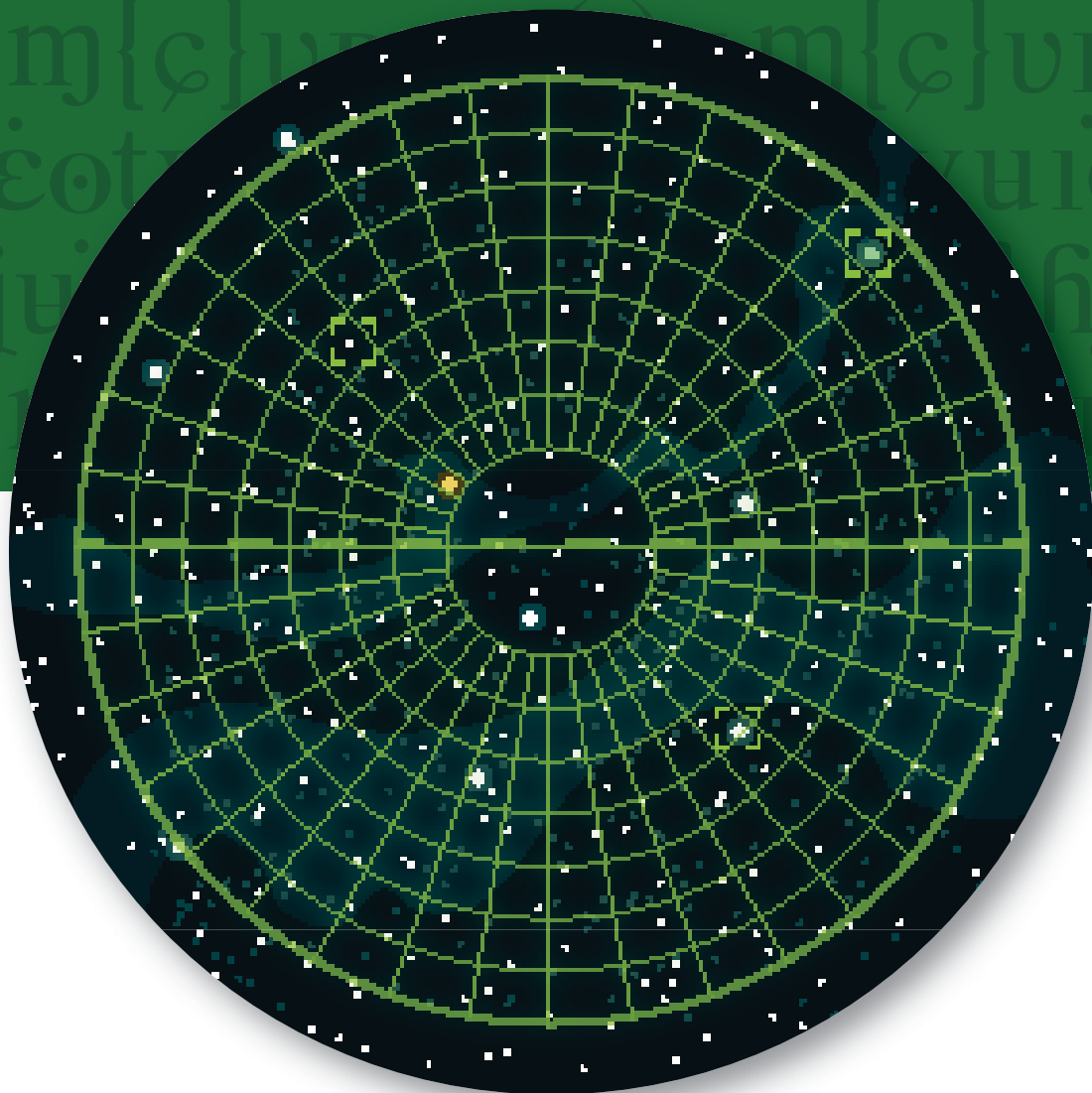




СИБИРСКИЙ ФЕДЕРАЛЬНЫЙ УНИВЕРСИТЕТ  
SIBERIAN FEDERAL UNIVERSITY



Б. С. Добронец

О. А. Попова

ВЫЧИСЛИТЕЛЬНЫЙ  
ВЕРОЯТНОСТНЫЙ АНАЛИЗ:  
МОДЕЛИ И МЕТОДЫ

УДК 519.676  
ББК 22.192.3  
Д564

**Р е ц е н з е н т ы:**

К. В. Сафонов, доктор физико-математических наук, профессор, заведующий кафедрой прикладной математики СибГУ им. М. Ф. Решетнёва;

Г. А. Доррер, доктор технических наук, профессор, профессор кафедры системотехники СибГУ им. М. Ф. Решетнёва

**Добронев, Б. С.**

Д564 Вычислительный вероятностный анализ: модели и методы : монография / Б. С. Добронев, О. А. Попова. – Красноярск : Сиб. федер. ун-т, 2020. – 236 с.

ISBN 978-5-7638-4232-6

Изложен подход к использованию вычислительного вероятностного анализа для решения задач с неопределенными входными данными. Основное внимание уделено процессу обработки, представления, моделирования и анализа информации для разных типов неопределенности. Рассмотрены различные математические модели и численные методы их обработки, вопросы надежности результатов численного моделирования для разнообразных задач в условиях ограниченного и большого объемов информации. Даны примеры применения рассматриваемого подхода для практических задач цифровой экономики, надежности технических систем и оборудования. Разработанные алгоритмы могут быть использованы для исследования сложных систем с входными данными, обусловленными различными типами неопределенности.

Предназначена для магистрантов, аспирантов и специалистов, занимающихся научными исследованиями и работающих в области решения задач с неточными входными данными.

