

**Учреждение образования
«Гомельский государственный университет имени Франциска Скорины»**

Факультет физики и
информационных технологий
Кафедра теоретической физики

СОГЛАСОВАНО
Заведующий кафедрой
теоретической физики

Г.Ю.Тюменков

_____ 2019 г.

СОГЛАСОВАНО
Декан
факультета физики и
информационных технологий

Д.Л.Коваленко

_____ 2019 г.

**ЭЛЕКТРОННЫЙ УЧЕБНО-МЕТОДИЧЕСКИЙ КОМПЛЕКС
ПО УЧЕБНОЙ ДИСЦИПЛИНЕ**

**СТАТИСТИЧЕСКИЕ МЕТОДЫ
ОБРАБОТКИ ДАННЫХ**

для специальности

1-31.04.08 Компьютерная физика

Составители: д.ф.-м.н., доцент Андреев В.В.

Рассмотрено и утверждено
на заседании научно-методического
совета университета
_____ 2019__ г., протокол № ____

Гомель 2019

Содержание
учебно-методического комплекса по дисциплине специализации
«Статистические методы обработки данных»
для специальности 1-31 04 08 Компьютерная физика

- 01 Титульный лист
- 02 Содержание
- 03 Пояснительная записка
- 1. Теоретический раздел
 - 1.1 Теоретический материал в виде лекций
 - 1.2 Презентации лекционного материала
- 2. Практический раздел
 - 2.1 Материал к лабораторным работам (лабораторный практикум)
- 3. Контроль знаний
 - 3.1 Перечень вопросов к зачету
- 4. Вспомогательный материал
 - 4.1 Учебная программа дисциплины
 - 4.2 Перечень рекомендуемой литературы

Репозиторий ГГУ им. Ф. Скорины

ПОЯСНИТЕЛЬНАЯ ЗАПИСКА

Электронный учебно-методический комплекс (ЭУМК) по дисциплине «Статистические методы обработки данных» представляет собой комплекс систематизированных учебных, методических и вспомогательных материалов, предназначенных для использования в образовательном процессе специальности 1-31 04 08 Компьютерная физика.

ЭУМК разработан в соответствии со следующими нормативными документами:

1. Положением об учебно-методическом комплексе на уровне высшего образования, утвержденном постановлением Министерства образования Республики Беларусь от 26.07.2011 №167.

2. Образовательный стандарт высшего образования ОСВО 1-31 04 08-2013 Компьютерная физика.

3. Учебный план учреждения высшего образования «ГГУ им. Ф. Скорины», регистрационный № G 31-01-16, дата утверждения 12.01.2016

4. Учебной программой по учебной дисциплине «Статистические методы обработки данных» для специальности 1-31 04 08 Компьютерная физика, регистрационный номер № УД-33-2016-248/уч.

Цель создания ЭУМК – овладение студентами основами математической статистики и теории вероятностей и их использование при обработке и анализе экспериментальных данных.

ЭУМК направлен на оказание помощи студентам в приобретение навыков и изучение способов, приёмов работы по анализу и статистической обработке информации на персональном компьютере, наиболее часто используемых в физических исследованиях (обратив внимание на смысл понятий критериев, возможности алгоритмизации).

Организация изучения дисциплины на основе ЭУМК предполагает продуктивную образовательную деятельность, позволяющую сформировать социально-личностные и профессиональные компетенции будущих специалистов, обеспечить развитие познавательных и созидательных способностей личности.

ЭУМК способствует успешному осуществлению учебной деятельности, дает возможность планировать и осуществлять самостоятельную работу обучающихся, обеспечивает рациональное распределение учебного времени по темам учебной дисциплины и совершенствование методики проведения занятий.

Учебная дисциплина «Статистические методы обработки данных» состоит из семи тем – «Вероятность и характеристики распределения вероятностей»; «Основы теории оценок»; «Гистограммы и операции над ними»; «Идентификация формы распределения экспериментальных данных»; «Регрессионный и корреляционный анализ»; «Точечные и интервальные оценки распределений в системе МАТНЕМАТИСА»; «Функции распределений случайных величин и средства статистического анализа в системе МАТНЕМАТИСА».

Общее количество часов – 96; аудиторное количество часов – 64, из них: лекции – 24 (из них УСП-6) , лабораторные занятия – 34. Форма отчётности – зачет. Дисциплина читается для студентов 2 курса.

ЭУМК включает в себя: титульный лист, пояснительную записку, теоретический раздел, который содержит тексты лекций и презентации лекций по дисциплине “Статистические методы обработки данных”, практическую часть, которая содержит описание лабораторных работ по данной учебной дисциплине и самостоятельные задания, вспомогательный раздел, который содержит учебную программу и список рекомендуемой литературы дисциплины. Все разделы ЭУМК в полной мере соответствуют содержанию учебной программы и требованиям образовательного стандарта.

