространенных заболеваний, однако для реальных данных это приведет к неоправданно резкому увеличению числа состояний и негативно повлияет на точность статистических выводов. Определив однозначно состояния, человеческую жизнь можно интерпретировать, как последовательность (цепь) блужданий по состояниям здоровья, включая предельное состояние «смерть».

Выбор и формирование состояний системы – важнейший этап исследования, предопределяющий ценность полученных впоследствии результатов. Например, для состояний  $E_j$  можно ввести [147, 148] следующую классификацию (см. табл. 7.1).

Таблица 7.1. Классификация состояний

Состояние системы	Класс
$E_0$	«ОТНОСИТЕЛЬНО ЗДОРОВЫЕ»
$E_1$	I00-I99, БОЛЕЗНИ СИСТЕМЫ КРОВООБРАЩЕНИЯ
$E_2$	S00-T98, ТРАВМЫ, ОТРАВЛЕНИЯ И НЕКОТОРЫЕ ДРУГИЕ ПОСЛЕДСТВИЯ ВОЗДЕЙСТВИЯ ВНЕШНИХ ПРИЧИН
$E_3$	C00-D48, НОВООБРАЗОВАНИЯ
$E_4$	J00-J99, БОЛЕЗНИ ОРГАНОВ ДЫХАНИЯ

<i>E</i> <sub>5</sub>	R00-R99, СИМПТОМЫ, ПРИЗНАКИ И ОТКЛОНЕНИЯ ОТ НОРМЫ, ВЫЯВЛЕННЫЕ ПРИ КЛИНИЧЕСКИХ И ЛАБОРАТОРНЫХ ИССЛЕДОВАНИЯХ, НЕ КЛАССИФИЦИРОВАННЫЕ В ДРУГИХ РУБРИКАХ
<i>E</i> <sub>6</sub>	K00-K93, БОЛЕЗНИ ОРГАНОВ ПИЩЕВАРЕНИЯ
<i>E</i> <sub>7</sub>	A00-B99, НЕКОТОРЫЕ ИНФЕКЦИОННЫЕ И ПАРАЗИТАРНЫЕ БОЛЕЗНИ
<i>E</i> <sub>8</sub>	G00-G99, БОЛЕЗНИ НЕРВНОЙ СИСТЕМЫ
<i>E</i> <sub>9</sub>	N00-N99, БОЛЕЗНИ МОЧЕПОЛОВОЙ СИСТЕМЫ
<i>E</i> <sub>10</sub>	F00-F99, ПСИХИЧЕСКИЕ РАССТРОЙСТВА И РАССТРОЙСТВА ПОВЕДЕНИЯ
<i>E</i> <sub>11</sub>	D50-D89, БОЛЕЗНИ КРОВИ, КРОВЕТВОРНЫХ ОРГАНОВ И ОТДЕЛЬНЫЕ НАРУШЕНИЯ, ВОВЛЕКАЮЩИЕ ИММУННЫЙ МЕХАНИЗМ
<i>E</i> <sub>12</sub>	E00-E90, БОЛЕЗНИ ЭНДОКРИННОЙ СИСТЕМЫ, РАССТРОЙСТВА ПИТАНИЯ И НАРУШЕНИЯ ОБМЕНА ВЕЩЕСТВ
<i>E</i> <sub>13</sub>	M00-M99, БОЛЕЗНИ КОСТНО-МЫШЕЧНОЙ СИСТЕМЫ И СОЕДИНИТЕЛЬНОЙ ТКАНИ
<i>E</i> <sub>14</sub>	Q00-Q99, ВРОЖДЕННЫЕ АНОМАЛИИ [ПОРОКИ РАЗВИТИЯ], ДЕФОРМАЦИИ И ХРОМОСОМНЫЕ НАРУШЕНИЯ
<i>E</i> <sub>15</sub>	L00-L99, БОЛЕЗНИ КОЖИ И ПОДКОЖНОЙ КЛЕТЧАТКИ
<i>E</i> <sub>16</sub>	O00-O99, БЕРЕМЕННОСТЬ, РОДЫ И ПОСЛЕРОДОВОЙ ПЕРИОД
<i>E</i> <sub>17</sub>	P00-P96, ОТДЕЛЬНЫЕ СОСТОЯНИЯ, ВОЗНИКАЮЩИЕ В ПЕРИНАТАЛЬНОМ ПЕРИОДЕ
<i>E</i> <sub>18</sub>	H60-H95, БОЛЕЗНИ УША И СОСЦЕВИДНОГО ОТРОСТКА
<i>E</i> <sub>19</sub>	H00-H59, БОЛЕЗНИ ГЛАЗА И ЕГО ПРИДАТОЧНОГО АППАРАТА
<i>E</i> <sub>20</sub>	«СМЕРТЬ»

Теоретической основой представленного разбиения по состояниям здоровья населения, включая смертность, является МКБ-10\*. В таблице 7.1 указаны принятые в системе здравоохранения коды соответствующих заболеваний.

\* Международная классификация болезней десятого пересмотра (МКБ-10) – система группировки болезней и патологических состояний, отражающая современный этап развития медицинской науки. Является основным нормативным документом при изучении состояния здоровья населения в странах-членах ВОЗ (Всемирной организации здравоохранения). Вступила в силу с 01.01.1993 г. (в России с 01.01.1999г.).

Поставим задачу – классифицировать (распределить по состояниям) все население, представленное в ПБД [147, 149]. Отметим, что каждый индивидуум в ПБД имеет свой уникальный идентификационный номер и «послужной список» заболеваемости, согласно которому и производится классификация.

Все зарегистрированные в ПБД на определенную дату, разбиваются на две группы – «относительно здоровые» и «заболевшие». При этом индивидуумы, не имеющие ранее зафиксированных тяжелых (неизлечимых) заболеваний, и ни разу не обратившиеся за медицинской помощью автоматически считаются относящимися к первой группе – «относительно здоровые» (состояние  $E_0$ ). Все остальные относятся к группе «заболевшие».

Каждому индивидууму из ПБД однозначно присваивается состояние  $E_i$  из классификации состояний, соответствующее его состоянию здоровья. Каждое состояние системы из  $E_1, E_2, \dots, E_{19}$  – это определенная группа заболеваний в соответствии с классами МКБ-10. Важнейшим условием применимости методики является несовместность состояний, т.е. один индивидуум в каждый момент времени может находиться лишь в одном состоянии. Решение проблемы однозначности обеспечивает алгоритм формирования требований для каждого состояния в соответствии с МКБ-10 и, при наличии нескольких заболеваний, определяющий приоритет (иерархию) этих заболеваний в соответствии с их тяжестью.

Отнесение состояния здоровья индивидуума к одной из групп тяжести, производится с учетом тяжести последствий, продолжительности и значимости конкретного заболевания. Всего определено три группы тяжести – «легкие заболевания», «заболевания средней тяжести», «тяжелые заболевания». Так как полностью автоматизировать процесс сортировки диагнозов по степени тяжести не представляется возможным, то сортировка проводится в несколько этапов с обязательным привлечением экспертов:

- все диагнозы сортируются по полю ПБД «исход лечения»; если хотя бы один исход лечения закончился смертью - то такое заболевание следует отнести к группе «тяжелые заболевания»;
- оставшиеся диагнозы вновь сортируются на две группы, «заболевания средней тяжести» и «легкие заболевания», следующим образом: если в результате лечения здоровье индивида не ухудшалось, то диагноз следует отнести к группе «легкие заболевания», оставшиеся - к группе «заболевания средней тяжести».

При формировании состояний системы «легкие заболевания» игнорируются, как не имеющие стойких последствий для здоровья. Например, человек, переболевший ОРЗ без дальнейших осложнений, не будет зафиксирован в состоянии «болезни органов дыхания».

На следующем этапе сформированные таким образом группы должны быть экспертно оценены. Суть экспертных оценок – проверка на достоверность выводов в малочисленных группах (статистически слабо достоверных) и разбор сомнительных (как правило, вызванных ошибками регистрации) и нестандартных случаев. Часть диагнозов после этого перераспределяется по другим группам, в том числе перераспределение происходит и в группу «относительно здоровые» -  $E_0$ . Таким образом, внутри каждого класса МКБ-10 выстраивается внутриклассовая иерархия диагнозов по степени тяжести.

Если индивидуум за предшествующий промежуток времени обращался несколько раз с разными заболеваниями (которые находятся в разных классах по МКБ-10), то в расчет берется заболевание, самое неблагоприятное по межклассовой иерархии. Если в течение года обращался с разными заболеваниями, которые находятся в одном классе МКБ-10, то используется внутриклассовая иерархия. Если два или более диагнозов имеют одинаковый ранг по тяжести, то необходимо использовать иерархию по причинам смерти (межклассовая иерархия).

Межклассовая иерархия по степени тяжести может быть построена в соответствии со структурой причин смертности населения региона. В табл. 7.1 состояния упорядочены в соответствии именно с этой структурой.

Таким образом, при наличии нескольких заболеваний классификация состояний осуществляется в соответствии со степенью тяжести заболевания и в соответствии с установленной иерархией состояний, внутрикласовой и межклассовой. Отметим в данном процессе классификации важную контролируемую функцию экспертов.

При обработке больших массивов реальных данных могут появиться неожиданные ошибки, вызванные неправильной регистрацией данных. Для их устранения в компьютерную программу обработки данных следует ввести некоторые ограничители. Например, в состояние  $E_{16}$  могут попасть только женщины фертильного возраста, а в состояние  $E_{17}$  – только дети в возрастной категории от 0 до 1 года. Соответствующие факты необходимо учесть не на экспертном, а на программном уровне.

### **7.3.2 Формирование стохастических матриц и вычисление безусловных вероятностей. Вероятностная модель жизни человека**

Возраст индивидуумов определяется как дискретная величина  $t_i$ , середина  $i$ -го возрастного интервала,  $i = 0, 1, 2, \dots, m$ . Естественно и весьма заманчиво при масштабном исследовании в качестве длины возрастного интервала использовать один год. Однако в реальной диагностике по обращаемости, при отсутствии ежегодной диспансеризации населения, которая могла бы гарантировать объективную картину заболеваемости, использование одногодичного промежутка не всегда оправдано. Во-первых, необращение за медицинской помощью в течение года часто вызвано не отличным здоровьем, а менталитетом населения, устоявшейся привычкой общения с врачом лишь в «крайнем случае». Во-вторых, наличие большого количества интервалов приведет к тому, что при реализации методики в некоторых группах (возраст, состояние здоровья) окажется недостаточно наблюдений, что, несомненно,

негативно отразится на точности вычислений. Исходя из сказанного, в реальном исследовании длина возрастного интервала должна составлять четыре или пять лет.

Статистические данные состояния здоровья населения, рассматриваемые в динамике, позволяют рассчитать конкретные значения переходных вероятностей из состояния в состояние для различных половозрастных групп. Отметим, что мужское и женское население в исследовании относятся к разным группам. Всё население группы на определенный начальный момент времени, например, в исследовании, проведенном в Новгородском медицинском центре СЗО РАМН, на 01.01.2001г., классифицируется по возрастным интервалам и состояниям системы [149]. Таким образом, каждый индивидуум оказывается в определенном возрастном интервале с однозначно определенным состоянием здоровья.

Далее в указанном исследовании при построении матриц те же индивидуумы, зарегистрированные в ПБД в 2001г., но со сдвигом возраста на пять лет (в нашем случае на момент времени 01.01.2006г), были вновь отсортированы согласно классификации состояний системы. Причем, все зарегистрированные в ПБД и умершие за период с 01.01.2001 г. по 31.12.2005 г., попали в состояние, соответствующие смертности,  $E_{20}$ . Индивидуумы, не попавшие в базу данных на начальную или конечную дату (в основном это касается мигрантов), из исследования исключались.

Рассмотренным методом формировались пары переходов для каждого индивидуума из состояния  $E_i$  в 2001 г. в состояние  $E_j$  на начало 2006 года, т.е. шаг процесса в нашем исследовании – пять лет. За 5 лет каждый индивидуум обязательно переходит в следующий возрастной интервал и в какое-то из состояний (может остаться и в том же состоянии). Для пятилетнего возрастного интервала, вычислив по статистическим данным переходные вероятности из каждого состояния  $E_i$  в состояние  $E_j$ , как отношение числа перешедших в  $E_j$  из  $E_i$  к числу всех, находящихся в  $E_i$  на начало исследования, создаем из них матрицу. А для всех возрастных интервалов получаем набор

стохастических матриц:  $\mathbf{P}(1), \mathbf{P}(2), \dots, \mathbf{P}(m)$ . Отметим, что в проведенном исследовании в силу специфики здоровья в младенческом возрасте первый возрастной интервал взят (0-1) полных лет, следующий – (2-4) полных лет. Поэтому для возраста (0-4) лет сформированы два интервала и соответственно вычислены две стохастические матрицы. Вследствие уменьшающегося с возрастом количества населения (соответственно и объема данных) последний возрастной интервал, 85 лет и более, также является не пятилетним, из этого интервала при любом состоянии здоровья переход возможен лишь в  $E_{20}$ . Остальные возрастные интервалы – пятилетние: (5-9), (10-14), ..., (80-84) полных лет. Следовательно, общее количество интервалов – 19.

**Итак**, после нахождения стохастических матриц (в рассмотренном случае  $\mathbf{P}(1), \mathbf{P}(2), \dots, \mathbf{P}(19)$ ) неоднородная цепь Маркова оказывается полностью определенной для каждого вектора начальных данных  $p_1(0), p_2(0), \dots, p_{20}(0)$ .

Человеческая жизнь моделируется [150], как последовательное объединение возрастных интервалов в совокупности с состояниями здоровья (табл. 7.2).

