



**Гришин, С. И.** Компьютерный анализ данных. Модели, алгоритмы, программы [Текст] : монография / С. И. Гришин, В. Х. Кириллов, А. К. Ширшков ; под ред. В. Х. Кириллова. - Одесса : ВМВ, 2014. – 304 с. : табл., рис. - Библиогр.: с. 302-303. - ISBN 978-966-413-468-9.

Монография представляет собой обзор теоретических и компьютерных методов по математической и прикладной статистике для технологических, технических и экономических вузов. Цель монографии — теоретические и справочные сведения по решению основных задач статистического анализа данных, используя компьютерные среды MatLab и SPSS. Для каждой из рассмотренных в книге задач даны: краткие теоретические сведения, описание математического метода решения, описание порядка и этапов выполнения задания в компьютерных средах.

Рассмотрены основы работы в операционных средах MatLab и SPSS и проводится обзор основных задач прикладной статистики и теории вероятностей по обработке данных. Подробно обсуждаются следующие основные задачи обработки экспериментальных данных: интерполяция и аппроксимация данных, основы проверки статистических гипотез, корреляционный и регрессионный анализ, дисперсный анализ, анализ временных рядов. На конкретных примерах проводится подробное описание порядка решения упомянутых задач в математическом пакете MatLab и статистическом пакете SPSS.

Монография адресована широкому кругу читателей: студентам, бакалаврам и магистрантам инженерных и экономических специальностей, преподавателям естественных дисциплин, научным работникам, пользователям компьютеров, применяющим статистические методы в практической работе.

## **Введение**

Все мы хорошо знакомы с закономерными явлениями и закономерными изменениями, они составляют главный объект научных исследований. Например, по астрономическим наблюдениям Тихо Браге Кеплер установил законы движения планет, а по этим законам, путём обобщения экспериментальных данных проведенных Галилеем для свободного падения тел, Ньютон установил знаменитый закон Всемирного тяготения. Знание этих

законов позволяет точно вычислять траектории спутников и комет, корректировать траекторию полёта снаряда и ракеты или предсказать когда и в какой местности произойдёт солнечное затмение.

Но отнюдь не во всех ситуациях интересующих нас результат полностью и жёстко определяется влияющими на него известными нам факторами. Например, мы не можем указать, сколько часов будет светить электрическая лампочка, или как долго будет служить жёсткий диск или DVD дисковод. Невозможно предвидеть число посетителей магазина и количество товаров, которое они купят, количество автомобилей, проходящих по улице за определённый промежуток времени, - и т. д. Явления, в которых результат испытания полностью определяется влияющими на него известными факторами, называются *детерминированными* или *закономерными*, а те, в которых это не выполняется — *недетерминированными* или *стохастическими*.

В большинстве явлений особенно в процессе измерений при проведении эксперимента присутствуют оба вида изменчивости — и закономерная и случайная (стохастическая), и для нахождения закономерностей нам приходится «отсеивать» возмущающие случайные факторы. Например, при одних и тех же начальных условиях стрельбы имеет место определённое рассеивание при падении снарядов, это рассеивание обусловлено многочисленными второстепенными факторами, влияющими на полёт снаряда в атмосфере (ветер, изменяющаяся плотность воздуха, турбулентность и др.). Другой пример, автомобиль движется равномерно по прямолинейному участку дороги, практически для того, чтобы обеспечить такое движение, водитель должен постоянно корректировать это движение с помощью рулевой колонки и педалью акселератора. Неизбежность корректировки связана с действием на автомобиль многочисленных возмущающих факторов (неровность полотна дороги, боковой ветер, турбулентность обтекаемого воздуха и др.).

Если требуется обработать большой и сложный набор данных, состоящий из небольших порций информации, статистика поможет классифицировать и проанализировать ситуацию, предоставив полезный обзор и резюме основных характеристик этих данных. Если данных недостаточно, статистика поможет собрать их и обеспечить получение ответов на заданные вопросы.