Репозиторий ГГУ им. Ф. Скорины

Учреждение образования
“ГОМЕЛЬСКИЙ ГОСУДАРСТВЕННЫЙ
УНИВЕРСИТЕТ
ИМЕНИ ФРАНЦИСКА СКОРИНЫ”
Кафедра теоретической физики

Статистические методы обработки данных

Конспект Лекций
Специальность
1-31 04 08 Компьютерная физика

Лекции: 30 часов
Практические занятия: -
Лабораторные занятия: 34 часа

Материал подготовил
Андреев
Виктор Васильевич
доктор физ.-мат. наук, доцент

Гомель, 2018

ОГЛАВЛЕНИЕ

1	Основные понятия теории вероятности и мат. статистики	3
1.1	Частота попадания случайной величины.	3
1.2	Интегральный закон распределения вероятности.	4
1.3	Дифференциальный закон распределения вероятности.	5
1.4	Квантили распределений	6
1.5	Примеры законов распределений.	7
2	Моменты распределения.	10
2.1	Математическое ожидание. Дисперсия.	10
2.2	Коэффициент асимметрии. Эксцесс.	13
3	Оценки параметров распределений	15
3.1	Точечные оценки	15
3.2	Интервальные оценки	18
3.3	Доверительный интервал для случайной величины	20
4	Взвешенное среднее значение	22
5	Ошибка косвенного измеряемой величины	23
6	Графическое представление эмпирических данных	24
7	Оптимальное число интервалов для получения гистограммы	26
8	Прوماхи и методы их исключения	31
8.1	Другие критерии для исключения промахов	34
9	Критерии согласия	35
10	Алгоритм предварительной обработки экспериментальных данных	40
11	Некоторые сведения о двумерных случайных величинах	41
11.1	Алгебраические и центральные моменты	43
11.2	Примеры многомерных функций распределения вероятностей	46
12	Корреляционный и регрессионный анализ	48
13	Линейный регрессионный анализ	52
13.1	Метод наименьших квадратов	52
13.2	Свойства метода наименьших квадратов	54
13.3	Точность оценок A и B	55
13.4	Редукция некоторых задач нелинейного регрессионного анализа к линейному	55

<small>Viktor Andreev</small>	<small>Viktor Andreev</small>	14 Элементы нелинейного регрессионного анализа	57
<small>Viktor Andreev</small>	<small>Viktor Andreev</small>	14.1 Метод максимального правдоподобия	58
<small>Viktor Andreev</small>	<small>Viktor Andreev</small>	14.2 Метод наименьших квадратов	63
<small>Viktor Andreev</small>	<small>Viktor Andreev</small>	15 Сравнение моделей	66
<small>Viktor Andreev</small>	<small>Viktor Andreev</small>	15.1 Проверка статистической значимости параметров уравнения регрессии	66
<small>Viktor Andreev</small>	<small>Viktor Andreev</small>	15.2 Проверка общего качества уравнения регрессии	71
<small>Viktor Andreev</small>	<small>Viktor Andreev</small>	15.3 Коэффициент детерминации	71
<small>Viktor Andreev</small>	<small>Viktor Andreev</small>	15.4 Недостаток R^2 и альтернативные показатели	74
<small>Viktor Andreev</small>	<small>Viktor Andreev</small>	16 Рекомендуемая литература	79

Репозиторий ГГУ им. Ф. Скорины

1 Основные понятия теории вероятности и мат. статистики

Для исследования свойств объекта проводят измерения, позволяющие количественно дать характеристики свойств этого объекта. На практике производится ограниченное число измерений n при одинаковых условиях. В результате получаем множество событий (значений) исследуемой величины X : $\{x_1, x_2, \dots, x_n\}$, которое представляет собой выборку объема n . Таким образом, возможные исходы некоторого эксперимента (измерения) представляют собой множество точек, которые можно сочетать разными способами. **Такие сочетания называют событиями.**

1.1 Частота попадания случайной величины.

В том случае, (т.е. когда одна и та же характеристика объекта принимает различные значения) говорят о том, что величина X является **случайной величиной**. Случайная величина - это действительное число (или набор действительных чисел), которое заключено между $-\infty$ и $+\infty$, которое сопоставляется каждой возможной точке из числа значений этой характеристики. При соответствующих условиях для каждое событие можно характеризовать **частотой появления** ν_n .

Определение 1.1

Частотой появления ν_n события A называется отношение числа m появления данного события к общему числу проведенных одинаковых испытаний, в каждом из которых могло появиться или не появиться данное событие:

$$\nu_n = \frac{m}{n}. \quad (1.1)$$

Определение 1.2

Если число испытаний n велико, то, как правило, частоты появления данного события A в различных сериях измерений отличаются мало друг от друга. Это утверждение записывают следующим образом:

$$\lim_{n \rightarrow \infty} \nu_n = p. \quad (1.2)$$

Число p называются вероятностью (англ.-probability) случайного события A .