Таблица 7.2. Модель человеческой жизни (возраст, состояние здоровья)

Номер интервала	0	1	2	...	$v$
Возраст и состояние	$(t_0, E_{j_0})$	$(t_1, E_{j_1})$	$(t_2, E_{j_2})$	...	$(t_v, E_{j_v})$
Вероятность	$a_{j_0}$	$p_{j_0 j_1}(1)$	$p_{j_1 j_2}(2)$	...	$p_{j_{v-1} j_v}(v)$

Состояние  $E_{j_v}$ , последнее в цепи состояний, в отличие от предыдущих, является поглощающим.

Для изучения здоровья населения используется модель условного поколения, в основе которой лежит упорядочивание по времени значений статистического показателя, взятых за короткий временной промежуток для населения всех возрастов. Значения показателя, представляющие один и тот же период времени, но относящиеся к разным возрастным группам, выстраиваются в виде временного ряда. При объединении значений показателей в раз-

ных возрастах предполагается, что наблюдаемые события, связанные с жизнедеятельностью и здоровьем населения, произошли не у разных поколений в одно время, а в разном возрасте у одного и того же поколения, которое и называется «условным».

Использование математического аппарата цепей Маркова позволяет по начальным данным вычислить наборы безусловных вероятностей и, следовательно, осуществить вероятностный прогноз состояния здоровья, включая смертность, как отдельного человека, так и целых групп населения на последующие периоды жизни. Для вычисления вектора безусловных вероятностей на  $k$ -м шаге используется формула (7.14).

$$\{p_j(k)\} = \{p_j(0)\} \mathbf{P}(1) \cdot \mathbf{P}(2) \cdot \dots \cdot \mathbf{P}(k), \text{ где } k = 1, 2, \dots, 19.$$

Любой из векторов безусловных вероятностей является вероятностным прогнозом для соответствующего возраста. Совокупность векторов безусловных вероятностей, найденных на каждом шаге вычисления, можно записать в виде матрицы, в которой каждая строка представляет возрастной интервал, а каждый столбец – состояние здоровья. Элементами, находящимися на пересечении  $i$ -й строки и  $j$ -го столбца, являются вероятности в  $i$ -м возрастном интервале оказаться в состоянии  $E_j$ .

Например, [147, 151] для «поколение 0 лет» мужского пола согласно реальным данным (поскольку у новорожденных многие из заболеваний выявляются лишь по прошествии некоторого времени, то фактически это данные первого возрастного интервала) полагаем начальный вектор (0,515; 0; 0; 0,001; 0,005; 0; 0,003; 0,001; 0,001; 0; 0; 0; 0; 0; 0,052; 0,003; 0; 0,419; 0; 0; 0). Тогда совокупности безусловных вероятностей будущих состояний здоровья для середины каждого последующего возрастного интервала записаны в виде векторов (строки), образующих матрицу (табл.7.3).

Таблица 7.3. Матрица безусловных вероятностей (мужское население).

Возраст	$E_0$	$E_1$	$E_2$	$E_3$	$E_4$	$E_5$	$E_6$	$E_7$	...	$E_{19}$	$E_{20}$
---------	-------	-------	-------	-------	-------	-------	-------	-------	-----	----------	----------

2-4	0,353	0,013	0,006	0,024	0,147	0,002	0,009	0,002	...	0,026	0,009
5-9	0,090	0,021	0,021	0,016	0,144	0,090	0,359	0,074	...	0,009	0,021
10-14	0,046	0,035	0,038	0,014	0,133	0,053	0,500	0,044	...	0,011	0,025
15-19	0,037	0,052	0,066	0,014	0,142	0,038	0,440	0,031	...	0,014	0,034
20-24	0,104	0,032	0,101	0,013	0,166	0,018	0,282	0,030	...	0,023	0,056
25-29	0,150	0,022	0,148	0,015	0,173	0,009	0,174	0,035	...	0,021	0,086
30-34	0,149	0,025	0,151	0,014	0,161	0,010	0,161	0,024	...	0,022	0,130
35-39	0,143	0,028	0,122	0,013	0,137	0,010	0,156	0,023	...	0,023	0,186
40-44	0,128	0,039	0,107	0,013	0,113	0,009	0,140	0,018	...	0,025	0,255
45-49	0,100	0,037	0,087	0,017	0,093	0,008	0,124	0,015	...	0,029	0,346
50-54	0,081	0,040	0,064	0,014	0,068	0,007	0,105	0,013	...	0,029	0,453
55-59	0,053	0,042	0,043	0,014	0,048	0,007	0,089	0,009	...	0,024	0,568
60-64	0,041	0,035	0,026	0,013	0,031	0,004	0,081	0,007	...	0,017	0,677
65-69	0,031	0,030	0,012	0,010	0,015	0,003	0,071	0,006	...	0,012	0,766
70-74	0,023	0,019	0,009	0,007	0,010	0,002	0,051	0,004	...	0,009	0,839
75-79	0,014	0,014	0,003	0,005	0,005	0,001	0,030	0,002	...	0,007	0,901
80-84	0,008	0,009	0,001	0,002	0,003	0,001	0,016	0,001	...	0,004	0,944
85-	0,005	0,004	0,000	0,001	0,001	0,001	0,007	0,001	...	0,002	0,974
-	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	...	0,000	1,000

Графическая иллюстрация матрицы безусловных вероятностей для мужского населения представлена на рис. 7.3.

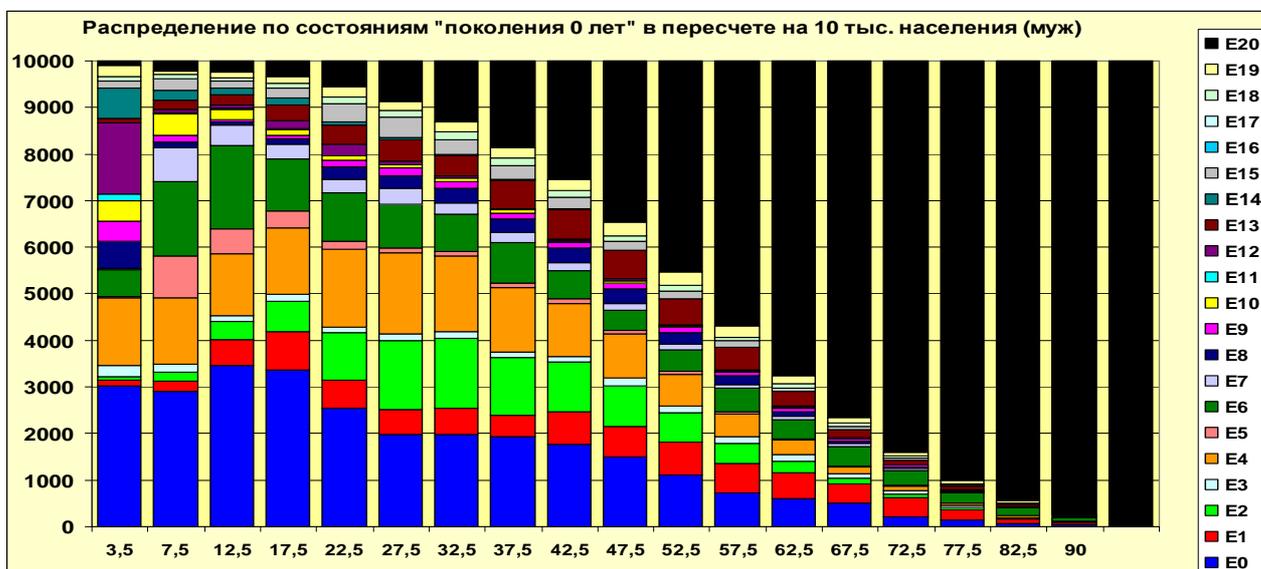
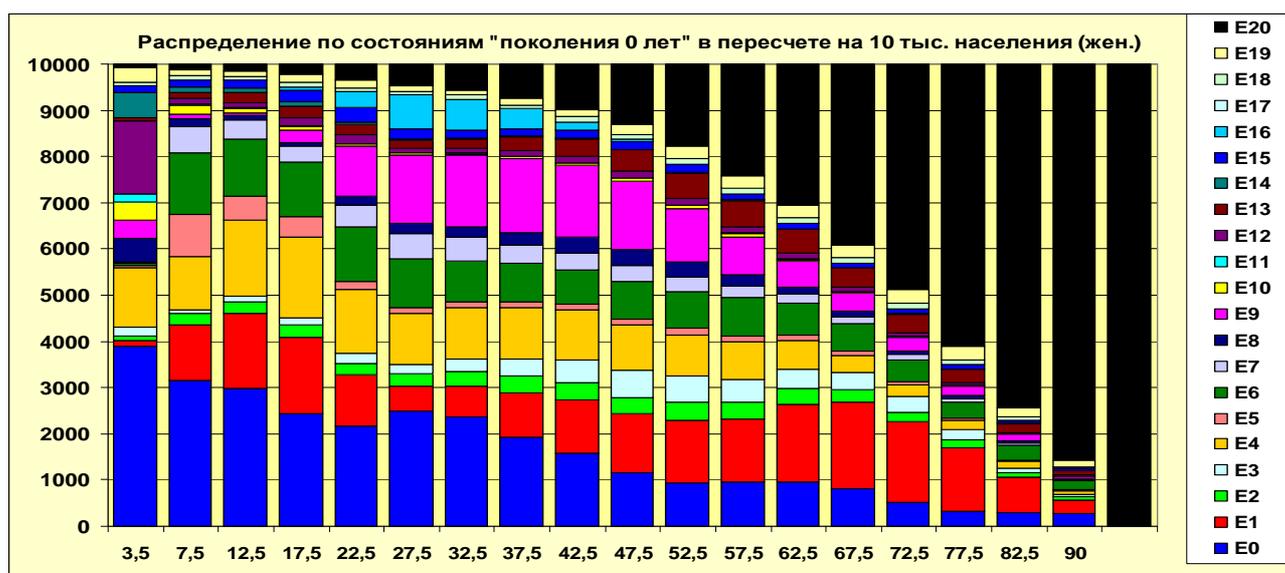


Рис. 7.3. Распределение «поколения 0 лет» в пересчете на 10 тыс. населения (мужчины)

Здесь распределение «поколения 0 лет» в пересчете на 10 тыс. населения обозначено по состояниям здоровья, включая состояние  $E_{20}$ , в зависимо-

сти от возраста. По оси абсцисс указаны середины возрастных интервалов. Каждому возрастному интервалу соответствует вертикальный столбец (10 тыс. мужского населения), который разбит на части в соответствии с долей каждого из возможных состояний в общей структуре заболеваемости и смертности для указанного возраста. В столбце состояния  $E_0, E_1, \dots, E_{19}, E_{20}$  упорядочены снизу вверх. Таким образом, для каждого возрастного интервала воссоздан спектр заболеваемости по состояниям с учетом и числа умерших в поколении к началу возрастного интервала. Середина предпоследнего на графике интервала «85 лет и более» обозначена, как 90 лет. Далее, вполне естественно, заболеваемость отсутствует, все 10 тыс. человек оказываются в состоянии  $E_{20}$ , что и зафиксировано на графике. Легко заметить, что в целом для мужского населения наиболее типичными являются состояния  $E_0$  (относительно здоров),  $E_4$  (болезни органов дыхания),  $E_6$  (болезни органов пищеварения),  $E_2$  (травмы и отравления),  $E_1$  (болезни системы кровообращения),  $E_{13}$  (болезни костно-мышечной системы).

Аналогично графическая иллюстрация матрицы безусловных вероятностей женского населения («поколение 0 лет») представлена на рис. 7.4.



**Рис. 7.4.** Распределение «поколения 0 лет» в пересчете на 10 тыс. населения (женщины)

В данном случае по возрастной спектр заболеваемости выглядит несколько иначе. В отличие от распределения на рис. 7.3, обращает на себя внимание в целом более низкий уровень смертности женщин. Также отметим иное количественное распределение заболеваемости:  $E_0$  (относительно здоров);  $E_1$  (болезни системы кровообращения);  $E_9$  (болезни мочеполовой системы);  $E_4$  (болезни органов дыхания);  $E_6$  (болезни органов пищеварения);  $E_{13}$  (болезни костно-мышечной системы). При этом показатели  $E_1$  и  $E_9$  существенно выше, чем у мужчин. Весьма заметны в разложении доли  $E_3$  (новообразования) и  $E_7$  (некоторые инфекционные и паразитарные болезни), превышающие соответствующие показатели для мужского населения. Значительно ниже, чем у мужского населения, доли состояния  $E_2$  (травмы и отравления). Практически во всех возрастах, как среди мужского, так и женского населения, отчетливо прослеживается доля состояния  $E_{19}$  (болезни глаза и его придаточного аппарата).

Отметим, что векторы безусловных вероятностей, характеризующие состояние здоровья поколения в будущем, можно найти и исследовать не только для «поколения 0 лет», но и для любых других возрастов и распределений по состояниям здоровья

Найденные в модели векторы безусловных вероятностей можно интерпретировать не только как вероятностный прогноз, но и, в первую очередь, как многомерные показатели текущего состояния здоровья. Полученные значения, хотя формально и являются «научно обоснованным предсказанием» на  $n$  лет вперед при неизменном комплексе условий окружающей среды, всё же научно-исследовательскую ценность представляют именно сейчас, а не через  $n$  лет. Например, прогноз показывает, что при сохранении текущего комплекса условий половина мужского населения в поколении новорожденных перейдет в состояние «смерть» через 54,5 лет. Сбудется этот прогноз или нет – покажет время. Однако значение 54,5 является более важным в настоящий момент времени, как оценка потенциала здоровья населения в поколе-

нии. И текущей задачей органов управления, органов здравоохранения и самого населения является не сохранение, а изменение комплекса условий (по сути, качества жизни) так, чтобы значение этого показателя существенно возросло.

### 7.3.3 Вычисление средней продолжительности жизни по группам.

#### Сравнение групп

Наборы безусловных вероятностей, находимых по вышеуказанной методике, данные органов здравоохранения и страховых кампаний о состоянии здоровья, а также половозрастные характеристики населения позволяют ввести ряд самостоятельных показателей, характеризующих здоровье и смертность населения. На основе созданной базы рассмотрим и некоторые новые показатели [147, 150, 151].