Вероятность — средство работы с риском и неопределенностью. Вероятность показывает возможность наступления в будущем каждого из потенциальных событий, рассчитанную на основе информации о некоторой ситуации. В то время как статистика позволяет переходить от наблюдений к обобщениям относительно рассматриваемой ситуации, вероятность имеет обратную направленность — исходя из характеристики можно выяснить, какие данные, скорее всего, удастся получить, и возможность получения каждого из вариантов таких данных.

С развитием и ростом сложности аппаратного и программного обеспечения все больше и больше внимания уделяется развитию методов прикладной статистики, резко ослабляющих допущения, на которых строятся используемые математические модели. Разработкой в этом секторе рынка программного обеспечения заняты как известные лидеры, так и новые развивающиеся компании. Инструменты прикладной статистики могут быть представлены либо как самостоятельное приложение, либо как дополнения к основному продукту.

Определенным импульсом для привлечения новых пользователей прикладной статистики явилось включение служб анализа данных в функциональность СУБД их поставщиками: это Microsoft (Microsoft SQL Server), Oracle, IBM (IBM Intelligent Miner for Data). Исследования показывают, что основным аргументом при выборе этих систем является то, что они поставляются совместно с сервером баз данных. Интеграция с системами оперативной обработки информации позволяет собрать оригинальные входные данные для анализа, а не «правильные» отчеты, подготовленные в отделах и филиалах корпорации.

Не только исследователи, но и менеджеры получили реальную возможность использовать в целях верного принятия решений очень сложный математический аппарат анализа данных, и знание этого аппарата является для них весьма существенным. По этой же причине люди, непосредственно использующие анализ своей деятельности (менеджеры и аналитики), должны понимать опасность выбора неадекватных методов анализа, так как некорректные результаты могут привести к принятию неверных решений.

Материалы по рассматриваемой теме входят в содержание обязательных дисциплин, преподаваемых студентам, обучающимся по направлениям областей знаний «Информатика и вычислительная техника» и «Системные науки и кибернетика». Однако наименования дисциплин и рекомендуемый объем материала, приведенные в стандартах высшего образования Украины для различных направлений подготовки, отличаются.

Наиболее полно материалы прикладной статистики представлены в направлении «Компьютерные науки», где они нашли отражение в дисциплинах «Теория вероятностей, вероятностные процессы и математическая статистика» и «Интеллектуальный анализ данных».

Достаточно подробно материалы по рассматриваемой теме освещены в направлении «Прикладная математика» - в дисциплинах «Теория вероятностей», «Математическая статистика» и «Анализ данных», а также в направлении «Информатика» - в дисциплинах «Теория вероятностей и математическая статистика» и «Распределенные информационно-аналитические системы». Из отсутствующих материалов в этих направлениях подготовки хотелось бы отметить популярную в последние годы технологию Data Mining.

## Оглавление

|  |    |
|--|----|
| Введение.....  | 6  |
| Глава 1. Знакомство с математическим пакетом MatLab..... | 10 |
| 1.1. Обработка данных в системе MatLab.....              | 10 |
| 1.1.1. Формирование данных случайным образом.....        | 10 |
| 1.1.2. Визуализация матриц.....                          | 10 |
| 1.1.3. Функции обработки данных для векторов.....        | 13 |
| 1.1.4. Функции обработки массивов данных.....            | 15 |
| 1.1.5. Блочные матрицы.....                              | 18 |
| 1.2. Формирование данных.....                            | 22 |