**Электронный вариант издания см.:**  
<http://catalog.sfu-kras.ru>

**УДК 519.676**  
**ББК 22.192.3**

ISBN 978-5-7638-4232-6

© Сибирский федеральный университет, 2020

# Оглавление

<b>Введение</b>	<b>6</b>
<b>1. Краткий обзор теории вероятностей</b>	<b>18</b>
1.1. Понятие измеримости . . . . .	18
1.2. Борелевские $\sigma$ -алгебры . . . . .	20
1.3. Вероятностные пространства и случайные величины . . . . .	21
1.4. Лемма Doob–Dynkin . . . . .	24
1.5. Интегрируемость и моменты случайных величин . . . . .	25
1.6. Случайные векторы и их вероятностные распределения . . . . .	26
1.7. Независимость и корреляция случайных величин . . . . .	27
1.8. Произведение вероятностных пространств . . . . .	28
1.9. Случайные поля . . . . .	29
1.10. Параметризация случайных коэффициентов . . . . .	32
<b>2. Непараметрические оценки функций плотности вероятности</b>	<b>34</b>
2.1. Гистограммы . . . . .	34
2.2. Частотные полигоны . . . . .	41
2.3. Ядерные оценки функции плотности вероятности . . . . .	43
2.4. Экстраполяция Ричардсона и правило Рунге . . . . .	45
<b>3. Функциональный анализ данных</b>	<b>54</b>
3.1. Введение . . . . .	54
3.2. Примеры функциональных данных . . . . .	58
3.3. Функциональные модели данных . . . . .	61
3.4. Цели функционального анализа данных . . . . .	65
3.5. Функциональная регрессия . . . . .	65
3.6. Прогноз плотности . . . . .	68

<b>4. Символьный анализ данных</b>	<b>74</b>
4.1. Символьные данные . . . . .	76
4.2. Типы переменных . . . . .	79
4.3. Классические переменные . . . . .	80
4.4. Новые типы переменных . . . . .	80
4.5. Категориальные многозначные переменные . . . . .	82
4.6. Квантильное представление . . . . .	83
4.7. Другие типы символьных данных . . . . .	84
4.8. Методы анализа символьных данных . . . . .	85
4.9. Символьная регрессия . . . . .	86
4.10. Анализ временных рядов . . . . .	87
<b>5. Функции случайных переменных</b>	<b>88</b>
5.1. Алгебра случайных переменных . . . . .	88
5.2. Вероятностные расширения . . . . .	90
5.3. Одномерный случай . . . . .	95
5.4. Случай двух переменных . . . . .	97
5.5. Многомерный случай . . . . .	101
5.6. Краевые задачи со случайными коэффициентами . . . . .	103
5.7. Надежные оценки эмпирических распределений . . . . .	105
<b>6. Алгебраические задачи с неопределенностями</b>	<b>116</b>
6.1. Интервальные СЛАУ . . . . .	116
6.2. Системы линейных алгебраических уравнений со случайными коэффициентами . . . . .	120
6.3. Использование вероятностных расширений . . . . .	124
6.4. Совместное использование метода Монте-Карло и вычислительного вероятностного анализа . . . . .	127
6.5. Решения нелинейных уравнений . . . . .	128
6.6. Системы нелинейных уравнений . . . . .	130
<b>7. Временные ряды распределений</b>	<b>134</b>
7.1. Основы временных рядов распределений . . . . .	137
7.2. Оценка погрешности для временных рядов распределений . . . . .	137
7.3. Прогноз временных рядов распределений . . . . .	138
7.4. Методы сглаживания для временных рядов распределений . . . . .	140
7.5. Метод расщепления . . . . .	146
7.6. Численный пример . . . . .	148

<b>8. Случайное программирование</b>	<b>152</b>
8.1. Постановка задачи . . . . .	155
8.2. Случайное линейное программирование . . . . .	156
8.3. Случайное нелинейное программирование . . . . .	160
<b>9. Регрессионный анализ</b>	<b>163</b>
9.1. Регрессионные модели над эмпирическими распределениями	164
9.2. Агрегация данных . . . . .	167
9.3. Регрессионное моделирование на основе агрегированных данных . . . . .	170
9.4. Классическая параметрическая регрессия . . . . .	171
9.5. Метрики в пространстве распределений . . . . .	172
9.6. Регрессия над эмпирическими распределениями . . . . .	173
9.7. Эмпирическая функциональная регрессия . . . . .	174
9.8. Применение регрессионного подхода к функциональным временным рядам . . . . .	177
<b>10. Приложения ВВА</b>	<b>181</b>
10.1. Проблемы цифровой экономики . . . . .	182
10.2. Методика построения гарантированных оценок показателей надёжности . . . . .	190
10.3. Оценка показателей надёжности . . . . .	196
10.4. Обработка и анализ гидрологических данных спутникового мониторинга . . . . .	202
10.5. Оптимизация выработки электроэнергии гидроэлектростан- цией в условиях неопределенности . . . . .	208
10.6. Технология извлечения и визуализации знаний . . . . .	212
10.7. Визуально-интерактивная анимация . . . . .	216
<b>Заключение</b>	<b>222</b>
<b>Список литературы</b>	<b>224</b>

# Введение

Монография посвящена вопросам исследования сложных систем на основе применения современных математических методов представления, численного моделирования и анализа в условиях различных видов неопределенности данных. Большинство компьютерных моделей для инженерных приложений разрабатываются для того, чтобы помочь оценить проектные или нормативные требования. В рамках этой задачи критически важна способность количественно оценить влияние изменчивости и неопределенности в контексте принимаемого решения. Вычислительная стоимость инженерных имитационных моделей довольно дорога: для моделирования с конечными элементами высокой точности может потребоваться несколько часов или дней, десятки процессоров. Таким образом, понимание того, как работают методы снижения уровня неопределенности и их относительные преимущества и затраты, очень важно.