Введём случайную величину  $X_{ij}$  для каждого наблюдаемого в возрасте  $t_i$ , где  $i$  - номер соответствующего возрастного интервала (например, в модели, рассматриваемой в 7.3.1, пятигодичного), и с состоянием здоровья  $E_j$ , где  $E_j$  - одно из возможных состояний здоровья, исключая состояние «смерть». Смысл  $X_{ij}$  - общее количество лет жизни человека, дожившего до возраста  $t_i$  и находящегося при этом возрасте в состоянии  $E_j$ . Например,  $X_{6,4}$  - общее количество лет жизни индивидуума, находящегося на момент исследования в 6-м возрастном интервале (20-24 полных лет), и с состоянием здоровья  $E_4$  (болезни органов дыхания). Всего получается  $mn$  случайных величин, где  $m$  - количество возрастных интервалов,  $n$  - количество состояний  $E_j$  без состояния  $E_n$ . Запишем их в матричной форме

$$\mathbf{X} = \begin{pmatrix} X_{1,0} & X_{1,1} & \dots & X_{1,n-1} \\ X_{2,0} & X_{2,1} & \dots & X_{2,n-1} \\ \dots & \dots & \dots & \dots \\ X_{m,0} & X_{m,1} & \dots & X_{m,n-1} \end{pmatrix}.$$

Для каждой случайной величины  $X_{ij}$  закон распределения, обладая своими персональными параметрами, имеет вид

$x$	$x_{i+1}$	$x_{i+2}$	$x_{i+3}$	.....	$x_{m-1}$	$x_m$
$p$	$p_{i+1}$	$q_{i+1}p_{i+2}$	$q_{i+1}q_{i+2}p_{i+3}$	.....	$q_{i+1}q_{i+2}\dots q_{m-2}p_{m-1}$	$q_{i+1}q_{i+2}\dots q_{m-1}1$

где  $x_{i+1}, x_{i+2}, \dots, x_m$  - середины каждого из последующих за  $i$ -м интервала<sup>\*)</sup>,  $p_s$  – вероятность умереть в возрасте  $x_s$ ,  $q_s = 1 - p_s$ . Для последнего возрастного интервала закон распределения  $X_{mj}$  тривиален:

$x$	$x_m$
$p$	1

Зная закон распределения  $X_{ij}$ , легко найти среднее и дисперсию:

$$M(X_{ij}) = \sum xp = M_{ij}, \quad D(X_{ij}) = M(X_{ij}^2) - M^2(X_{ij}) = D_{ij}.$$

Найденные для всех случайных величин характеристики удобно записать в матричном виде

$$\mathbf{M} = \begin{pmatrix} M_{1,0} & M_{1,1} & \dots & M_{1,n-1} \\ M_{2,0} & M_{2,1} & \dots & M_{2,n-1} \\ \dots & \dots & \dots & \dots \\ M_{m,0} & M_{m,1} & \dots & M_{m,n-1} \end{pmatrix}; \quad \mathbf{D} = \begin{pmatrix} D_{1,0} & D_{1,1} & \dots & D_{1,n-1} \\ D_{2,0} & D_{2,1} & \dots & D_{2,n-1} \\ \dots & \dots & \dots & \dots \\ D_{m,0} & D_{m,1} & \dots & D_{m,n-1} \end{pmatrix}.$$

<sup>\*)</sup> В последнем интервале, имеющем длину большую, чем другие, значение  $x_m$  определяется статистически и не совпадает с серединой интервала.

Смысл величины  $M_{ij}$  - средняя ожидаемая продолжительность жизни индивидуумов, находящихся в настоящее время в  $i$ -м возрастном интервале в состоянии здоровья  $E_j$ .

В матрице  $\mathbf{M}$  номер строки – номер возрастного интервала. А каждая из строк представляет собой набор значений средней продолжительности жизни при различных состояниях здоровья в этом возрастном интервале.

Таким образом, матрица  $\mathbf{M}$ , по существу, является матрицей значений средней продолжительности жизни в зависимости от возраста и текущего состояния здоровья. Поскольку найдены выборочные дисперсии  $D_{ij}$ , то для каждого параметра  $\Theta_{ij}$  генеральной совокупности, являющегося «истинным значением» средней продолжительности жизни, можно построить доверительный интервал с надёжностью  $\gamma$ :

$$M_{ij} - t_{\gamma} \frac{\sqrt{D_{ij}}}{\sqrt{n_{ij}}} < \Theta_{ij} < M_{ij} + t_{\gamma} \frac{\sqrt{D_{ij}}}{\sqrt{n_{ij}}},$$

где  $n_{ij}$  - количество соответствующих наблюдений,  $t_{\gamma}$  - коэффициент доверия. При больших объёмах статистических данных доверительные интервалы оказываются вполне приемлемыми.

Матрица значений средней продолжительности жизни, рассчитанная по реальным данным здоровья и смертности населения Новгородской области, представлена в табл. 7.4. Отметим, что в последнем из интервалов «85 и старше» значения показателя по состояниям не рассчитывались, и значение «90, 00 лет» является оценочным. Также отметим, что полученные значения, вычисленные как математическое ожидание случайной величины, отличаются в меньшую сторону от показателя, вычисляемого, как медиана, что вполне естественно для распределений с ярко выраженной левосторонней асимметрией.

**Таблица 7.4. Матрица средней продолжительности жизни (мужское население)**

Возраст	m0	m1	m2	m3	m4	m5	m6	m7	m8	...	m18	m19
0-1	51,01	51,87	51,44	51,56	51,32	51,84	51,76	51,83	50,95	...	51,64	51,83
2-4	51,44	51,87	51,44	51,80	51,48	51,84	51,85	51,77	51,10	...	51,64	51,83
5-9	51,76	52,04	51,93	51,39	51,83	51,71	51,89	51,90	51,99	...	52,00	51,75
10-14	52,12	51,98	51,66	51,93	52,00	52,13	52,08	52,09	51,44	...	51,92	52,04
15-19	52,73	52,50	52,39	51,35	52,47	52,10	52,38	52,46	52,13	...	52,48	52,51
20-24	53,62	53,28	52,62	52,82	53,39	52,65	53,22	52,78	53,20	...	53,29	53,42
25-29	54,69	54,11	53,66	54,04	54,51	53,54	54,27	54,15	54,26	...	54,96	54,87
30-34	56,21	55,61	55,04	54,05	55,91	55,89	55,46	55,32	55,34	...	56,08	56,08
35-39	57,84	56,52	56,77	55,94	57,26	57,82	57,40	57,52	56,33	...	57,63	56,88
40-44	59,94	58,21	58,36	57,95	59,28	58,98	59,29	59,01	58,55	...	59,18	59,18
45-49	62,22	60,70	61,15	59,70	61,46	62,21	61,76	61,12	60,66	...	61,14	61,75
50-54	65,05	63,06	63,43	61,60	64,95	64,60	64,84	62,57	63,34	...	63,92	64,57
55-59	68,42	66,63	67,21	65,06	67,63	68,09	67,88	66,35	67,24	...	67,84	67,37
60-64	71,45	70,16	69,82	69,40	71,28	71,18	71,78	71,39	70,57	...	71,96	71,21
65-69	74,68	73,64	74,48	73,49	74,30	74,98	75,33	74,00	74,01	...	75,61	74,71
70-74	77,78	77,36	77,17	77,15	78,21	78,17	78,72	78,85	78,47	...	78,36	78,31
75-79	81,10	81,59	80,87	81,33	81,91	81,20	82,55	81,16	83,19	...	82,29	82,42
80-84	85,09	85,57	85,84	84,81	85,18	85,56	86,49	82,00	86,44	...	86,69	85,90
85 и старше	90,00	90,00	90,00	90,00	90,00	90,00	90,00	90,00	90,00	...	90,00	90,00

Исследование и сравнение конкретных значений матрицы **M** представляет несомненный интерес для медицинских работников, занимающихся изучением, как отдельных видов заболеваемости, так и охраной общественного здоровья в целом. При этом вновь подчеркнем, что значения показателя  $M_{ij}$  и связанных с ним следует трактовать не столько как прогноз, который может оправдаться лишь при сохранении неизменными условий жизни на прогнозируемый период (что маловероятно), сколько как показатель нынешнего состояния заболеваемости и общественного здоровья в целом.

#### 7.3.4 Вычисление показателей продолжительности жизни фактического населения

При известной средней продолжительности жизни  $M_{ij}$  для каждой группы  $(t_i, E_j)$ , опираясь на реальные данные половозрастной структуры населения, можно вычислять [150, 152] средние показатели, связанные с про-

должительностью жизни. Введём матрицу, представляющую структуру возрастного состава населения в соответствии с состоянием здоровья:

$$\mathbf{Y} = \begin{pmatrix} y_{1,0} & y_{1,1} & \cdots & y_{1,n-1} \\ y_{2,0} & y_{2,1} & \cdots & y_{2,n-1} \\ \cdots & \cdots & \cdots & \cdots \\ y_{m,0} & y_{m,1} & \cdots & y_{m,n-1} \end{pmatrix},$$

где  $y_{ij}$  - количество человек в пересчете на 10000 населения одного пола, проживающих на определённой территории в возрасте  $t_i$  с состоянием здоровья  $E_j$ . Данные рассматриваются отдельно для мужчин и женщин.

Найдём произведения матриц  $\mathbf{Y}\mathbf{M}^T$  и  $\mathbf{Y}^T\mathbf{M}$  :

$$\begin{aligned} \mathbf{Y}\mathbf{M}^T &= \begin{pmatrix} y_{1,0} & y_{1,1} & \cdots & y_{1,n-1} \\ y_{2,0} & y_{2,1} & \cdots & y_{2,n-1} \\ \cdots & \cdots & \cdots & \cdots \\ y_{m,0} & y_{m,1} & \cdots & y_{m,n-1} \end{pmatrix} \begin{pmatrix} \mathbf{M}_{1,0} & \mathbf{M}_{2,0} & \cdots & \mathbf{M}_{m,0} \\ \mathbf{M}_{1,1} & \mathbf{M}_{2,1} & \cdots & \mathbf{M}_{m,1} \\ \cdots & \cdots & \cdots & \cdots \\ \mathbf{M}_{1,n-1} & \mathbf{M}_{2,n-1} & \cdots & \mathbf{M}_{m,n-1} \end{pmatrix} = \\ &= \begin{pmatrix} \tau_{1\bullet} & \cdots & \cdots & \cdots \\ \cdots & \tau_{2\bullet} & \cdots & \cdots \\ \cdots & \cdots & \ddots & \cdots \\ \cdots & \cdots & \cdots & \tau_{m\bullet} \end{pmatrix}; \\ \mathbf{Y}^T\mathbf{M} &= \begin{pmatrix} y_{1,0} & y_{2,0} & \cdots & y_{m,0} \\ y_{1,1} & y_{2,1} & \cdots & y_{m,1} \\ \cdots & \cdots & \cdots & \cdots \\ y_{1,n-1} & y_{2,n-1} & \cdots & y_{m,n-1} \end{pmatrix} \begin{pmatrix} \mathbf{M}_{1,0} & \mathbf{M}_{1,1} & \cdots & \mathbf{M}_{1,n-1} \\ \mathbf{M}_{2,0} & \mathbf{M}_{2,1} & \cdots & \mathbf{M}_{2,n-1} \\ \cdots & \cdots & \cdots & \cdots \\ \mathbf{M}_{m,0} & \mathbf{M}_{m,1} & \cdots & \mathbf{M}_{m,n-1} \end{pmatrix} = \\ &= \begin{pmatrix} \zeta_{\bullet 0} & \cdots & \cdots & \cdots \\ \cdots & \zeta_{\bullet 1} & \cdots & \cdots \\ \cdots & \cdots & \ddots & \cdots \\ \cdots & \cdots & \cdots & \zeta_{\bullet n-1} \end{pmatrix}. \end{aligned}$$

На главной диагонали каждой из получившихся матриц находятся значения количества лет жизни 10000 населения, причём в первом случае,  $\mathbf{Y}\mathbf{M}^T$ , разложенные по возрастным группам, а во втором случае,  $\mathbf{Y}^T\mathbf{M}$ , - по группам

состояния здоровья. Сумма элементов главной диагонали – общее количество лет жизни, приходящееся на 10000 населения:

$$\sum_{i=1}^m \tau_{i\bullet} = \sum_{j=0}^{n-1} \zeta_{\bullet j}$$

Поделив количество лет жизни на соответствующее количество населения, можно вычислить средние продолжительности жизни в зависимости от возраста или состояния здоровья. В частности,

$$\bar{M}_{i\bullet} = \frac{\tau_{i\bullet}}{\sum_{j=0}^{n-1} y_{ij}} - \text{средняя ожидаемая продолжительность жизни населения,}$$

находящегося в  $i$ -м возрастном интервале (вне зависимости от состояния здоровья);

$$\bar{M}_{\bullet j} = \frac{\zeta_{\bullet j}}{\sum_{i=0}^m y_{ij}} - \text{средняя ожидаемая продолжительность жизни населения,}$$

находящегося в  $j$ -м состоянии здоровья (вне зависимости от возраста);

$$\bar{M} = \frac{1}{10000} \sum_{i=1}^m \tau_{i\bullet} - \text{средняя ожидаемая продолжительность жизни всего}$$

населения.