Отметим, что существуют такие события у которых частота появления может сильно отличаться от вероятности, даже при большом числе испытаний.

Как видим из вышеприведенного примера, для описания случайного поведения величины необходима совокупность, содержащая неограниченное число значений измеряемой величины ($n = \infty$). Такая выборка называется генеральной совокупностью. Генеральная совокупность часто используется как важное абстрактное понятие, необходимое в теоретических расчетах, связанных с исследованием поведения физической величины как случайной величины.

Различают два основных типа случайных величин: дискретные случайные величины и непрерывные случайные величины. Если величина X имеет **конечное** число (счетное множество) из последовательности возможных значений $\{x_1, x_2, \dots, x_k, \dots\}$, то такая величина называется **дискретной случайной величиной**. Если случайная величина может принимать любое значение из интервала возможных значений, то такая величина называется **непрерывной случайной величиной**.

1.2 Интегральный закон распределения вероятности.

Для характеристики частоты появления различных значений случайной величины X теория вероятностей предлагает пользоваться указанием **закона распределения вероятностей** различных значений этой величины.

При этом различают два вида описания законов распределения:

1. интегральный закон распределения вероятности;
2. дифференциальный закон распределения вероятности.

Определение 1.3

Интегральным законом, или функцией распределения вероятностей $F(x)$ случайной величины X , называют функцию, значения которой представляют вероятность того, что значения x_k случайной величины X меньше некоторого значения x (x – некоторое произвольное число).

Данное утверждение символически записывается в виде

$$F(x) = \text{Prob}(x_k < x), \quad (1.3)$$

где $\text{Prob}(x_k < x)$ и представляет собой вероятность события в вышеприведенном определении.

Очевидно, что

$$F(a) \leq F(b), \quad \text{при } a \leq b \quad (\text{неубывающая функция}) \quad (1.4)$$

$$F(-\infty) = 0, \quad F(\infty) = 1. \quad (1.5)$$

Для дискретной одномерной случайной величины X закон распределения удобно представить в виде таблицы

$$X = \begin{pmatrix} x_1, x_2, \dots, x_n \\ p_1, p_2, \dots, p_n \end{pmatrix} \quad (1.6)$$

где x_1, x_2, \dots, x_n – значения случайной величины X , а p_1, p_2, \dots, p_n – вероятности появления этих значений.

Другими словами $\text{Prob}(X = x_i) = p_i$. В этом случае интегральный закон вероятности в соответствии с законом распределения вероятности имеет вид:

$$F(x) = \text{Prob}(x_k < x) = \sum_{i=1}^k p_i. \quad (1.7)$$

1.3 Дифференциальный закон распределения вероятности.

Для случайной величины с непрерывной и дифференцируемой функцией распределения $F(x)$ можно найти **дифференциальный закон распределения вероятностей**:

$$p(x) = \frac{dF(x)}{dx}. \quad (1.8)$$

$p(x)$ называют кривой плотности распределения вероятностей (или просто **плотность вероятности.**)

Свойства $p(x)$:

1. $p(x) \geq 0$
2. Из (1.8) следует, что

$$F(x) = \int_{-\infty}^x p(\xi) d\xi . \quad (1.9)$$

3. Условие нормировки:

$$\int_{-\infty}^{\infty} p(x) dx = 1 . \quad (1.10)$$

Для дискретной случайной величины из определения (1.6) следует, что

$$p(x) = \sum_{i=1}^n \delta(x - x_i) p_i , \quad (1.11)$$

где одномерная δ -функция Дирака определена соотношениями:

$$\delta(x) = \begin{cases} +\infty , & \text{если } x = 0 \\ 0 , & \text{если } x \neq 0 \end{cases} ,$$

$$\int_{-\infty}^{\infty} \delta(x) dx = 1 . \quad (1.12)$$

1.4 Квантили распределений

Определение 1.4

Квантилем случайной величины x вероятности q (или уровня значимости $\alpha = 1 - q$), называется величина x_q , для которой имеем:

$$F(x_q) = \text{Prob}(x < x_q) = \int_{-\infty}^{x_q} p(x) dx = q \quad (1.13)$$

или

$$\text{Prob}(x > x_q) = \int_{x_q}^{\infty} p(x) dx = \alpha = 1 - q . \quad (1.14)$$

Отметим, что с помощью законов распределений вероятности можно найти вероятность нахождения случайной величины в некотором интервале $[a, b]$:

$$\text{Prob}(a < x < b) = \int_a^b p(x) dx = F(b) - F(a) . \tag{1.15}$$

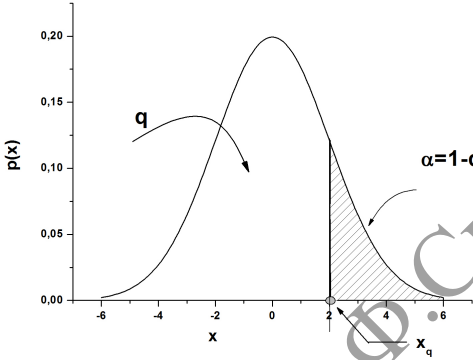


Рисунок 1– Иллюстрация квантиля распределения $p(x)$

1.5 Примеры законов распределений.