В работе обсуждаются и находят дальнейшее развитие идеи, представленные в монографии [17], рассматриваются новые, активно развивающиеся направления анализа данных, такие как вероятностный анализ (probabilistic analysis), функциональный (functional analysis) и символьный анализ (symbolic analysis). Изучаются новые аспекты повышения точности и организации вычислительного процесса обработки и анализа данных, связанные с разработкой технологии быстрых и надежных вычислений.

Предлагаются новые методы и алгоритмы, учитывающие такие виды информационной неопределенности, как элиторная (aleatory uncertainty) и эпистемическая (epistemic uncertainty). Теория вероятностей предназначена для моделирования, оценки и оперирования именно элитерными неопределенностями. Элиторная неопределенность характеризует присутствующую случайность в поведении системы или в стадии ее изучения. Она включает в себя: изменчивость, стохастическую неопределенность. Примерами случайной неопределенности являются отказы компонентов си-

стемы, полученные в результате статистически значимых испытаний в условиях, относящихся к применению. Элиторные неопределенности характеризуются частотными распределениями.

В свою очередь, неопределённость самих вероятностных оценок называют эпистемической. Эпистемическая неопределённость прямо связана с объёмом и достоверностью информации, на основании которой получаются эти оценки [68].

Эпистемические неопределенности могут быть устранены путем более глубокого понимания (исследования), на основе увеличения объема данных или с помощью более новых достоверных предположений.

Проблема надежных вычислений сегодня выходит на передний план среди проблем вычислительной математики. Следует отметить, что значительную часть производимых сегодня в мире вычислений нельзя назвать надежными, поскольку методы обеспечения надежности еще не получили должного распространения, а после выполнения обычных вычислений пользователи не всегда могут получить убедительные аргументы относительно важнейших свойств полученного решения, в том числе и его точности.

Надежные вычисления (reliable computing) достигаются с учетом многих факторов, прежде всего оценками погрешности вычислительных алгоритмов и учетом неопределенностей входных данных. В этой связи важное значение приобретают апостериорные оценки погрешностей результатов численного моделирования [35]. Для практической реализации идеи повышения надежности вычислений важную роль сыграли достижения интервальной математики. Корректные интервальные вычисления гарантируют выполнение важнейших свойств численного решения и прежде всего — его локализацию.

В настоящее время актуализировалась проблема применения и разработки вычислительных технологий, реализующих технику быстрых и надежных вычислений для решения разнообразных практических задач, имеющих отношение к исследованию состояний и процессов функционирования сложных систем. Например, использование систем искусственного интеллекта в технике и других областях неизбежно приводит к необходимости обработки огромных массивов информации, поступающих в устройства. В этой связи специалисты по созданию интеллектуальных систем столкнулись с проблемой обработки данных объемов (big data). Отметим также задачи, которые решаются в рамках бизнес-аналитики, дистанционного мониторинга распределенных систем, робо-

тотехники, гидро- и атомной энергетики, при анализе отказов технических систем ответственного назначения, оценки и прогнозирования техногенных, экологических, экономических и других видов рисков и т. д. Информация, которая составляет основу подобных задач, характеризуется имеющимся объемом данных, неоднородностью, динамичностью, уровнем и различными видами неопределенности.

Специфика сложности исследования таких систем обуславливается как объективными, так и субъективными аспектами. К объективным аспектам можно отнести следующие три группы факторов. Первая группа обуславливается внутренней сложностью системы как таковой. Вторая — внешней сложностью, непредсказуемостью, неопределенностью явлений и процессов, влияющих на систему и взаимодействующих с ней. Третья группа факторов связана с особенностями имеющейся у исследователя эмпирической информации и возможностями для ее обработки и анализа. Субъективный аспект связан прежде всего с тем, что практикам необходимо иметь определенный уровень доверия к применяемым математическим моделям и методам. Для них важно иметь убедительный ответ на вопрос, суть которого заключается в возможности получить достоверные, обоснованные результаты исследований, позволяющие установить с помощью численных расчетов достаточно полезную и реалистичную картину последствий принимаемых управленческих решений, несмотря на тот факт, что информация, на основе которой принимается решение, носит существенно неопределенный характер.

Обеспечение необходимой надежности и сложность исследования таких систем требует привлечения большого объема материальных, финансовых, интеллектуальных, временных, информационных и других ресурсов. При этом практика показывает, что привлекаемые ресурсы и вложения их в исследования не всегда пропорциональны требуемому уровню надежности и качеству функционирования систем в условиях различных видов неопределенности. Поэтому изучение способов и разработка новых моделей и методов представления информационной неопределенности в данных, обоснованное применение известных методов моделирования и разработка новых, реализующих перечисленные выше аспекты, представляет собой актуальную задачу.

Существующая неопределенность информации отражается в данных. Можно выделить три типа «неопределенных» данных: случайные, нечеткие и интервальные. Случайные числа задаются некоторыми вероятностными распределениями их возможных значений, нечеткие данные зада-



ются лингвистически сформулированными распределениями их возможных значений, интервальные данные задаются интервалами их возможных значений без указания какого-либо распределения внутри заданного интервала [101, 10, 50]. Изучение интервальной неопределенности способствовало созданию интервального анализа. Для случайной неопределенности знание законов распределения случайных величин позволяет оценивать параметры стохастических систем, используя метод Монте-Карло. Теория нечетких множеств широко используется для моделирования систем и принятия решений. В настоящее время для ряда задач в условиях стохастической неопределенности используется вычислительный вероятностный анализ [17, 31, 33, 71, 76].

В ряде случаев он успешно заменяет метод Монте-Карло [27, 39, 51], обладая значительно более высокой скоростью сходимости. В отличие от метода Монте-Карло он направлен на непосредственное построение распределений вероятности выходных переменных. Это существенно повышает качество полученных численных решений.