Зная все значения средней продолжительности жизни,  $M_{ij}$ , и текущий возраст  $t_i$ , легко найти ожидаемое количество лет *оставшейся* жизни в зависимости от возраста и состояния здоровья. Соответствующие значения запишем в матричном виде:

$$\mathbf{Q} = \begin{pmatrix} M_{1,0} - t_1 & M_{1,1} - t_1 & \dots & M_{1,n-1} - t_1 \\ M_{2,0} - t_2 & M_{2,1} - t_2 & \dots & M_{2,n-1} - t_2 \\ \dots & \dots & \dots & \dots \\ M_{m,0} - t_m & M_{m,1} - t_m & \dots & M_{m,n-1} - t_m \end{pmatrix} = \begin{pmatrix} q_{1,0} & q_{1,1} & \dots & q_{1,n-1} \\ q_{2,0} & q_{2,1} & \dots & q_{2,n-1} \\ \dots & \dots & \dots & \dots \\ q_{m,0} & q_{m,1} & \dots & q_{m,n-1} \end{pmatrix},$$

где  $q_{ij} = M_{ij} - t_i$ . Используя матрицу  $\mathbf{Y}$  структуры реального населения в пересчете на 10000 населения, можно получить матрицы  $\mathbf{YQ}^T$  и  $\mathbf{Y}^T\mathbf{Q}$ , в которых

на главных диагоналях представлены разложения ожидаемого количества лет будущей жизни соответственно по возрастам и состояниям здоровья:

$$\mathbf{YQ}^T = \begin{pmatrix} \gamma_{1\bullet} & \dots & \dots & \dots \\ \dots & \gamma_{2\bullet} & \dots & \dots \\ \dots & \dots & \ddots & \dots \\ \dots & \dots & \dots & \gamma_{m\bullet} \end{pmatrix}; \quad \mathbf{Y}^T\mathbf{Q} = \begin{pmatrix} \delta_{\bullet 0} & \dots & \dots & \dots \\ \dots & \delta_{\bullet 1} & \dots & \dots \\ \dots & \dots & \ddots & \dots \\ \dots & \dots & \dots & \delta_{\bullet n-1} \end{pmatrix},$$

где  $\gamma_{i\bullet}$  - суммарное количество лет оставшейся жизни населения в возрасте  $t_i$  на 10000 реального населения,  $\delta_{\bullet j}$  - суммарное количество лет будущей жизни населения в состоянии здоровья  $E_j$  на 10000 реального населения. Разумеется,

$$\sum_{i=1}^m \gamma_{i\bullet} = \sum_{j=0}^{n-1} \delta_{\bullet j}.$$

Аналогично могут быть найдены величины

$$\bar{L}_{i\bullet} = \frac{\gamma_{i\bullet}}{\sum_{j=0}^{n-1} y_{ij}} - \text{средняя ожидаемая продолжительность оставшейся жизни}$$

населения, находящегося в  $i$ -м возрастном интервале (вне зависимости от состояния здоровья);

$$\bar{L}_{\bullet j} = \frac{\delta_{\bullet j}}{\sum_{i=0}^m y_{ij}} - \text{средняя ожидаемая продолжительность оставшейся жизни}$$

населения, находящегося в  $j$ -м состоянии здоровья (вне зависимости от возраста);

$$\bar{L} = \frac{1}{10000} \sum_{i=1}^m \gamma_{i\bullet} - \text{средняя ожидаемая продолжительность оставшейся}$$

жизни для всего населения.

### 7.3.5 Показатели, характеризующие состояния здоровья

Показатель средней ожидаемой продолжительности жизни в зависимости от возраста и состояния здоровья является весьма информативным при

исследовании здоровья сообщества. Однако наборы значений  $M_{ij}$  можно использовать [150-152] и для характеристики самих заболеваний (состояний  $E_j$  при определенных возрастах).

Очевидно, что наибольшим значением  $M_{ij}$  в  $i$ -й строке матрицы  $\mathbf{M}$  должно быть число  $M_{i0}$ , соответствующее состоянию  $E_0$  (т.е. «относительно здоров») в  $i$ -м интервале. Для индивидуума возраста  $t_i$  введем показатель

$$\pi_{ij} = M_{i0} - M_{ij},$$

равный потере лет оставшейся жизни вследствие наличия в настоящий момент какого-либо заболевания, представляющего состояние  $E_j$ . Таким образом, показатель  $\pi_{ij}$  является характеристикой состояния  $E_j$  для возраста  $t_i$ . Сравнение проводится с состоянием  $E_0$  при том же возрасте. Как легко заметить,  $\pi_{i0} = 0$  для всех  $i$ . Значения  $\pi_{i0}$  вычисляются в абсолютных единицах (годах).

Множество показателей потерь  $\pi_{ij}$  запишем в матричном виде:

$$\mathbf{\Pi} = \begin{pmatrix} 0 & M_{1,0} - M_{1,1} & \dots & M_{1,0} - M_{1,n-1} \\ 0 & M_{2,0} - M_{2,1} & \dots & M_{2,0} - M_{2,n-1} \\ \dots & \dots & \dots & \dots \\ 0 & M_{m,0} - M_{m,1} & \dots & M_{m,0} - M_{m,n-1} \end{pmatrix} = \begin{pmatrix} \pi_{1,0} & \pi_{1,1} & \dots & \pi_{1,n-1} \\ \pi_{2,0} & \pi_{2,1} & \dots & \pi_{2,n-1} \\ \dots & \dots & \dots & \dots \\ \pi_{m,0} & \pi_{m,1} & \dots & \pi_{m,n-1} \end{pmatrix},$$

где  $\mathbf{\Pi}$  - матрица потерь (лет жизни). Далее с использованием матрицы  $\mathbf{Y}$ , представляющей структуру реального населения в пересчете на 10000 населения, находим значения реальных потерь.

На главных диагоналях матриц  $\mathbf{Y}\mathbf{\Pi}^T$  и  $\mathbf{Y}^T\mathbf{\Pi}$ , представлены разложения будущих суммарных потерь, рассредоточенных соответственно по возрастам и состояниям здоровья:

$$\mathbf{Y}\mathbf{\Pi}^T = \begin{pmatrix} \lambda_{1\bullet} & \dots & \dots & \dots \\ \dots & \lambda_{2\bullet} & \dots & \dots \\ \dots & \dots & \ddots & \dots \\ \dots & \dots & \dots & \lambda_{m\bullet} \end{pmatrix}; \quad \mathbf{Y}^T\mathbf{\Pi} = \begin{pmatrix} \omega_{\bullet 0} & \dots & \dots & \dots \\ \dots & \omega_{\bullet 1} & \dots & \dots \\ \dots & \dots & \ddots & \dots \\ \dots & \dots & \dots & \omega_{\bullet n-1} \end{pmatrix}.$$

Будущие общие потери равняются сумме  $\sum_{i=1}^m \lambda_{i\bullet}$  (или равной ей сумме  $\sum_{j=0}^{n-1} \omega_{\bullet j}$ ).

Вновь введем ряд показателей типа средних:

$$\bar{\pi}_{i\bullet} = \frac{\lambda_{i\bullet}}{\sum_{j=0}^{n-1} y_{ij}} - \text{средние ожидаемые потери оставшейся жизни населения,}$$

находящегося в  $i$ -м возрастном интервале (вне зависимости от состояния здоровья);

$$\bar{\pi}_{\bullet j} = \frac{\omega_{\bullet j}}{\sum_{i=0}^m y_{ij}} - \text{средние ожидаемые потери будущей жизни населения, на-}$$

ходящегося в  $j$ -м состоянии здоровья (вне зависимости от возраста);

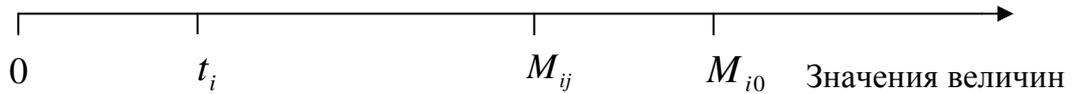
$$\bar{\pi} = \frac{1}{10000} \sum_{i=1}^m \lambda_{i\bullet} - \text{средние ожидаемые потери оставшейся жизни для все-}$$

го населения.

Для характеристики влияния на среднюю продолжительность жизни состояния здоровья  $E_j$  в возрасте  $t_i$  введем относительный показатель

$$u_{ij} = \frac{\pi_{ij}}{q_{i0}} = \frac{M_{i0} - M_{ij}}{M_{i0} - t_i}.$$

Данный показатель назовем индексом потерь. Индекс потерь представляет собой отношение количества лет, не дожитых вследствие наличия заболевания из  $E_j$  на настоящий момент времени, к количеству оставшихся лет жизни при условии отсутствия этих заболеваний в настоящий момент времени (см. рис. 7.5).



**Рис.7.5.** Иллюстрация к структуре показателя  $u_{ij}$ .

Как легко заметить,  $u_{ij}$  не имеет размерности,  $0 \leq u_{ij} \leq 1$ ,  $u_{i0} = 0$ . Чем большее число лет отнимает состояние здоровья  $E_j$  по отношению к состоянию  $E_0$  (относительно здоров) в конкретном возрасте  $t_i$ , тем больше значение  $u_{ij}$ . Следовательно, показатель  $u_{ij}$  является характеристикой (индексом) состояния  $E_j$  применительно к возрасту  $t_i$ . Значения индексов  $u_{ij}$  для всех возможных пар  $(t_i, E_j)$  удобно оформить в виде матрицы

$$\mathbf{U} = \begin{pmatrix} u_{1,0} & u_{1,1} & \dots & u_{1,n-1} \\ u_{2,0} & u_{2,1} & \dots & u_{2,n-1} \\ \dots & \dots & \dots & \dots \\ u_{m,0} & u_{m,1} & \dots & u_{m,n-1} \end{pmatrix}$$

Назовем матрицу  $\mathbf{U}$  индексной матрицей состояний.

Аналогично индексу  $u_{ij}$ , в основе которого лежат потерянные годы жизни, можно ввести индекс  $v_{ij}$ , по смыслу противоположный индексу  $u_{ij}$ , а именно:

$$v_{ij} = \frac{q_{i0} - \pi_{ij}}{q_{i0}} = \frac{M_{ij} - t_i}{M_{i0} - t_i} = 1 - u_{ij}.$$

Указанный индекс, так же как и  $u_{ij}$ , является характеристикой состояния  $E_j$  применительно к возрасту  $t_i$ , но представляет долю оставшихся лет жизни, а не потерянных вследствие  $E_j$  (рис.7.5).

Анализируя реальные значения средней продолжительности жизни  $M_{ij}$ , можно заметить, что в каждом столбце матрицы  $\mathbf{M}$  сверху вниз числа, как правило, возрастают. Т.е. ожидаемая средняя продолжительность жизни, рас-

считываемая при конкретном возрасте, для более старших возрастов выше, что естественно, так как смертность в младших возрастах при расчете уже не учитывается. Исходя из этого, для каждого состояния  $E_j$  и каждого  $i$ -го возрастного интервала введем показатели прироста средней продолжительности жизни:

$\rho_{ij} = M_{ij} - M_{1j}$  ( $i=1,2,\dots,m$ ) - базисный прирост средней продолжительности жизни за счет реально прожитого количества лет;

$\theta_{ij} = M_{ij} - M_{i-1,j}$ , ( $i=2,3,\dots,m$ ) – цепной прирост средней продолжительности жизни за счет реально прожитого количества лет.

Базисные и цепные приросты являются показателями, характеризующими влияние возраста на степень тяжести состояния  $E_j$ . Множества базисных и цепных приростов удобно записать в виде соответствующих матриц приростов.

В качестве относительного показателя влияния возраста на степень тяжести состояния  $E_j$  можно предложить отношение продолжительности будущей жизни при одном и том же состоянии здоровья в разных возрастных интервалах, например,

$$\varphi_{ij}^{(1)} = \frac{M_{ij} - t_i}{M_{1j} - t_1} \quad \text{или} \quad \varphi_{ij}^{(i-1)} = \frac{M_{ij} - t_i}{M_{i-1j} - t_{i-1}}.$$

Каждый из этих показателей, как правило, не превосходит 1. Сравнение однотипных показателей производится для разных состояний здоровья, а также разных групп населения (в частности, мужчин и женщин).

В заключение отметим, что введенные выше показатели имеют очевидную прикладную направленность и могут быть существенно задействованы при изучении общественного здоровья [128].

## ПРИЛОЖЕНИЯ

Таблица П.1. Значения функции Лапласа  $\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt$