|   |     |
|---|-----|
| 1.2.1. Excel Link.....  | 22  |
| 1.2.2. Конфигурирование Excel.....  | 23  |
| 1.2.3. Обмен данными между MatLab и Excel.....  | 24  |
| Глава 2. Обработка данных детерминированных систем.....                                 | 27  |
| 2.1. Интерполяция.....  | 28  |
| 2.1.1. Интерполяция по Лагранжу.....  | 29  |
| 2.1.2. Метод разделенных разностей.....   | 31  |
| 2.1.3. Итерационные методы интерполяции.....  | 33  |
| 2.2. Аппроксимация данных.....  | 34  |
| 2.2.1. Аппроксимация данных сглаживающей<br>кривой - линейный регрессионный анализ..... | 35  |
| 2.2.2. Аппроксимация данных с помощью<br>ортгональных полиномов.....                    | 38  |
| 2.3. Сплайн-интерполяция.....   | 40  |
| 2.4. Интерполяция и аппроксимация в MatLab.....   | 42  |
| Глава 3. Обработка данных стохастических систем.....                                    | 55  |
| 3.1. Data Mining и перегрузка информацией.....  | 56  |
| 3.2. Системы поддержки принятия решений.....  | 58  |
| 3.3. Базы данных — основа СППР.....   | 60  |
| 3.4. Хранилище данных.....  | 64  |
| Глава 4. Первое знакомство с SPSS.....  | 67  |
| 4.1. Начало работы в среде SPSS.....  | 68  |
| 4.2. SPSS для Windows — обзор.....  | 87  |
| 4.3. Редактор данных - SPSS Data Editor.....  | 99  |
| Глава 5. Основные положения теории вероятностей.....                                    | 114 |
| 5.1. Виды случайных событий.....  | 114 |
| 5.2. Понятие вероятности.....   | 115 |
| 5.3. Умножение вероятностей.....  | 117 |
| 5.4. Обобщение умножения и сложения<br>вероятностей.....                                | 120 |
| Глава 6. Случайные величины.....  | 125 |
| 6.1. Функции распределения случайной величины.....                                      | 125 |
| 6.2. Числовые характеристики непрерывных<br>случайных величин.....                      | 129 |
| Глава 7. Основные распределения случайных величин.....                                  | 138 |
| 7.1. Дискретные случайные величины.....   | 138 |
| 7.2. Непрерывные случайные величины.....  | 142 |
| Глава 8. Решение вероятностных задач в SPSS.....  | 151 |
| Глава 9. Основные задачи математической статистики.....                                 | 174 |
| 9.1. Выборки и их описание.....   | 174 |
| 9.1.1. Выборки.....   | 174 |
| 9.1.2. Выборочные характеристики.....   | 176 |
| 9.1.3. Ранги и ранжирование.....  | 181 |
| 9.1.4. Методы описательной статистики.....  | 183 |
| 9.1.5. Наглядные методы описательной статистики.....                                    | 186 |
| 9.1.6. Методы описательной статистики в пакете SPSS.....                                | 190 |

|   |     |
|---|-----|
| 9.2. Основы проверки статистических гипотез.....                          | 191 |
| 9.2.1. Статистические гипотезы.....                                       | 191 |
| 9.2.2. Проверка статистических гипотез о<br>функциях распределения.....   | 195 |
| 9.2.3. Критерий согласия Пирсона.....                                     | 198 |
| 9.2.4. Критерий Колмогорова — Смирнова.....                               | 205 |
| 9.3. Корреляционный и регрессионный анализ.....                           | 213 |
| 9.3.1. Корреляционная зависимость.....                                    | 214 |
| 9.3.2. Регрессионная зависимость.....                                     | 224 |
| 9.4. Дисперсионный анализ.....  | 239 |
| 9.4.1. Однофакторный дисперсионный анализ.....                            | 239 |
| 9.4.2. Двухфакторный дисперсионный анализ.....                            | 251 |
| 9.5. Анализ временных рядов.....  | 262 |
| 9.5.1. Цели, этапы и методы анализа временных рядов.....                  | 265 |
| 9.5.2. Детерминированная и случайная<br>составляющие временного ряда..... | 267 |
| 9.5.3. Тренд, сезонная и циклическая компоненты.....                      | 270 |
| 9.5.4. Модели тренда.....   | 274 |
| 9.5.5. Модели случайной компоненты.....                                   | 275 |
| 9.5.6. Порядок анализа временных рядов.....                               | 278 |
| 9.6. Анализ временных рядов в SPSS.....                                   | 280 |
| 9.6.1. Обзор возможностей.....  | 280 |
| 9.6.2. Подбор тренда и прогнозирование.....                               | 281 |
| 9.6.3. Устранение сезонной компоненты.....                                | 297 |
| Литература.....   | 302 |