Для оценки качества решений важное значение имеет надежность полученных результатов. Любое измерение и методы его обработки содержат неточности. Рассмотрим последовательно этапы «эпохи» развития надежных вычислений. До «эпохи» надежных вычислений использовали «сырые данные» без предварительной обработки. Первый этап надежных вычислений заключался в статистической обработке и приближенном вычислении различных статистических характеристик. Ошибки численных методов приближенно оценивались с помощью двусторонних методов, например, правило Рунге, машинные арифметики не учитывали ошибки округления на компьютерах [21].

Второй этап — эра интервального анализа (ИА) началась с 50-х годов прошлого века. На этом этапе неопределенные данные представлялись в виде интервальных данных. Машинные арифметики, используемые в ИА, уже учитывали ошибки округления, а ошибки численных методов оценивались с помощью интервалов. ИА дает полностью гарантированные оценки, при этом значительно увеличивая время работы алгоритмов. К недостаткам интервального анализа можно отнести значительную ширину интервальных оценок по сравнению с оптимальными. ИА не использует информацию о возможных распределениях входных данных и соответственно не дает внутреннего распределения результатов вычислений, которые часто оказывались сосредоточенными только в небольших областях. Интервальные данные можно отнести к эпистимическому

типу неопределенности. Несмотря на указанные недостатки интервальный анализ позволяет эффективно решать многие практические задачи и широко используется при численном моделировании, например [24].

Третий этап — использование распределенных данных, в частности функций плотностей вероятности. Понятие распределенных данных — достаточно новое и появилось в научной литературе совсем недавно. Начало было положено разработкой численных операций над плотностями случайных величин, включая гистограммную арифметику. Одно из интересных представлений распределенных данных — символьные данные. Символьные данные были описаны Edwin и Diday в 1987 году [58]. Символьные переменные позволяют описывать группы индивидов и понятия. Символьные переменные включают списки значений (с весами или без них), интервальные переменные и даже гистограммы. Символьные представления могут включать внутреннюю структуру (иерархии) и логическую зависимость (правила).

Другой подход, при котором данные представляются в агрегированном виде, получил название Granular Computing (см., например, [112]). Информационные гранулы определяются, как группы отдельных наблюдений, которые отражают семантику абстрактных объектов, представляющих интерес. Как правило, с учетом набора данных  $D$ , в результате грануляции получается набор гранул, образованных на основе сходства или близости, которая может быть достигнута, например, с помощью алгоритмов кластеризации. Когда данные числовые, гранулы часто принимают форму гиперкубов. Информационные гранулы, описанные в теории нечетких множеств, представляются с помощью функции принадлежности. Распределенные переменные позволяют описывать каждую группу переменных посредством распределений. Распределения не используют статистические данные, такие как среднее, дисперсия, минимум и максимум и т. д. На практике сосредотачиваются на представлениях, которое лучше подходит для решения проблемы. Методы для распределенных данных включают в себя следующие разделы: описательная статистика, регрессия, кластеризация, уменьшение размерности, прогнозирование временных рядов, методы визуализации. Параллельно с символьным анализом развивается вычислительный вероятностный анализ (ВВА).

Вычислительный вероятностный анализ разработан как новое направление в вычислительной статистике (Computational Statistics) и предназначен для решения практических задач, связанных с исследованиями сложных систем в условиях различных видов неопределенности и типов

эмпирических данных. Основой ВВА являются численные операции над плотностями случайных величин. В ВВА используются различные типы представления плотностей случайных величин: дискретные, гистограммы, полигоны, кусочно-полиномиальные модели и аналитическое представление. Использование порядковых статистик и случайных интерполяционных полиномов позволяет строить достоверные оценки функций распределения [36].

Одним из основных разделов вычисленного вероятностного анализа являются арифметики над данными, представленными в виде кусочно-полиномиальных функций.

В рамках ВВА реализуется подход, созвучный тезису «распределения — числа будущего» (Distributions are the Numbers of the Future), сформулированному в 1984 году В. Schweizer [126].

Суть данного подхода реализует идею представления эмпирических данных в виде функций распределений на основе применения кусочно-полиномиальных моделей. Разработанные в рамках ВВА численные арифметики и использование нового понятия «вероятностное расширение» позволили авторам разработать методы численного моделирования и анализа распределений, которые можно рассматривать как особый вид переменных, над которыми выполняются соответствующие операции и процедуры.

Отметим, что ВВА оперирует в первую очередь с понятиями функции плотности вероятности (ФПВ) и для изучения свойств изменчивости данных и разработки численных операций над ними использует ФПВ-представление в виде кусочно-полиномиальных моделей.

Применение его методов и процедур позволяет представить выходное распределение вероятностей как функцию входных распределений и использовать методы анализа неопределенностей, чтобы оценить влияние входных неопределённостей, привносимых входными характеристиками, на неопределенность выходных параметров модели.

Для построения распределений используются специальные способы агрегации данных [73, 74], рассматриваются задачи интерполяции [36], задачи оценки надежности.

В рамках ВВА решаются задачи случайной оптимизации [69, 113], повышения точности и оценок погрешности получаемых решений [35, 71].

Как показал анализ литературы, проблема снижения уровня неопределенности в исходных данных и повышения эффективности численных методов представления, обработки, моделирования и анализа в течение

многих десятилетий находится в центре внимания и остается предметом многих научных исследований. Наиболее значимые результаты в данной предметной области были получены учеными: В. Liu [95], S. Ferson [85, 86], А. Neumaier [107, 108], Н. Schjaer-Jacobsen [125], D. Dubois [22], О. И. Ужга-Ребровым [42, 43, 44, 45, 46].