x	Φ(x)												
0,00	0,5000	0,41	0,6591	0,82	0,7939	1,23	0,8907	1,64	0,9495	2,10	0,9821	2,92	0,9982
0,01	0,5040	0,42	0,6628	0,83	0,7967	1,24	0,8925	1,65	0,9505	2,12	0,9830	2,94	0,9984
0,02	0,5080	0,43	0,6664	0,84	0,7995	1,25	0,8944	1,66	0,9515	2,14	0,9838	2,96	0,9985
0,03	0,5120	0,44	0,6700	0,85	0,8023	1,26	0,8962	1,67	0,9525	2,16	0,9846	2,98	0,9986
0,04	0,5160	0,45	0,6736	0,86	0,8051	1,27	0,8980	1,68	0,9535	2,18	0,9854	3,00	0,99865
0,05	0,5199	0,46	0,6772	0,87	0,8078	1,28	0,8997	1,69	0,9545	2,20	0,9861	3,20	0,99931
0,06	0,5239	0,47	0,6808	0,88	0,8106	1,29	0,9015	1,70	0,9554	2,22	0,9868	3,40	0,99966
0,07	0,5279	0,48	0,6844	0,89	0,8133	1,30	0,9032	1,71	0,9564	2,24	0,9875	3,60	0,999841
0,08	0,5319	0,49	0,6879	0,90	0,8159	1,31	0,9049	1,72	0,9573	2,26	0,9881	3,80	0,999928
0,09	0,5359	0,50	0,6915	0,91	0,8186	1,32	0,9066	1,73	0,9582	2,28	0,9887	4,00	0,999968
0,10	0,5398	0,51	0,6950	0,92	0,8212	1,33	0,9082	1,74	0,9591	2,30	0,9893	4,50	0,999997
0,11	0,5438	0,52	0,6985	0,93	0,8238	1,34	0,9099	1,75	0,9599	2,32	0,9898	5,00	0,999997
0,12	0,5478	0,53	0,7019	0,94	0,8264	1,35	0,9115	1,76	0,9608	2,34	0,9904		
0,13	0,5517	0,54	0,7054	0,95	0,8289	1,36	0,9131	1,77	0,9616	2,36	0,9909		
0,14	0,5557	0,55	0,7088	0,96	0,8315	1,37	0,9147	1,78	0,9625	2,38	0,9913		
0,15	0,5596	0,56	0,7123	0,97	0,8340	1,38	0,9162	1,79	0,9633	2,40	0,9918		
0,16	0,5636	0,57	0,7157	0,98	0,8365	1,39	0,9177	1,80	0,9641	2,42	0,9922		
0,17	0,5675	0,58	0,7190	0,99	0,8389	1,40	0,9192	1,81	0,9649	2,44	0,9927		
0,18	0,5714	0,59	0,7224	1,00	0,8413	1,41	0,9207	1,82	0,9656	2,46	0,9931		
0,19	0,5753	0,60	0,7257	1,01	0,8438	1,42	0,9222	1,83	0,9664	2,48	0,9934		
0,20	0,5793	0,61	0,7291	1,02	0,8461	1,43	0,9236	1,84	0,9671	2,50	0,9938		
0,21	0,5832	0,62	0,7324	1,03	0,8485	1,44	0,9251	1,85	0,9678	2,52	0,9941		
0,22	0,5871	0,63	0,7357	1,04	0,8508	1,45	0,9265	1,86	0,9686	2,54	0,9945		
0,23	0,5910	0,64	0,7389	1,05	0,8531	1,46	0,9279	1,87	0,9693	2,56	0,9948		
0,24	0,5948	0,65	0,7422	1,06	0,8554	1,47	0,9292	1,88	0,9699	2,58	0,9951		
0,25	0,5987	0,66	0,7454	1,07	0,8577	1,48	0,9306	1,89	0,9706	2,60	0,9953		
0,26	0,6026	0,67	0,7486	1,08	0,8599	1,49	0,9319	1,90	0,9713	2,62	0,9956		
0,27	0,6064	0,68	0,7517	1,09	0,8621	1,50	0,9332	1,91	0,9719	2,64	0,9959		
0,28	0,6103	0,69	0,7549	1,10	0,8643	1,51	0,9345	1,92	0,9726	2,66	0,9961		
0,29	0,6141	0,70	0,7580	1,11	0,8665	1,52	0,9357	1,93	0,9732	2,68	0,9963		
0,30	0,6179	0,71	0,7611	1,12	0,8686	1,53	0,9370	1,94	0,9738	2,70	0,9965		
0,31	0,6217	0,72	0,7642	1,13	0,8708	1,54	0,9382	1,95	0,9744	2,72	0,9967		
0,32	0,6255	0,73	0,7673	1,14	0,8729	1,55	0,9394	1,96	0,9750	2,74	0,9969		
0,33	0,6293	0,74	0,7703	1,15	0,8749	1,56	0,9406	1,97	0,9756	2,76	0,9971		
0,34	0,6331	0,75	0,7734	1,16	0,8770	1,57	0,9418	1,98	0,9761	2,78	0,9973		
0,35	0,6368	0,76	0,7764	1,17	0,8790	1,58	0,9429	1,99	0,9767	2,80	0,9974		
0,36	0,6406	0,77	0,7794	1,18	0,8810	1,59	0,9441	2,00	0,9772	2,82	0,9976		
0,37	0,6443	0,78	0,7823	1,19	0,8830	1,60	0,9452	2,02	0,9783	2,84	0,9977		
0,38	0,6480	0,79	0,7852	1,20	0,8849	1,61	0,9463	2,04	0,9793	2,86	0,9979		
0,39	0,6517	0,80	0,7881	1,21	0,8869	1,62	0,9474	2,06	0,9803	2,88	0,9980		
0,40	0,6554	0,81	0,7910	1,22	0,8883	1,63	0,9484	2,08	0,9812	2,90	0,9981		

Таблица П.2. Значения  $t_\gamma$  при распределении Стьюдента, удовлетворяю-

щие условию  $2 \int_0^{t_\gamma} s_{n-1}(t) dt = \gamma$  в зависимости от  $\gamma$  и  $n-1$

$\gamma$ $n-1$	0,1	0,2	0,3	0,4	0,5	0,6	0,7	0,8	0,9	0,95	0,98	0,99	0,999	$\gamma$ $n-1$
1	0,158	0,325	0,510	0,727	1,000	1,376	1,963	3,03	6,31	12,71	31,8	63,7	636,6	1
2	0,142	0,289	0,445	0,617	0,816	1,061	1,336	1,886	2,92	4,30	6,96	9,92	31,6	2
3	0,137	0,277	0,424	0,584	0,765	0,978	1,250	1,638	2,35	3,18	4,54	5,84	12,94	3
4	0,134	0,271	0,414	0,569	0,741	0,941	1,190	1,533	2,13	2,77	3,75	4,60	8,61	4
5	0,132	0,267	0,408	0,559	0,727	0,920	1,156	1,476	2,02	2,57	3,36	4,03	6,86	5
6	0,131	0,265	0,404	0,553	0,718	0,906	1,134	1,440	1,943	2,45	3,14	3,71	5,96	6
7	0,130	0,263	0,402	0,549	0,711	0,896	1,119	1,415	1,895	2,36	3,00	3,50	5,40	7
8	0,130	0,262	0,399	0,546	0,706	0,889	1,108	1,397	1,860	2,31	2,90	3,36	5,04	8
9	0,129	0,261	0,398	0,543	0,703	0,883	1,100	1,383	1,833	2,26	2,82	3,25	4,78	9
10	0,129	0,260	0,397	0,542	0,700	0,879	1,093	1,372	1,812	2,23	2,76	3,17	4,59	10
11	0,129	0,260	0,396	0,540	0,697	0,876	1,088	1,363	1,796	2,20	2,72	3,11	4,44	11
12	0,128	0,259	0,395	0,539	0,695	0,873	1,083	1,356	1,782	2,18	2,68	3,06	4,32	12
13	0,128	0,259	0,394	0,538	0,694	0,870	1,079	1,350	1,771	2,16	2,65	3,01	4,22	13
14	0,128	0,258	0,393	0,537	0,692	0,868	1,076	1,345	1,761	2,14	2,62	2,98	4,14	14
15	0,128	0,258	0,393	0,536	0,691	0,866	1,074	1,341	1,753	2,13	2,60	2,95	4,07	15
16	0,128	0,258	0,392	0,535	0,690	0,865	1,071	1,337	1,746	2,12	2,58	2,92	4,02	16
17	0,128	0,257	0,392	0,534	0,689	0,863	1,069	1,333	1,740	2,11	2,57	2,90	3,96	17
18	0,127	0,257	0,392	0,534	0,688	0,862	1,067	1,330	1,734	2,10	2,55	2,88	3,92	18
19	0,127	0,257	0,391	0,533	0,688	0,861	1,066	1,328	1,729	2,09	2,54	2,86	3,88	19
20	0,127	0,257	0,391	0,533	0,687	0,860	1,064	1,325	1,725	2,09	2,53	2,84	3,85	20
21	0,127	0,257	0,391	0,532	0,686	0,859	1,063	1,323	1,721	2,08	2,52	2,83	3,82	21
22	0,127	0,256	0,390	0,532	0,686	0,858	1,061	1,321	1,717	2,07	2,51	2,82	3,79	22
23	0,127	0,256	0,390	0,532	0,685	0,858	1,060	1,319	1,714	2,07	2,50	2,81	3,77	23
24	0,127	0,256	0,390	0,531	0,685	0,857	1,059	1,318	1,711	2,06	2,49	2,80	3,74	24
25	0,127	0,256	0,390	0,531	0,684	0,856	1,058	1,316	1,708	2,06	2,48	2,79	3,72	25
26	0,127	0,256	0,390	0,531	0,684	0,856	1,058	1,315	1,706	2,06	2,48	2,78	3,71	26
27	0,127	0,256	0,389	0,531	0,684	0,855	1,057	1,314	1,703	2,05	2,47	2,77	3,60	27
28	0,127	0,256	0,389	0,530	0,683	0,855	1,056	1,313	1,701	2,05	2,47	2,76	3,67	28
29	0,127	0,256	0,389	0,530	0,683	0,854	1,055	1,311	1,699	2,04	2,46	2,76	3,66	29
30	0,127	0,256	0,389	0,530	0,683	0,854	1,055	1,310	1,697	2,04	2,46	2,75	8,65	30
40	0,126	0,255	0,388	0,529	0,681	0,851	1,050	1,303	1,684	2,02	2,42	2,70	3,55	40
60	0,126	0,254	0,387	0,527	0,679	0,848	1,046	1,296	1,671	2,00	2,39	2,66	3,46	60
120	0,126	0,254	0,386	0,526	0,677	0,845	1,041	1,289	1,658	1,980	2,36	2,62	3,37	120
$\infty$	0,126	0,253	0,385	0,524	0,674	0,842	1,036	1,282	1,645	1,960	2,33	2,58	3,29	$\infty$
$n-1$ $\gamma$	0,1	0,2	0,3	0,4	0,5	0,6	0,7	0,8	0,9	0,95	0,98	0,99	0,999	$n-1$ $\gamma$

Таблица П.3. Критические значения распределения Стьюдента  
в зависимости от числа степеней свободы  $\nu$  и уровня значимости  $\alpha$

$\nu$	Уровень значимости $\alpha$ (двусторонняя критическая область)								
	0,50	0,20	0,10	0,05	0,02	0,01	0,005	0,002	0,001
1	1,000	3,078	6,314	12,706	31,821	63,657	127,321	318,309	636,619
2	0,816	1,886	2,920	4,303	6,965	9,925	14,089	22,327	31,599
3	0,765	1,638	2,353	3,182	4,541	5,841	7,453	10,215	12,924
4	0,741	1,533	2,132	2,776	3,747	4,604	5,598	7,173	8,610
5	0,727	1,476	2,015	2,571	3,365	4,032	4,773	5,893	6,869
6	0,718	1,440	1,943	2,447	3,143	3,707	4,317	5,208	5,959
7	0,711	1,415	1,895	2,365	2,998	3,499	4,029	4,785	5,408
8	0,706	1,397	1,860	2,306	2,896	3,355	3,833	4,501	5,041
9	0,703	1,383	1,833	2,262	2,821	3,250	3,690	4,297	4,781
10	0,700	1,372	1,812	2,228	2,764	3,169	3,581	4,144	4,587
11	0,697	1,363	1,796	2,201	2,718	3,106	3,497	4,025	4,437
12	0,695	1,356	1,782	2,179	2,681	3,055	3,428	3,930	4,318
13	0,694	1,350	1,771	2,160	2,650	3,012	3,372	3,852	4,221
14	0,692	1,345	1,761	2,145	2,624	2,977	3,326	3,787	4,140
15	0,691	1,341	1,753	2,131	2,602	2,947	3,286	3,733	4,073
16	0,690	1,337	1,746	2,120	2,583	2,921	3,252	3,686	4,015
17	0,689	1,333	1,740	2,110	2,567	2,898	3,222	3,646	3,965
18	0,688	1,330	1,734	2,101	2,552	2,878	3,197	3,610	3,922
19	0,688	1,328	1,729	2,093	2,539	2,861	3,174	3,579	3,883
20	0,687	1,325	1,725	2,086	2,528	2,845	3,153	3,552	3,850
21	0,686	1,323	1,721	2,080	2,518	2,831	3,135	3,527	3,819
22	0,686	1,321	1,717	2,074	2,508	2,819	3,119	3,505	3,792
23	0,685	1,319	1,714	2,069	2,500	2,807	3,104	3,485	3,768
24	0,685	1,318	1,711	2,064	2,492	2,797	3,091	3,467	3,745
25	0,684	1,316	1,708	2,060	2,485	2,787	3,078	3,450	3,725
26	0,684	1,315	1,706	2,056	2,479	2,779	3,067	3,435	3,707
27	0,684	1,314	1,703	2,052	2,473	2,771	3,057	3,421	3,690
28	0,683	1,313	1,701	2,048	2,467	2,763	3,047	3,408	3,674
29	0,683	1,311	1,699	2,045	2,462	2,756	3,038	3,396	3,659
30	0,683	1,310	1,697	2,042	2,457	2,750	3,030	3,385	3,646
31	0,682	1,309	1,696	2,040	2,453	2,744	3,022	3,375	3,633
32	0,682	1,309	1,694	2,037	2,449	2,738	3,015	3,365	3,622
33	0,682	1,308	1,692	2,035	2,445	2,733	3,008	3,356	3,611
34	0,682	1,307	1,691	2,032	2,441	2,728	3,002	3,348	3,601
35	0,682	1,306	1,690	2,030	2,438	2,724	2,996	3,340	3,591
36	0,681	1,306	1,688	2,028	2,434	2,719	2,990	3,333	3,582
37	0,681	1,305	1,687	2,026	2,431	2,715	2,985	3,326	3,574
$\nu$	0,25	0,10	0,05	0,025	0,01	0,005	0,0025	0,001	0,0005
	Уровень значимости $\alpha$ (односторонняя критическая область)								

v	Уровень значимости $\alpha$ (двусторонняя критическая область)								
	0,50	0,20	0,10	0,05	0,02	0,01	0,005	0,002	0,001
<b>38</b>	0,681	1,304	1,686	2,024	2,429	2,712	2,980	3,319	3,566
<b>39</b>	0,681	1,304	1,685	2,023	2,426	2,708	2,976	3,313	3,558
<b>40</b>	0,681	1,303	1,684	2,021	2,423	2,704	2,971	3,307	3,551
<b>42</b>	0,680	1,302	1,682	2,018	2,418	2,698	2,963	3,296	3,538
<b>44</b>	0,680	1,301	1,680	2,015	2,414	2,692	2,956	3,286	3,526
<b>46</b>	0,680	1,300	1,679	2,013	2,410	2,687	2,949	3,277	3,515
<b>48</b>	0,680	1,299	1,677	2,011	2,407	2,682	2,943	3,269	3,505
<b>50</b>	0,679	1,299	1,676	2,009	2,403	2,678	2,937	3,261	3,496
<b>52</b>	0,679	1,298	1,675	2,007	2,400	2,674	2,932	3,255	3,488
<b>54</b>	0,679	1,297	1,674	2,005	2,397	2,670	2,927	3,248	3,480
<b>56</b>	0,679	1,297	1,673	2,003	2,395	2,667	2,923	3,242	3,473
<b>58</b>	0,679	1,296	1,672	2,002	2,392	2,663	2,918	3,237	3,466
<b>60</b>	0,679	1,296	1,671	2,000	2,390	2,660	2,915	3,232	3,460
<b>62</b>	0,678	1,295	1,670	1,999	2,388	2,657	2,911	3,227	3,454
<b>64</b>	0,678	1,295	1,669	1,998	2,386	2,655	2,908	3,223	3,449
<b>66</b>	0,678	1,295	1,668	1,997	2,384	2,652	2,904	3,218	3,444
<b>68</b>	0,678	1,294	1,668	1,995	2,382	2,650	2,902	3,214	3,439
<b>70</b>	0,678	1,294	1,667	1,994	2,381	2,648	2,899	3,211	3,435
<b>72</b>	0,678	1,293	1,666	1,993	2,379	2,646	2,896	3,207	3,431
<b>74</b>	0,678	1,293	1,666	1,993	2,378	2,644	2,894	3,204	3,427
<b>76</b>	0,678	1,293	1,665	1,992	2,376	2,642	2,891	3,201	3,423
<b>78</b>	0,678	1,292	1,665	1,991	2,375	2,640	2,889	3,198	3,420
<b>80</b>	0,678	1,292	1,664	1,990	2,374	2,639	2,887	3,195	3,416
<b>90</b>	0,677	1,291	1,662	1,987	2,368	2,632	2,878	3,183	3,402
<b>100</b>	0,677	1,290	1,660	1,984	2,364	2,626	2,871	3,174	3,390
<b>120</b>	0,677	1,289	1,658	1,980	2,358	2,617	2,860	3,160	3,373
<b>140</b>	0,676	1,288	1,656	1,977	2,353	2,611	2,852	3,149	3,361
<b>160</b>	0,676	1,287	1,654	1,975	2,350	2,607	2,846	3,142	3,352
<b>180</b>	0,676	1,286	1,653	1,973	2,347	2,603	2,842	3,136	3,345
<b>200</b>	0,676	1,286	1,653	1,972	2,345	2,601	2,839	3,131	3,340
$\infty$	0,6745	1,2816	1,6449	1,9600	2,3263	2,5758	2,8070	3,0902	3,2905
v	0,25	0,10	0,05	0,025	0,01	0,005	0,0025	0,001	0,0005
	Уровень значимости $\alpha$ (односторонняя критическая область)								

Таблица П.4. Наибольшие случайные значения модуля коэффициента корреляции

Число степеней свободы $t$	Уровень значимости $\alpha$			
	0,05	0,01	0,0027	0,001
5	0,75	0,87	0,93	0,95
10	0,58	0,71	0,78	0,82
15	0,48	0,61	0,68	0,72
20	0,42	0,53	0,61	0,65
25	0,38	0,49	0,55	0,60
30	0,35	0,45	0,51	0,55
35	0,32	0,42	0,48	0,52
40	0,30	0,39	0,45	0,49
50	0,27	0,35	0,41	0,44
60	0,25	0,33	0,37	0,41
70	0,23	0,30	0,35	0,38
80	0,22	0,28	0,33	0,36
90	0,21	0,26	0,31	0,34
100	0,19	0,25	0,29	0,32
120	0,18	0,23	0,27	0,30
150	0,16	0,21	0,24	0,26
200	0,14	0,18	0,21	0,23
300	0,11	0,15	0,17	0,19
400	0,10	0,13	0,15	0,16
500	0,09	0,11	0,13	0,15
700	0,07	0,10	0,11	0,12
900	0,06	0,09	0,10	0,11
1000 и более	< 0,06	< 0,09	< 0,10	< 0,11

Таблица П.5. Критические точки  $F$  распределения Фишера в зависимости от числа степеней свободы большей дисперсии  $k_1$  и числа степеней свободы меньшей дисперсии  $k_2$

Значения  $F$  при  $\alpha = 0,05$

$k_2 \backslash k_1$	1	2	3	4	5	6	8	12	24	$\infty$
1	161,4	199,5	215,7	224,6	230,2	234,0	238,9	243,9	249,0	254,3
2	18,51	19,00	19,16	19,25	19,30	19,33	19,37	19,41	19,45	19,50
3	10,13	9,55	9,28	9,12	9,01	8,94	8,84	8,74	8,64	8,53
4	7,71	6,94	6,59	6,39	6,26	6,16	6,04	5,91	5,77	5,63
5	6,61	5,79	5,41	5,19	5,05	4,95	4,82	4,68	4,53	4,36
6	5,99	5,14	4,76	4,53	4,39	4,28	4,15	4,00	3,84	3,67
7	5,59	4,74	4,35	4,12	3,97	3,87	3,73	3,57	3,41	3,23
8	5,32	4,46	4,07	3,84	3,69	3,58	3,44	3,28	3,12	2,93
9	5,12	4,26	3,86	3,63	3,48	3,37	3,23	3,07	2,90	2,71
10	4,96	4,10	3,71	3,48	3,33	3,22	3,07	2,91	2,74	2,54
11	4,84	3,98		3,36	3,20	3,09	2,95	2,79	2,61	2,40
12	4,75	3,88	3,59	3,26	3,11	3,00	2,85	2,69	2,50	2,30
13	4,67	3,80	3,49	3,18	3,02	2,92	2,77	2,60	2,42	2,21
14	4,60	3,74	3,41	3,11	2,96	2,85	2,70	2,53	2,35	2,13
15	4,54	3,68	3,34	3,06	2,90	2,79	2,64	2,48	2,29	2,07
16	4,49	3,63	3,29	3,01	2,85	2,74	2,59	2,42	2,24	2,01
17	4,45	3,59	3,24	2,96	2,81	2,70	2,55	2,38	2,19	1,96
18	4,41	3,55	3,20	2,93	2,77	2,66	2,51	2,34	2,15	1,92
19	4,38	3,52	3,16	2,90	2,74	2,63	2,48	2,31	2,11	1,88
20	4,35	3,49	3,13	2,87	2,71	2,60	2,45	2,28	2,08	1,84
21	4,32	3,47	3,10	2,84	2,68	2,57	2,42	2,25	2,05	1,81
22	4,30	3,44	3,07	2,82	2,66	2,55	2,40	2,23	2,03	1,78
23	4,28	3,42	3,05	2,80	2,64	2,53	2,38	2,20	2,00	1,76
24	4,26	3,40	3,03	2,78	2,62	2,51	2,36	2,18	1,98	1,73
25	4,24	3,38	3,01	2,76	2,60	2,49	2,34	2,16	1,96	1,71
26	4,22	3,37	2,99	2,74	2,59	2,47	2,32	2,15	1,95	1,69
27	4,21	3,35	2,98	2,73	2,57	2,46	2,30	2,13	1,93	1,67
28	4,20	3,34	2,96	2,71	2,56	2,44	2,29	2,12	1,91	1,65
29	4,18	3,33	2,95	2,70	2,54	2,43	2,28	2,10	1,90	1,64
30	4,17	3,32	2,93	2,69	2,53	2,42	2,27	2,09	1,89	1,62
40	4,08	3,23	2,92	2,61	2,45	2,34	2,18	2,00	1,79	1,52
60	4,00	3,15	2,84	2,52	2,37	2,25	2,10	1,92	1,70	1,39
120	3,92	3,07	2,76	2,45	2,29	2,17	2,02	1,83	1,61	1,25
$\infty$	3,84	2,99	2,68	2,37	2,21	2,09	1,94	1,75	1,52	1,00

Значения  $F$  при  $\alpha = 0,01$ 

$k_1 \backslash k_2$	1	2	3	4	5	6	8	12	24	$\infty$
1	4052	4999	5403	5625	5764	5859	5981	6106	6234	6366
2	98,49	99,00	99,17	99,25	99,30	99,33	99,36	99,42	99,46	99,50
3	34,12	30,81	29,46	28,71	28,24	27,91	27,49	27,05	26,60	26,12
4	21,20	18,00	16,69	15,98	15,52	15,21	14,80	14,37	13,93	13,46
5	16,26	13,27	12,03	11,39	10,97	10,67	10,29	9,89	9,47	9,02
6	13,74	10,92	9,78	9,15	8,75	8,47	8,10	7,72	7,31	6,88
7	12,25	9,55	8,45	7,85	7,46	7,19	6,84	6,47	6,07	5,65
8	11,26	8,65	7,59	7,01	6,63	6,37	6,03	5,67	5,28	4,86
9	10,56	8,02	6,99	6,42	6,06	5,80	5,47	5,11	4,73	4,31
10	10,04	7,56	6,55	5,99	5,64	5,39	5,03	4,71	4,33	3,91
11	9,65	7,20	6,22	5,67	5,32	5,07	4,74	4,40	4,02	3,60
12	9,33	6,93	5,95	5,41	5,06	4,82	4,50	4,16	3,78	3,36
13	9,07	6,70	5,74	5,20	4,86	4,62	4,30	3,96	3,59	3,16
14	8,86	6,51	5,56	5,03	4,69	4,46	4,14	3,80	3,43	3,00
15	8,68	6,36	5,42	4,89	4,56	4,32	4,00	3,67	3,29	2,87
16	8,53	6,23	5,29	4,77	4,44	4,20	3,89	3,55	3,18	2,75
17	8,40	6,11	5,18	4,67	4,34	4,10	3,79	3,45	3,08	2,65
18	8,28	6,01	5,09	4,58	4,25	4,01	3,71	3,37	3,00	2,57
19	8,18	5,93	5,01	4,50	4,17	3,94	3,63	3,30	2,92	2,49
20	8,10	5,85	4,94	4,43	4,10	3,87	3,56	3,23	2,86	2,42
21	802	5,78	4,87	4,37	4,04	3,81	3,51	3,17	2,80	2,36
22	7,94	5,72	4,82	4,31	3,99	3,76	3,45	3,12	2,75	2,31
23	7,88	5,66	4,76	4,26	3,94	3,71	3,41	3,07	2,70	2,26
24	7,82	5,61	4,72	4,22	3,90	3,67	3,36	3,03	2,66	2,21
25	7,77	5,57	4,68	4,18	3,86	3,63	3,32	2,99	2,62	2,17
26	7,72	5,53	4,64	4,14	3,82	3,59	3,29	2,96	2,58	2,13
27	7,68	5,49	4,60	4,11	3,78	3,56	3,26	2,93	2,55	2,10
28	7,64	5,45	4,57	4,07	3,75	3,53	3,23	2,90	2,52	2,06
29	7,60	5,42	4,54	4,04	3,73	3,50	3,20	2,87	2,49	2,03
30	7,56	5,39	4,51	4,02	3,70	3,47	3,17	2,84	2,47	2,01
40	7,31	5,18	4,31	3,83	3,51	3,29	2,99	2,66	2,29	1,80
60	7,08	4,98	4,13	3,65	3,34	3,12	2,82	2,50	2,12	1,60
120	6,85	4,79	3,95	3,48	3,17	2,96	2,66	2,34	1,95	1,38
$\infty$	6,64	4,60	3,78	3,32	3,02	2,80	2,51	2,18	1,79	1,00

## Литература

1. Айвазян С.А. Интегральные индикаторы качества жизни населения: их построение и использование в социально-экономическом управлении и межрегиональных сопоставлениях.- М.: РАН, Центральный экономико-математический институт, 2000. – 32 с.
2. Айвазян С.А., Мхитарян В.С. Прикладная статистика и основы эконометрики. – М.: ЮНИТИ, 1998. – 40 Мгб.
3. Амосов Н.М. Раздумья о здоровье.- М.: Физкультура и спорт, 1987. – 18 с.
4. Амосова Н.Н., Куклин Б.А., Макарова С.Б., Максимов Ю.Д и др. вероятностные разделы математики / Под ред. Ю.Д. Максимова. – СПб: Иван Федоров, 2001. -592 с.
5. Андерсен Т. Статистический анализ временных рядов. – М.: Мир, 1976. – 736 с.
6. Аникин И.В. Модели нечётких нейронных сетей // Сб.: Эволюционное моделирование. / Под ред. В.А. Райхлина. – Казань: АЭН (Наука), 2004. – С. 111 ÷ 136.
7. Антропология - медицине / Под ред. Т.И.Алексеева. М.: Изд-во МГУ, 1989. 235 с.
8. Афанасьев В.Н., Юзбашев М.М. Анализ временных рядов. – М.: Мир, 2001. –
9. Балантер Б.М., Ханин М.А., Чернавский Д.С. Введение в математическое моделирование патологических процессов. – М.: Медицина, 1980. – 174 с.
10. Бахвалов Н.С. Численные методы. – М.: Наука, 1975. – 632 с.
11. Бачманов А. А. Математические модели интегрального показателя оценки здоровья населения. // Материалы междунар. научно-метод. конф. “Математика в вузе”. – Петрозаводск, 2003. – С. 112 – 113.
12. Бачманов А.А, Прохорова А.В. Информационные ресурсы и потоки в здравоохранении // Материалы научной сессии Новгородского научного центра СЗО РАМН, Том 2. – М.: 2003. – С. 82 ÷ 87.
13. Бачманов А.А., Рязанцев П.П. Некоторые вопросы формирования единой базы данных «Здоровье населения Новгородской области» // Охрана здоровья населения – национальный приоритет государственной политики. – Сб. научных тр. ННЦ СЗО РАМН, Т. 5, 2006. С. 81 – 85.
14. Беллман Р., Заде Л. Принятие решений в расплывчатых условиях. // Вопросы анализа и процедуры принятия решений: Сб. статей / Пер. с англ.; Под ред. И.Ф. Шахнова. – М.: Ил., 1976. – С. 172 ÷ 215.
15. Бокс Дж., Дженкинс Г.Д. Анализ временных рядов. Прогноз и управление. Вып. 1 / Пер. с англ. Под ред. В.Ф. Писаренко. – М.: Мир, 1974. – 408 с.
16. Болотханов Э.Б. К проблеме оценки характеристик нестационарных случайных процессов / Труды междунар. научно-метод. конф. ”Математика в вузе“. – Псков: ППИ, 2001. – С. 146 – 149.
17. Большая советская энциклопедия. Т. 16. 3-е издание. – М.: Советская энциклопедия, 1974. – 616 с.
18. Большев Л.Н., Смирнов Н.В. Таблицы математической статистики – М.: Наука, 1983 – 416 с.
19. Борисов Б.М., Примаков В.И., Мартирова Т.А.. Экологические подходы в оценке состояния подростков // Военно-медицинский журнал.- 1996.- № 2.- С. 51-55.
20. Борисов Ю.П., Цветнов В.В. Математическое моделирование радиотехнических систем и устройств. – М.: Радио и связь, 1985. – 296 с.
21. Боровиков В.П., Ивченко Г.И. Прогнозирование в системе STATISTICA в среде WINDOWS. – М.: Финансы и статистика, 2006. – 368 с.
22. Бриллинджер Д. Временные ряды. – М.: Мир, 1980. – 536 с.
23. Васильев Ф.П. Численные методы решения экстремальных задач. – М.: Наука, 1980. – 519 с.