Актуализировалась проблема вычисления и анализа неопределённости выходных характеристик системы, индуцированных неопределённостями на её входах [19]. Среди публикаций, посвященных вопросам оценивания, анализа, управления неопределенностями следует указать на работы О. И. Ужга-Реброва. Можно выделить три основных группы методов, направленных на исследование и решение данной проблемы [42]: методы представления и оценивания неопределённости на выходе (выходах) модели, в зависимости от вида и уровня неопределённости на её входах (методы распространения неопределённостей); методы расчёта эффекта изменений на входах на предсказания модели, т. е. анализ чувствительности; методы сравнения важности входных неопределённостей в терминах их относительных вкладов на неопределённость на выходе (выходах), т. е. анализ неопределённостей.

Наиболее общий и распространенный метод включения неопределенности в моделирование состоит в том, чтобы предположить определенные распределения неопределенных входных значений, произвести выборку из этих распределений, запустить модель с выбранными значениями и делать это многократно, чтобы создать распределение выходных данных. Это классическое распространение неопределенности.

Для анализа и распространения элиторных и эпистемических неопределенностей в инженерных моделях используются также подходы: методы с использованием выборки на основе латинского гиперкуба (Latin Hypercube sampling), аналитические методы надежности (Analytic Reliability Methods) и методы разложения по полиномиальному хаосу (Polynomial Chaos Expansions) [78]. Эти методы являются альтернативными методами статистических испытаний, но опираются на аналитические вычисления и требуют от входных данных определенных свойств гладкости и принадлежности к гауссовым распределениям.

Методы построения выборки могут быть различными, например, можно использовать метод статистических испытаний (метод Монте-Карло), включая построение стратифицированной выборки (Latin Hypercube sampling), которая распределяет выборки по пространству, или квазивыборку, построенную методом Монте-Карло, который является способом гене-

рации последовательностей, приближающихся к равномерному распределению.

Отметим, что при решении проблемы снижения уровня выходной неопределенности важное значение имеет способ получения дополнительных оснований (знаний), снижающих уровень неопределенности во входных данных. Чтобы получить необходимые основания для оценки или восстановления неизвестного входного распределения на основе неполной, неточной информации, можно использовать различные процедуры и способы представления неопределенностей. Например, P-boxes [86], облака [108], интервальные гистограммы, гистограммы второго порядка [19].

Как один из подходов к идее «распространения неопределенности» использовался метод вероятностных границ (Probability bounds). Его основная идея в том, что функция неизвестного распределения вероятностей (Cumulative Distribution Function) должна лежать в области — ящике (box), ограниченная нулем и единицей по вертикали и от минимума и максимума горизонтально. Истинная функция распределения, какой бы она ни была, должна находиться в этой области. Облака Неймайера (Neumaier's clouds) являются еще одним способом представления неопределенности, выступая посредником между понятием нечеткого множества и вероятностным распределением [108].

В рамках основных подходов к распространению неопределенности следует указать также на математическую теорию очевидностей (свидетельств) Демпстера–Шафера [128], основанную на функции доверия (belief functions) и функции правдоподобия (plausible reasoning), которые используются, чтобы скомбинировать отдельные части информации (свидетельства) для вычисления вероятности события. Данная теория позволяет построить необходимые основания в условиях неопределенности путем оценки верхней и нижней границы интервала возможностей.

Среди подходов к распространению неопределенностей следует особенно выделить метод, который опирается на понятие «вероятность второго порядка», и известен как second-order probability. Данный подход представляет собой метод, позволяющий строить вероятностные оценки в случае эпистемистической неопределенности. Концепция вероятностей второго порядка была изложена в 1996 году в работах А. Mosleh и V. M. Bier. Анализ публикаций показал, что, несмотря на то, что данное направление достаточно активно развивается за рубежом, понятие «вероятность второго порядка» еще находится в стадии определения [103].

В монографии [17] предлагается новый способ представления данных в виде гистограмм второго порядка. Такой способ преобразования данных можно рассматривать в контексте проблемы распространения неопределенности и эффективно применять, когда законы распределения вероятностей зависят от неопределенных параметров.

Далее рассмотрим более подробно содержание глав монографии. Отметим, что при составлении содержания монографии и написания ее текста авторы исходили из того, что при работе с эмпирическими данными процесс их исследования представляет собой последовательность взаимосвязанных этапов, включающую методы предобработки, обработки и постобработки данных. Многообразие применяемых на каждом этапе методов, актуализирует проблему оценки точности полученных результатов. Ее решение во многом определяется надежностью тех вычислительных алгоритмов и методов, которые были выбраны для обработки, численного моделирования и анализа данных. Очевидно, что уже на стадии подготовки и преобразования данных необходимо применять процедуры представления данных в зависимости от имеющегося объема информации и типа неопределенности.

В главе 1 приводятся основные сведения из теории вероятностей.

В главе 2 рассматривается непараметрический подход, применяемый в настоящее время для оценки эмпирических функций распределений. Он имеет свои плюсы и минусы. Так, в отличие от параметрических методов не требует предположений о виде закона распределения наблюдаемых величин. Заметим, что во многих случаях достаточно сложно найти убедительные доказательства, по которым конкретное распределение результатов наблюдений должно входить в то или иное параметрическое семейство. В работе для решения задач анализа статистической информации развиваются методы ядерного оценивания и повышения их точности.