24. Венедиктов Д.Д. Очерки системной теории и стратегии здравоохранения. – М.: 2008. – 336 с.
25. Вентцель Е.С. Теория вероятностей. – М: Физматлит, 1998. – 576 с.
26. Вентцель Е.С., Овчаров Л.А. Теория вероятностей и её приложения. – М.: Наука, 1988. – 480 с.
27. Вероятностные разделы математики / Под общей ред. Ю.Д. Максимова. – СПб: «Иван Федоров», 2001. – 592 с.
28. Геронтология in silico: становление новой дисциплины : Математические модели, анализ данных и вычислительные эксперименты : сборник науч. тр. / Под ред. Г. И. Марчука, В. Н. Анисимова, А. А. Романюхи, А. И. Яшина. – М.: БИНОМ. Лаборатория знаний, 2007. – 535 с. : ил.
29. Гмурман В.Е. Теория вероятностей и математическая статистика. – М.: Высшая школа, 1998. – 479 с.
30. Голяндина Н.А., Некруткин В.В., Браунов К.А. Главные компоненты временных рядов: метод "Гусеница". – СПб.: СПб-унив-т, 1997. – 7 Кб.
31. Девяткова Г.И. Математическое моделирование в управлении деятельностью хирургической клиники (монография) // Препринт.- Пермь: Изд-во Пермского ун-та. 2000. – 110 с.
32. Девяткова Г.И., Михеев В.В. Моделирование технологий оказания медицинской помощи в хирургическом стационаре на этапе операционной // Стоматология XXI века: Новейшие технологии и материалы: Сб. науч. трудов III Всероссийского симпозиума. – Пермь: Перм. гос. мед. Академия, 2000. – С. 23 ÷ 24.
33. Девяткова Г.И., Радионова М.В., Попов А.В. Использование математических моделей в диагностике причин желтухи у хирургического больного // Здоровье и образование. Медико-социальные и экономические проблемы: Материалы междунар. научно-практич. конф. / Пермь-Париж, 2004. С. 66 ÷ 69.
34. Девяткова Г.И., Суслонов В.М., Радионова М.В. Математическое моделирование синдромов ЖКБ (монография) // Пермь: Изд-во Перм. ун-та, 2005. – 122 с.
35. Дуброва Т.А. Статистические методы прогнозирования. – М.: ЮНИТИ-ДАНА, 2003. – 208 с.
36. Егоренков Д.Л., Фрадков А.Л., Харламов В.Ю. Основы математического моделирования. Построение и анализ моделей с примерами на языке MATLAB. – СПб: БГТУ, 1994. – 192 с.
37. Емалетдинова Л.Ю., Галлиев Ш.И., Разина М.А и др. Использование непрерывных моделей для оптимизации расположения станций скорой помощи // Вычислительная механика и современные прикладные программные системы. – Докл. XII междунар. конф., Владимир, 2003. – С. 266 ÷ 267.
38. Ермаков С.М., Михайлов Г.А. Курс статистического моделирования. – М.: Наука, 1976. – 320 с.
39. Ермаков С.П. Современные возможности интегральной оценки медико-демографических процессов. – М.: Медицина, 1996. – 16 с.
40. Ефимова М.Р., Бычкова С.Г. Социальная статистика. – М.: Финансы и статистика, 2003. – 560 с.
41. Зайцева Н.В., Землянова М.А., Кирьянов Д.А. Определение критических параметров загрязнения атмосферного воздуха по критерию обращаемости за медицинской помощью // Гигиена и санитария. 2002. № 2. С. 18-20.
42. Зинин Н.А. Вариант комплексного измерителя здоровья и методы его расчёта / Тезисы докл. XIV обл. научно-практич. конф. по вопросам курортного лечения. – Самара, 1988. – С. 97 – 99.
43. Ильеня А.М. Моделирование случайных векторов с произвольными распределениями координат // Вестник Новгородского гос. ун-та. Сер.: Математика и информатика. № 22, 2002. – С. 32 ÷ 35.

44. Ильин В.П. Численный анализ. Ч. 1. – Новосибирск: ИВМ И МГ СО РАН, 2004. – 335 с.
45. Кант В.И. Математические методы и моделирование в здравоохранении. – М.: Медицина, 1987. – 224 с.
46. Кендэл М.Д. Временные ряды / Пер. с англ. Под ред. Ю.П. Лукашина. – М.: Финансы и статистика, 1981. – 199 с.
47. Кендэл М., Стьюарт А. Многомерный статистический анализ и временные ряды / Пер. с англ. Под ред. А.Н. Колмогорова и Ю.В. Прохорова. – М.: Наука, 1976. 736 с.
48. Кеслер Г.Дж. Теория моделей / Пер. с англ.; Под ред. С.Г. Гончарова. – М.: Мир, 1977. – 614 с.
49. Кирьянов Б.Ф. К проблеме определения весовых коэффициентов параметров линейных моделей интегральных показателей качества систем // Вестник Новгородского гос-го унив-та, № 44, 2007. – С. 33 – 37.
50. Кирьянов Б.Ф. Простой метод моделирования случайных векторов // Вестник Новгородского гос. ун-та. Сер.: Естеств. и технич. науки. № 13, 1999. – С. 88 ÷ 90.
51. Кирьянов Б.Ф. Исследование динамики интегрального показателя здоровья населения России. -
52. Кирьянов Б.Ф. Методика определения значений параметров моделей интегрального показателя общественного здоровья // Охрана здоровья населения – национальный приоритет государственной политики. – Сб. научных тр. ННЦ СЗО РАМН, Т. 5, 2006. С. 125 – 130.
53. Кирьянов Б.Ф. К теории построения интегральных показателей здоровья населения / Роль медицинской науки и здравоохранения в реализации демографической политики государства. – Сб. научных тр. ННЦ СЗО РАМН, Т. 6, 2007. – С. 198 – 203.
54. Кирьянов Б.Ф. Анализ распределений и моделирование показателей здоровья / Роль медицинской науки и здравоохранения в реализации демографической политики государства. – Сб. научных тр. ННЦ СЗО РАМН, Т. 6, 2007. – С. 203 – 206.
55. Кирьянов Б.Ф. Интегральные показатели качества систем // Труды XX-й международной. научно-методической конференции «Математика в вузе» . – СПб: 2008. – С. 8 – 9.
56. Кирьянов Б.Ф., Бачманов А.А. Прогнозирование показателей здоровья населения // Материалы научной сессии ННЦ СЗО РАМН, Т. 2. - М.: Медицина, 2003. – С. 76 ÷ 82.
57. Кирьянов Б.Ф., Болотханов Э.Б. Усечение распределений в стохастических моделях. Реализация квазинормального распределения/Материалы докладов междунар. научно-практич. конф. “Математическое моделирование в науке, промышленности и образовании”. Тирасполь: РИО ТГУ, 2001. С. 83 – 85.
58. Кирьянов Б.Ф, Кирьянов Д.В. К теории построения интегральных показателей качества систем на основе линейных математических моделей // М.: Современные наукоёмкие технологии, № 4, 2008. – С. 73 – 74.
59. Кирьянов Б.Ф., Кознов А.В. Процессы авторегрессии со случайными коэффициентами и их применение при моделировании радиотехнических систем // Прикладная математика, Под ред. Б.Ф. Кирьянова, Новгород: НПИ, 1994, вып. 1. – С 3÷8.
60. Кирьянов Б.Ф., Лысенко В.А. Интерполирование кусочно-дифференцируемых функций локальными кубическими сплайнами. – Новгород: НовГУ, 1995. – 14 с.
61. Кирьянов Б.Ф., Майоров В.В. Математическая модель ишемической болезни сердца // Вестник Новгородского гос. ун-та. Сер.: Медиц. науки. № 7, 1998. – С. 78 ÷ 80.
62. Кирьянов Б.Ф., Медик В.А. Интегральная оценка здоровья населения по средней продолжительности жизни // Сборник научных трудов Новгородского научного центра СЗО РАМН. Том 4. – М.: Медицина, 2005. – С. 42 ÷ 46.
63. Кирьянов Б.Ф., Медик В.А. Усовершенствованные многопараметрические модели

- интегрального показателя общественного здоровья населения // Охрана здоровья населения – национальный приоритет государственной политики. – Сб. научных тр. НИЦ СЗО РАМН, Т. 5, 2006. С. 67 – 73.
64. Кирьянов Б.Ф., Медик В.А., Токмачёв М.С. Чувствительность интегральных показателей многопараметрических систем // Вестник Новгородского гос. ун-та. Сер.: Технич. науки. № 26, 2004. – С. 114 ÷ 116.
  65. Кирьянов Б.Ф., Петрова Ю.Ю. Прогнозирование временных рядов с особыми значениями // Вестник НовГУ, Сер. «Технич. науки», № 28, 2004. – С. 92 ÷ 96.
  66. Клементьев А.А. Использование методов математического моделирования для изучения общественного здоровья. – М.: Ин-т проблем управления АН СССР, 1989. – 313 с. – Рук. деп. в ВИНТИ 31.07.89, № Д-18154.
  67. Кобринский Б.А. Медико-экологический мониторинг как основа профилактики хронической патологии у детей. // Российский вестник перинатологии и педиатрии. – 1994. - № 5. – С. 2-5.
  68. Котова Т.Е., Бачманов А.А. Анализ уровня и структуры исчерпанной заболеваемости населения Новгородской области // Охрана здоровья населения – национальный приоритет государственной политики. – Сб. научных тр. НИЦ СЗО РАМН, Т. 5, 2006. С. 18 – 24.
  69. Корчагин В.П. Финансовое обеспечение здравоохранения. – М.: Эпидавр, 1997. – 42 с.
  70. Кудрявцев В.А. Краткий курс математического анализа.- М.: Наука, 1989. – 736 с.
  71. Кузин Л.Т. Основы кибернетики. Т.2. Основы кибернетических моделей.– М.: Энергия, 1979. – 584 с.
  72. Лебедев А.Н. Моделирование в научно-технических исследованиях. – М.: Радио и связь, 1989. – 294 с.
  73. Лисицин В.И. Смертность как интегральный показатель здоровья населения // Изучение факторов риска и прогнозирование здоровья населения на региональном уровне. - Сборник научных работ Новгородского научного центра СЗО РАМН. Том 1. – М.: Медицина, 2003. – С.105 ÷ 115.
  74. Лисицын Ю.П. “Атрибуты” здоровья (К вопросу характеристики обусловленности здоровья) // Русский медицинский сервер. – Здравология, 2002. – 5 с.
  75. Лисицын Ю.П. Здоровье населения и современные теории медицины. – М.: Медицина, 1982. – 287 с.
  76. Лисицын Ю.П. “Модус” здоровья россиян // Экономика здравоохранения, № 2, 2001. – с. 32 – 37.
  77. Лисицын Ю.П. Образ жизни и здоровье населения – М., 1982. – 40 с.
  78. Лоу А,М, Кельтон В.Д. Имитационное моделирование. 3. – СПб: ПИТЕР, 2004. – 847 с.
  79. Лукашин Ю.П. Адаптивные методы краткосрочного прогнозирования. – М.: Статистика, 1979. – 254 с.
  80. Марченко А.Г. Групповые оценки здоровья населения при использовании различных источников информации // Итоги комплексного изучения здоровья населения в 1969 – 1971 гг./ Под ред. А.Ф. Серенко. – М.: Медицина. – С. 148 ÷ 150.
  81. Матвеев Н.М. Обыкновенные дифференциальные уравнения. – СПб: Специальная Литература. 1996. – 372 с.
  82. Медик В.А. Заболеваемость населения: история, современное состояние и методология изучения. – М.: Медицина, 2003. – 512 с.
  83. Медик В.А. Показатели комплексной оценки здоровья// Советское здравоохранение”, № 2, 1991. – С. 24 ÷ 29.
  84. Медик В.А., Кирьянов Б.Ф. Подходы к прогнозированию показателей здоровья населения / Проблемы социальной гигиены, здравоохранения и истории медицины. – М.: Медицина, 2005, № 6. – С.3 – 5.