Главы 3, 4, 7 посвящены вопросам представления различных типов данных в виде математических моделей, приводятся примеры таких данных, рассматриваются вопросы построения и исследования функциональных (глава 3), символьных (глава 4) и временных рядов распределений (глава 7). В монографии большое внимание уделяется представлению и исследованию различных типов данных и выбору соответствующих процедур их обработки и анализа. Например, широко распространены типами данных, которые в процессе наблюдения за объектом или объектами фиксируются непрерывно в течение определенного про-

межутка времени или периодически в дискретные моменты времени. Высокая внутренняя размерность этих данных создает проблемы как для теории, так и для вычислений, а их исследование требует применения специальных методов и подходов. Эти проблемы зависят во многом от того, как были собраны данные, какова структура и размерность данных, каковы их источники. Ответы на эти вопросы позволяют выявить новые направления исследований и анализа данных с целью разработки моделей и методов, учитывающих их особенности и повышающих надежность численных процедур обработки и анализа данных. Отмечается, что для изучения таких данных можно применять функциональный анализ данных (ФАД) (Functional Data Analysis) [102, 116, 117, 118, 119], который занимается анализом и теорией данных, представленных в виде некоторых функций, изображений или более общих объектов. Одним из основных понятий ФАД является понятие функциональных данных, которые представлены так, что для каждого субъекта в случайной выборке записывается одна или несколько функций.

Важно отметить, что идея представления эмпирических данных на основе применения математических моделей на этапе предобработки данных и последующего их использования в виде входных и выходных факторов для моделирования способствовала появлению особого вида переменных. Например, использование гистограммных моделей данных в виде входных переменных для регрессионного моделирования способствовало появлению нового понятия гистограммно-значные переменные, которые представляют собой особый вид переменных, где каждому такому объекту (признаку) соответствует распределение, которое может быть представлено в виде гистограммы. Такие переменные изучаются, например, в символьном анализе [90, 106, 122, 137] (глава 4). В последнее время наблюдается растущий интерес к моделированию и анализу интервально-значных и гистограммно-значных [63, 105, 106]. Однако анализ публикаций по данной теме исследований показал, что существующие методы и подходы к регрессионному моделированию на гистограммно-значных переменных встречают ряд трудностей [64]. Например, для линейных моделей регрессии для этого типа данных отмечается, что ее параметры не могут быть отрицательными. Для определения параметров этой модели необходимо решить квадратичную задачу оптимизации, при условии неотрицательности ограничений на неизвестных. Определенную проблему составляет задача выбора и вычисления меры погрешности между предсказанными и наблюдаемыми распределения-

ми. Избежать этих трудностей можно, используя численные операции над функциями плотностей, что как раз может быть успешно реализовано в рамках ВВА.

Глава 5 посвящена функциям от случайных аргументов, где определяется новое и одно из основных понятий ВВА — вероятностное расширение, здесь также рассматриваются вопросы построения надежных оценок для функций распределения.

В главе 6 обсуждаются вопросы использования численных вероятностных арифметик, реализующих численные операции над кусочно-полиномиальными представлениями функций плотности вероятности и вероятностных расширения для решения систем линейных и нелинейных алгебраических уравнений со случайными параметрами.

В главе 8 рассматривается применение ВВА к решению задач оптимизации со случайными входными параметрами (случайная оптимизация). В результате решения подобных задач методами математического программирования получаются оптимальные решения, зависящие от этих параметров. В тех случаях, когда известны плотности вероятности входных параметров, на основе вычислительного вероятностного анализа возможно построение совместной функции плотности вероятности оптимального решения. В отличие от стохастического программирования [129], где оптимальное решение представляет собой некоторое фиксированное решение, данный подход позволяет построить все множество решений оптимизационной задачи, определяемое построенной совместной функцией плотности вероятности. Методы, позволяющие строить множество решений оптимизационной задачи со случайными входными параметрами на основе применения численного вероятностного анализа, назовем случайным программированием.

В главе 9 исследуются вычислительные проблемы построения регрессионных моделей над эмпирическими распределениями. Исследуются вопросы агрегированного представления данных и методы построения регрессионных моделей с агрегированными входными параметрами и функциональными временными рядами. Изучаются различные метрики в пространстве распределений.

Глава 10 посвящена приложениям ВВА для практических задач. Например, рассматриваются основные вычислительные аспекты, характерные для задач цифровой экономики. Первый аспект связан с необходимостью обработки данных больших объемов. Для его реализации предлагается использовать процедуры агрегирования данных, основанные на



применении математических моделей представления данных. Переход к более обобщенному представлению с помощью агрегирования необходим по нескольким причинам. Во-первых, агрегация существенным образом может снизить объем данных. Во-вторых, детализированные данные часто оказываются очень изменчивыми из-за воздействия различных случайных факторов, разброса значений и поэтому слабо отражают общие тенденции и свойства исследуемого множества. Агрегация в этом случае позволяет увидеть имеющиеся тенденции и закономерности. Второй аспект связан с организацией вычислительного процесса, обеспечивающей необходимую для решения соответствующей практической задачи оперативность получения необходимой информации. Для преодоления этой проблемы предлагается использовать рекурсивную схему организации вычислительного процесса. Третий аспект отражает требование к достоверности полученных результатов моделирования, обеспеченных надежными вычислительными процедурами, адекватными тем типам неопределенности, которые содержатся в сырых данных.

Рассматривается задача построения достоверных оценок показателей надежности оборудования в условиях малых выборок статистических данных об отказах. Применение вычисленного вероятностного анализа (ВВА) позволяет получить гарантированные оценки показателей надежности функционирования технических объектов в условиях неопределенности и ограниченного объема информации.

# Глава 1

## Краткий обзор теории вероятностей

Основные понятия и определения теории вероятностей, необходимые для дальнейшего изложения, рассматриваются в этом разделе. Опираясь на работы Rudin [123], Loéve [96] и Rao и Swift [120], представлено краткое введение к основам теории вероятностей, а затем исследовано несколько важных понятий, таких как вещественные случайные величины и векторы, понятие моментных операторов и случайных процессов. Дальнейшие концепции теории вероятностей можно найти, например, в [134], [96] и [92].