85. Медик В.А., Кирьянов Б.Ф., Бачманов А.А. Линейные модели интегрального показателя оценки здоровья населения // Сборник научных трудов Новгородского научного центра СЗО РАМН. Том 4. – М.: Медицина, 2005. – С. 72 ÷ 78.
86. Медик В.А., Кирьянов Б.Ф., Бачманов А.А., Петрова Ю.Ю. Прогнозирование показателей здоровья населения на основе статистических данных с нетиповыми значениями // Состояние здоровья населения и методология его изучения. – Сборник научных трудов ННЦ СЗО РАМН, Т. 3, 2004. – С. 96 ÷ 101.
87. Медик В.А., Кирьянов Б.Ф., Токмачёв М.С., Бачманов А.А. Моделирование интегральных показателей для оценки общественного здоровья // Сборник научных трудов Новгородского научного центра СЗО РАМН. Том 1. – М.: Медицина, 2003. – С. 90 ÷ 95.
88. Медик В. А., Кирьянов Б. Ф., Токмачёв М. С., Бачманов А. А. Моделирование интегральных показателей для комплексной оценки здоровья населения. // Материалы XI междунар. конф. и дискус-го научного клуба “Новые информационные технологии в медицине, биологии, фармакологии и экологии”, Украина, Крым, Ялта-Гурзуф, 2003. – С. 181 – 183.
89. Медик В.А., Токмачев М.С. Математическая статистика в медицине и биологии. – Новгород: НовГУ, 1998. – 417 с.
90. Медик В.А., Токмачев М.С. Моделирование интегральных показателей оценки здоровья населения // Здравоохранение РФ, № 3, 2003. – С. 17 ÷ 20.
91. Медик В.А., Токмачёв М.С. Математическая статистика в медицине: учебное пособие. – М.: Финансы и статистика, 2007. – 800 с.
92. Медик В.А., Токмачев М.С. Соотношения параметров физического развития детей. - Сб. науч. тр. Новгородского научного центра Северо-Западного отделения РАМН - М.: Медицина, 2005.- Т.4.- С. 78 – 83.
93. Медик В.А., Токмачёв М. Руководство по статистике здоровья и здравоохранения. – М.: Медицина, 2006. – 528 с
94. Медик В.А., Токмачев М.С., Фишман Б.Б. Статистика в медицине и биологии. Том 1. Теоретическая статистика. – М.: Медицина, 2000. – 456 с.
95. Медик В.А., Токмачев М.С., Фишман Б.Б. Статистика в медицине и биологии. Том 2. Прикладная статистика здоровья – М.: Медицина, 2001. – 352 с.
96. Медик В.В., Швецов А.Г., Бачманов А.А. Современные подходы к определению статуса здоровья индивида // Здоровье населения и приоритеты здравоохранения. – Сборник научных трудов ННЦ СЗО РАМН, Т. 4, 2005. – С. 92 – 97.
97. Методические указания по изучению здоровья населения / О.П. Щепин, В.А. Медик, В.И. Стародубов и др. Утв. МЗ РФ и РАМН 15.12.2005 г. – М., 2005. – 71 с.
98. Моделиране и симулиране на човешката памет «МНЕМО-89» // Материалы междунар. научно-технич. конф. «Моделирование человеческой памяти». – Болгария, Стара Загора, 1989. – 136 с.
99. Нейман Ю. Вводный курс теории вероятностей и математической статистики // Пер. с англ.; Под ред. Ю.В. Линника. – М.: Наука, 1968. – 448 с.
100. Оконенко Т.И., Токмачев М.С., Вебер В.Р. Экологические подходы к оценке влияния загрязнения атмосферного воздуха на детей с заболеваниями дыхательной системы. - Экология человека. № 4. 2006. С. 6-9.
101. Основные показатели состояния здоровья населения и деятельности организаций здравоохранения Новгородской области за 2005 год. – Вел. Новгород: Облздрав, ГУЗ «Мед. информ.-аналит. центр», 2006. – 50 с.
102. Петленко В.П. Интегральная медицина XXI века // РИА Медицина, № 7, 1988. – 9 с.
103. Полляк Ю.Г. Вероятностное моделирование на ЭВМ. – М.: Сов. радио, 1971.- 400 с.
104. Поляков И.В., Петрова Н.Г. Комплексная характеристика качества диагностики и лечения тяжелых болезней // Сов. здравоохранение, № 11, 1985. – С. 32 ÷ 34.
105. Попов Л.А. Анализ и моделирование трудовых показателей. Изд. 2-е, доп. и пере-

- раб. – М.: Финансы и статистика, 1999. – 208 с.
106. Прицкер А. Введение в имитационное моделирование и язык СЛАМ II./ Пер. с англ. под ред. А.Д. Цвиркуна и В.А. Филиппова. – М.: Мир, 1987. – 646 с
  107. Пугачёв В.С. Теория вероятностей и математическая статистика. – М.: Наука, 1979. – 496 с.
  108. Пугачёв В.С., Сеницын И.Н. Теория стохастических систем. – М.: Логос, 2000. – 1000 с.
  109. Робертс Ф. Дискретные модели и их приложения в технике, биологии, экологии. – М.: Наука, 1986. – 494 с
  110. Рязанцев П.П. Основы программного обеспечения для статистического анализа комплексного медицинского обследования населения // Роль медицинской науки и здравоохранения в реализации демографической политики государства. – Сб. научных тр. ННЦ СЗО РАМН, Т. 7, 2007. – С. 211 – 214.
  111. Рязанцев П.П., Токмачёв М.С. Разработка программного комплекса для расчёта новых показателей здоровья // Роль медицинской науки и здравоохранения в реализации демографической политики государства. – Сб. научных тр. ННЦ СЗО РАМН, Т. 7, 2007. – С. 214 – 219.
  112. Славин М.Б. Марковская модель заболеваемости, выживаемости, повторения заболеваемости и смертности. ИСА РАН. – М.: Депонир. в ВИНТИ 03.08.92, № 2525 – В92.
  113. Славин М.Б. Практика системного моделирования в медицине. – М.: Медицина, 2002. – 168 с.
  114. Славин М.Б. Системное моделирование патологических процессов.– М: Медицина, 1983. – 182 с. // число стр. – приближённое.
  115. Советов Б.Я., Яковлев С.А. Моделирование систем. – М.: ВШ, 1985. – 271 с.
  116. Соломонов А.И., Вялков А.И. Мониторинг здоровья населения как основа развития здравоохранения. – М.: ГЭОТАР Медицина, 1998. – 38 с.
  117. Токмачев М.С. Цепи Маркова в прогнозировании медико-социальных показателей // Обозрение прикладной и промышленной математики. - Т.10. - Вып.2. - М., 2003. – С. 517 – 518; Токмачев М. С.
  118. Токмачев М. С. Регрессионные модели заболеваемости и смертности населения. - Сб. науч. тр. Новгородского научного центра Северо-Западного отделения РАМН; М.: Медицина, 2004. - Т.3.- С. 151 – 155.
  119. Токмачев М. С. Разработка ряда показателей общественного здоровья на основе цепей Маркова. Приложение к: Вестник НовГУ. Сер.: Технич. науки. № 28. 2004. Препринт. –5 с.
  120. Токмачёв М.С. Временные ряды и прогнозирование. – Великий Новгород: НовГУ, 2005. – 192 с.
  121. Токмачев М. С. Разработка новых показателей общественного здоровья на основе статистических данных. Сб. науч. тр. Новгородского научного центра Северо-Западного отделения РАМН - М.: Медицина, 2005.- Т.4.- С. 119-127.
  122. Токмачев М.С. Вычисление площади поверхности тела человека. - Сб. науч. тр. Новгородского научного центра Северо-Западного отделения РАМН - М.: Медицина, 2005.- Т.4.– С. 127 – 132.
  123. Токмачев М. С. Разработка математико-статистических методов оценки здоровья населения.- Сб. материалов Всероссийской научно-практической конференции и трудов Новгородского научного центра Северо-Западного отделения РАМН. – М.: Медицина, Т. 6. 2007.– С. 219 – 229.
  124. Токмачёв М.С. Разработка математико-статистических методов оценки здоровья населения / Роль медицинской науки и здравоохранения в реализации демографической политики государства. – Сб. научных тр. ННЦ СЗО РАМН, Т. 7, 2007. – С.

- 219 – 229.
125. Токмачев М. С. Математическая модель процесса здоровья населения региона. Четвертая международная конференция по проблемам управления (26-30 января 2009г.): Сб. трудов. – М: Учреждение Российской академии наук Институт проблем управления им. В.А. Трапезникова РАН, 2009. – 2030 с. (С. 893– 906).
  126. Токмачев М. С., Бачманов А.А. Особенности построения алгоритма формирования стохастических матриц / Сб. науч. тр. Новгородского научного центра Северо-Западного отделения РАМН. – М.: Медицина, 2004. – Т.3.- С. 156 – 161.
  127. Токмачев М. С., Рязанцев П.П. Разработка программного комплекса для расчета новых показателей здоровья населения. – Сб. материалов Всероссийской научно-практической конференции и трудов Новгородского научного центра Северо-Западного отделения РАМН. М.: Медицина, 2007. – Т.6. С. 219 – 229
  128. Токмачёв М. С., Фишман Б. Б. Оценка математической зависимости «пылевой фактор – заболеваемость верхних дыхательных путей и лёгких рабочих» на производстве огнеупоров // Медицина труда и промышленная экология. №7. 2003. – С. 38-43.
  129. Тюрин Ю.Н., Макаров А.А. Статистический анализ данных на компьютере. – М.: ИН-ФРА, 1998. – 528 с.
  130. Феллер В. Введение в теорию вероятностей и ее приложения. – М.: “Мир”, Т. 2, 1984. – 752 с.
  131. Хорафас Д.Н. Системы и моделирование / Пер. с англ.; Под ред. И.Н. Коваленко. – М.: Мир, 1967. – 420 с.
  132. Шалыгин А.С., Палагин Ю.И. Прикладные методы статистического моделирования. – Л.: Машиностроение, 1986. – 320 с.
  133. Швецов А.Г. Авторский взгляд на понятия “здоровье” и “физическое состояние” индивида // Здоровье населения и приоритеты здравоохранения. –Сборник научных трудов ННЦ СЗО РАМН, Т. 4, 2005. – С. 145 – 150.
  134. Швецов А.Г., Калишев М.Г., Приз В.Н., Кабиева С.М. Новый подход к классификации и оценке уровней здоровья человека // Клиническая медицина.: вопросы клини-ки, диагностики, профилактики и лечения. – Межвуз. сб. стран СНГ, Т. 7. – Великий Новгород, Алматы, 2001. – С. 33 – 41.
  135. Шеннон Р. Имитационное моделирование систем – искусство и наука / Пер. с англ.; Под. ред. Е.К. Масловского. – М.: Мир, 1978. – 420 с.
  136. Шор Я.Б. Статистические методы анализа и контроля качества и надёжности. – М.: Советское ради, 1962. – 552 с.
  137. Щепин О.П., Купеева И.А., Щепин В.О. Какорина Е.П. Современные региональные особенности здоровья населения и здравоохранения России. – М.: Медицина, Шико, 2007. – 360 с.
  138. Banks J. Interpreting Simulation Software Checklists. – OR/MS Today, 1996, Num. 23. – P. 74 ÷ 78.
  139. Biles W.E. Experimental Design in Computer Simulation. / Proc. Winter Simulation Conference, San Diego, 1979. – P. 3 ÷ 9.
  140. Boyd E. Experimental error inherent in measuring growing human body. Am J Physiol. 1930.
  141. Chatfield C. The Analysis of Time Series: 4th ed. – Chapman and Hall, 1989 – 242 p.
  142. Chen M.K. The G-index for program priority // Health Status Indexes, 1973. – P. 29 ÷ 39.
  143. Chiang Ch.L. Life table and mortality analysis. – Geneva: World Health Organization, 1978. – 399 p.
  144. Exploring Health Policy Development in Europe / Ed. A. Risatakiset et al. – Copenhagen: WHO, RP, European Series, № 86, 2000. – 536 p.
  145. Fanshel S., Bush J.A health status index and its application to health services sateames. Operations Research, 1979, vol. 18, № 6. – P. 1021 ÷ 1056.
  146. Fisher R.A. The Design of Experiments. – Edinburgh: Oliver and Boyd, 1942. ÷ 466 p.

147. Fujimoto S., Watanabe T., Sakamoto A., Yukawa K., Morimoto K. Studies on the physical surface area of Japanese 18. Calculation formulae in three stages over all ages. Nippon Eiseigaku Zasshi. 1968
148. Gehan E.A., George S.L. Estimation of human surface area from height and weight. Cancer Chemother Rep part 1. 1970.
149. Granger C.W.J., Newbold P. Forecasting Economic Time Series, 2nd ed. – Academic Press, Inc., 1986. – 338 p.
150. Haycock G.B., Schwartz G.J., Winstock D.H. Geometric method measuring body surface area: A height-weight formula validated in infants, children, and adults. J Pediatr 1978.
151. Kaplan R.M. New health promotion Indicators: the general health policy model // Health Promotion. – V. 3, № 1, 1996. – P. 35 ÷ 49.
152. Medic V.A., Kirianov B.F., Bachmanov A. A. Model of an integral evaluation of health of the population // The Novgorod University Scientific Papers, 2002. – 6 p.
153. Miller J.E. An indicator to and management in assisting program priorities // Pub. Health Rep. – V. 85, № 8. P. 573 ÷ 600.
154. Morris M.D. Measuring the condition of the worlds poor: the physical quality of life index. – New York: Pergamon Press, 1979. – 24 p.
155. Mosteller R.D. Simplified calculation of body surface area. N Engl. J Med. 1987.
156. N. Colundina, V. Nekrutkin, A. Zhigjavaky. Analysis of Time Series Structure SSA and Related Techniques. – Chapman and Hall / CRC, 2002. – 303 p.
157. Sullivan D.A. Single index of mortality and morbidity. – HSMHA Health Reports, V. 86, № 4, 1971. – P. 347 ÷ 354.
158. Torrance Y.W. Health status index models: a unified mathematical view // Masnag. Sei., V. 22, № 9, 1976. – P. 990 ÷ 1001.
159. World Health Organization: Health: Services Research Methodology Core Library Recommendations, Switzerland, 2007.
160. World Health Organization: Statistics of World Health Organization, 2001.