### 1.1. Понятие измеримости

Класс непрерывных функций играет фундаментальную роль в топологической теории. Он имеет несколько элементарных свойств, общих с измеримыми функциями, которые играют важную роль в теории интегрирования. Далее будет представлен материал, подчеркивающий аналогии между понятиями топологического пространства, открытого множества и непрерывных функций с измеримыми пространствами, измеримыми множествами и измеримыми функциями. Здесь  $\Omega$  определяется как непустое множество с конечным или бесконечным (счетным или несчетным) количеством элементов  $\omega$ .

**Определение 1 (топологическое пространство).** *А топологией  $\mathcal{F}$  на непустом множестве  $\Omega$  является коллекция подмножеств  $\Omega$  такая, что*

1)  $\emptyset \in \mathcal{F}, \Omega \in \mathcal{F};$

2) если  $A_i \in \mathcal{F}, i = 1, 2, \dots, n$   $\bigcap_{i=1}^n A_i \in \mathcal{F};$

3) если  $A_\alpha \in \mathcal{F}$  для  $\alpha \in \mathcal{A}$ , для произвольного набора индексов  $\mathcal{A}$ , то  $\bigcup_{\alpha \in \mathcal{A}} A_\alpha \in \mathcal{F}$ ,

где члены  $\mathcal{F}$  называются открытыми множествами  $\Omega$ , а упорядоченная пара  $(\Omega, \mathcal{F})$  называется топологическое пространство.

**Определение 2** ( $\sigma$ -алгебра и измеримое пространство). Коллекция  $\mathcal{F}$  подмножества непустого множества  $\Omega$  называется  $\sigma$ -алгеброй  $\Omega$ , если  $\mathcal{F}$  удовлетворяет

1)  $\Omega \in \mathcal{F}$ ;

2) если  $A \in \mathcal{F}$ , то  $\Omega \setminus A \in \mathcal{F}$ ;

3) если  $\{A_n\}_{i=1}^n \subset \mathcal{F}$ , то  $\bigcup_{i=1}^n A_n \subset \mathcal{F}$ ,

в этом случае упорядоченная пара  $(\Omega, \mathcal{F})$  называется измеримым пространством, а члены  $\mathcal{F}$  называются измеримые множества в  $\Omega$ .

**Определение 3** (измеримая функция). Пусть  $(\Omega, \mathcal{F})$  и  $(\Upsilon, \Sigma)$  — измеримые пространства. Тогда функция  $\mu : \Omega \rightarrow \Upsilon$  измерима, если для каждого  $A \in \Sigma$ , прообраз  $A$  под  $\mu$  находится в  $\mathcal{F}$ , т. е.

$$\mu^{-1}(A) \equiv \{\omega \in \Omega | \mu(\omega) \in A\} \subset \mathcal{F}.$$

**Определение 4** (положительная мера и пространство меры). Пусть  $(\Omega, \mathcal{F})$  измеримое пространство. Функция  $\mu : \mathcal{F} \rightarrow [0, \infty]$  называется положительной мерой, если  $\mu$  удовлетворяет следующему.

1) Неотрицательность: для всех  $A \in \mathcal{F}$ ,  $\mu(A) \geq 0$ .

2) Пустое множество:  $\mu(\emptyset) = 0$ .

3) Счетная аддитивность: если  $A_1, A_2, \dots \in \mathcal{F}$  и  $A_i \cap A_j = \emptyset$  для  $i \neq j$ , то

$$\mu\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \mu(A_i).$$

Тройка  $(\Omega, \mathcal{F}, \mu)$  называется измеримым пространством.

**Замечание 1.** Измеримое пространство часто называют «упорядоченными тройками»  $(\Omega, \mathcal{F}, \mu)$ , где  $\Omega$  — множество,  $\mathcal{F}$  —  $\sigma$ -алгебра в  $\Omega$ , а  $\mu$  — мера, определенная на  $\mathcal{F}$ . Аналогично измеримые пространства часто называют «упорядоченными парами»  $(\Omega, \mathcal{F})$ .

Эти соглашения имеют здравый смысл и являются логически правильными, несмотря на то, что они несколько избыточны. Например, с учетом вышеупомянутого упорядоченная пара, множество  $\Omega$  является просто наибольшим членом  $\mathcal{F}$ ; следовательно, учитывая  $\mathcal{F}$ , мы можем построить  $\Omega$ . Более того, по определению каждая мера принимает  $\sigma$ -алгебру как его область, так что, учитывая меру  $\mu$ , мы можем вывести  $\sigma$ -алгебру  $F$ , в которой  $\mu$  определено, и мы также знаем множество  $\Omega$ , в котором  $\mathcal{F}$  является  $\sigma$ -алгеброй, поэтому допустимо использовать выражения «пусть  $\mu$  будет мерой» или «пусть  $\mu$  будет мерой на  $\Omega$ , если мы хотим подчеркнуть множество, или даже пусть  $\mu$  будет мерой на  $\mathcal{F}$ , если мы хотим подчеркнуть  $\sigma$ -алгебру». Обычный подход, который логически довольно бессмысленен, это сказать «пусть  $\Omega$  будет пространством меры», даже когда понятно, что есть мера, определенная на  $\mathcal{F}$  в  $\Omega$ , и это мера, которая нас математически интересует.

## 1.2. Борелевские $\sigma$ -алгебры

$\sigma$ -алгебра Бореля является важным примером  $\sigma$ -алгебры, которая используется в теории функций, интеграла Лебега и теории вероятности. Дадим определение и сформулируем классическую теорему о  $\sigma$ -алгебрах.

**Теорема 1.** Пусть  $\Omega$  — множество, а  $\mathcal{V}$  — непустая совокупность подмножества  $\Omega$ . В  $\Omega$  существует наименьшая  $\sigma$ -алгебра, обозначаемая  $\sigma(\mathcal{V})$  такой, что  $\mathcal{V} \subset \sigma(\mathcal{V})$ , а именно

$$\sigma(\mathcal{V}) = \bigcap \{ \mathcal{F} : \mathcal{F} \text{ является } \sigma\text{-алгеброй } \Omega, \mathcal{V} \subset \mathcal{F} \},$$

которая также называется  $\sigma$ -алгеброй, порожденной  $\mathcal{V}$ .

Теперь пусть  $\Omega$  — топологическое пространство. По теореме 1 если  $\mathcal{V}$  является набором открытых множеств (или, что то же самое, все замкнутые множества)  $\Omega$ , то наименьшая  $\sigma$ -алгебра  $\mathcal{B} = \sigma(\mathcal{V})$  называется борелевской  $\sigma$ -алгеброй на  $\Omega$ . Элементы  $B \in \mathcal{B}$  называются борелевскими множествами, которые могут быть сформированы из открытых множеств (или, что то же самое, из закрытых множеств) через операции счетного пересечения, счетного объединения и относительного дополнений. Поскольку  $\mathcal{B}$  является  $\sigma$ -алгеброй, мы можем рассматривать  $(\Omega, \mathcal{B})$  как измеримое пространство, где борелевские множества играют роль

измеримых множеств. Если  $\mu : \Omega \rightarrow \Upsilon$  является непрерывным отображением  $\Omega$ , где  $\Upsilon$  — другое топологическое пространство, тогда из определений получаем, что  $\mu^{-1}(A) \in \mathcal{B}$  для каждого открытого множества  $A \in \Upsilon$ .

В заключение, каждое непрерывное отображение  $\Omega$  измеримо по Борелю. Измеримые по Борелю отображения часто называют борелевскими отображениями, или борелевскими функциями.

### 1.3. Вероятностные пространства и случайные величины

#### Вероятностная мера

По сути, вероятность является числовой мерой неопределенности результатов действия или эксперимента. Фактическое присвоение этих значений должно быть основано на опыте и поддаваться проверке при проведении эксперимента, если это возможно, повторяться при практически одинаковых условиях. Чтобы построить аксиоматическое представление, мы сначала представляем все возможные результаты эксперимента как отдельные точки непустого множества. С момента сбора все такие возможности могут быть бесконечно большими, различные комбинации их, полезные для экспериментов, необходимо учитывать. Затем мы определяем комбинации таких результатов как *события* и рассматриваем алгебру событий как первичные данные, которые включают в себя все мыслимое использование для эксперимента. Наконец, каждому событию присваивается числовая мера, которая соответствует «количеству» неопределенности таким образом, что эта неопределенность обладает аддитивными свойствами. Математически эта аксиоматическая формулировка была создана Колмогоровым.

**Определение 5 (вероятностная мера и вероятностное пространство).** Пусть  $(\Omega, \mathcal{F})$  — измеримое пространство, представляющее все возможные результаты эксперимента, где члены  $\sigma$ -алгебры  $\mathcal{F}$ , называемые событиями, являются коллекциями результатов эксперимента.

$P : \mathcal{F} \rightarrow [0, 1]$  называется вероятностная мера, или просто вероятность, если мера на  $(\Omega, \mathcal{F})$  удовлетворяет условиям

$$P(A) > 0 \text{ для всех } A \in \mathcal{F}$$

и

$$P(\Omega) = 1.$$

*А пространство вероятностей является тройкой  $(\Omega, \mathcal{F}, P)$ .*

Таким образом, вероятностное пространство — это пространство с конечной мерой, у которого функция меры нормируется так, чтобы мера всего пространства была равна единице. Пространство  $(\Omega, \mathcal{F}, P)$  называется полное пространство вероятностей, если  $\mathcal{F}$  содержит все подмножества  $A$  в  $\Omega$  с внутренней мерой  $P$ :

$$P^*(A) = \inf\{P(F) : F \in \mathcal{F}, A \subset F\} = 0.$$

Подмножества  $A$  в  $\Omega$ , которые принадлежат  $\mathcal{F}$ , называются  $\mathcal{F}$ -измеримыми. Однако в контексте вероятности интерпретации эти событий разные. Например, когда мы пишем  $P(A)$ , что означает «вероятность того, что событие  $A$  произойдет». В частности, если  $P(A) = 1$ , мы говорим, что « $A$  происходит с вероятностью 1», или «почти наверняка».

## Условная вероятность

Пусть  $(\Omega, \mathcal{F}, P)$  обозначает вероятностное пространство, и пусть  $A_1, A_2 \in \mathcal{F}$  события с  $P(A_1) > 0$  и  $P(A_2) > 0$ . Обозначим пересечение  $A_1 \cap A_2$   $A_1 A_2$ . Тогда отношение  $P(A_1 A_2) / P(A_1)$  называется *условная вероятность  $A_2$  при заданном  $A_1$* , или просто вероятность  $A_2$ , заданная  $A_1$ , и обозначается через  $P(A_2|A_1)$ , так что

$$P(A_1 A_2) = P(A_1) P(A_2|A_1). \quad (1.1)$$

Тогда по индукции для  $A_1, A_2, \dots, A_N \in \mathcal{F}$  получаем правило цепочки

$$P\left(\bigcap_{i=1}^N A_i\right) = P(A_1)P(A_2|A_1) \dots P(A_1 A_2 \dots A_{N-1}|A_N). \quad (1.2)$$

Более того, если  $\cup_{i=1}^N A_i = \Omega$  с  $A_i \cap A_j = \emptyset$  и  $A_i$  и  $B \in \mathcal{F}$ ,

$$P(B) = P(\Omega B) = \sum P(A_i B). \quad (1.3)$$

Тогда правило полной вероятности следует из (A.1), а именно

$$P(B) = \sum P(A_i)P(B|A_i). \quad (1.4)$$

Наконец, используя (1.1)–(1.4), мы приходим к *теореме Байеса*:

$$P(A_j|B) = \frac{P(A_j)P(B|A_j)}{\sum P(A_i)P(B|A_i)}.$$