Нормальное распределение:

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left[-\frac{(x - a)^2}{2\sigma^2} \right] , \tag{1.16}$$

a и σ -параметры нормального распределения. **Распределение χ^2 с n степенями свободы:**

$$p(\chi^2) = \frac{(\chi^2)^{(\frac{n}{2}-1)} e^{-\frac{\chi^2}{2}}}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} , \quad (\chi^2 > 0) , \tag{1.17}$$

n – параметры χ^2 распределения.

Равномерное распределение:

$$p(x) = \frac{1}{b - a} , \quad (b > a) , \tag{1.18}$$

a, b - параметры распределения.

Рисунок 2– Нормальное распределение с $a = 0$ и $\sigma = 1$

Рисунок 3– Распределение χ^2 с $n = 3,4,5,6$ степенями свободы

Рисунок 4– Равномерное распределение с $a = 0$ и $b = 1,2,4$

Дискретные распределения:

Распределение Пуассона:

$$p(x) = \frac{e^{-\mu} \mu^k}{k!} k \geq 0, \quad (1.19)$$

μ - параметр распределения, а переменная $k = 0, 1, \dots$

Репозиторий ГГУ им. Ф. Скорины

2 Моменты распределения.

Моменты k -того порядка для непрерывной случайной величины записываются в виде:

$$\mu_k = \int_{-\infty}^{\infty} x^k p(x) dx, \quad (2.1)$$

где μ_k – алгебраический момент k -того порядка.

$$\nu_k = \int_{-\infty}^{\infty} (x - \mu_1)^k p(x) dx, \quad (2.2)$$

где ν_k – центральный момент k -того порядка.

2.1 Математическое ожидание. Дисперсия.

Первый алгебраический момент называют **математическим ожиданием**:

Определение 2.1

Математическим ожиданием непрерывной случайной величины с плотностью вероятности $p(x)$ называется величина $E\{x\}$, определяемая соотношением:

$$\mu_1 = E\{x\} = \int_{-\infty}^{\infty} xp(x) dx. \quad (2.3)$$

Для функции случайной величины $g(x)$ соотношение (2.3) обобщается следующим образом:

$$E\{g(x)\} = \int_{-\infty}^{\infty} g(x)p(x) dx. \quad (2.4)$$

Иногда, для сокращения записи, используем следующее обозначение для математического ожидания:

$$E\{g(x)\} \equiv \langle g(x) \rangle. \quad (2.5)$$

Используя (1.11), (2.3) и свойство δ -функции

$$\int_{-\infty}^{\infty} f(x)\delta(x-a)dx = f(a) \quad (2.6)$$

для дискретной случайной величины получаем, что математическое ожидание вычисляется по формуле:

$$E\{x\} = \sum_{i=1}^n p_i x_i . \quad (2.7)$$

Свойства математического ожидания (2.4):

1.

$$E\{c\} = c , \text{ если } c = \text{const} ; \quad (2.8)$$

2.

$$E\left\{\sum_{k=1}^n c_k g_k(x)\right\} = \sum_{k=1}^n c_k E\{g_k(x)\} , \text{ если } c_k = \text{const} ; \quad (2.9)$$

3.

$$E\left\{\sum_{k=1}^n c_k \xi_k\right\} = \sum_{k=1}^n c_k E\{\xi_k\} , \text{ если } \xi_k \text{ случайные величины} . \quad (2.10)$$

Математическое ожидание характеризует центр распределения $p(x)$. Однако следует отметить, что не для всех распределений существует математическое ожидание.

Наиболее общей характеристикой центра распределения следует считать медиану.

Определение 2.2

Медиана это такое значение случайной величины x_m для которой вероятности появления различных значений случайной величины X $p_1 = \text{Prob}(X < x_m)$ и $p_2 = \text{Prob}(X > x_m)$ равны между собой, т. е. $p_1 = p_2 = 0,5$.

Другими словами можно сказать, что медиана это квантиль распределения вероятности $q = 0,5$.

Также центр распределения может характеризоваться модой распределения.

Определение 2.3

Мода распределения это такое значение случайной величины x_{mod} для которой плотность вероятности появления значения случайной величины X максимальна, т. е. $p(x_{mod}) = \max p(x)$.

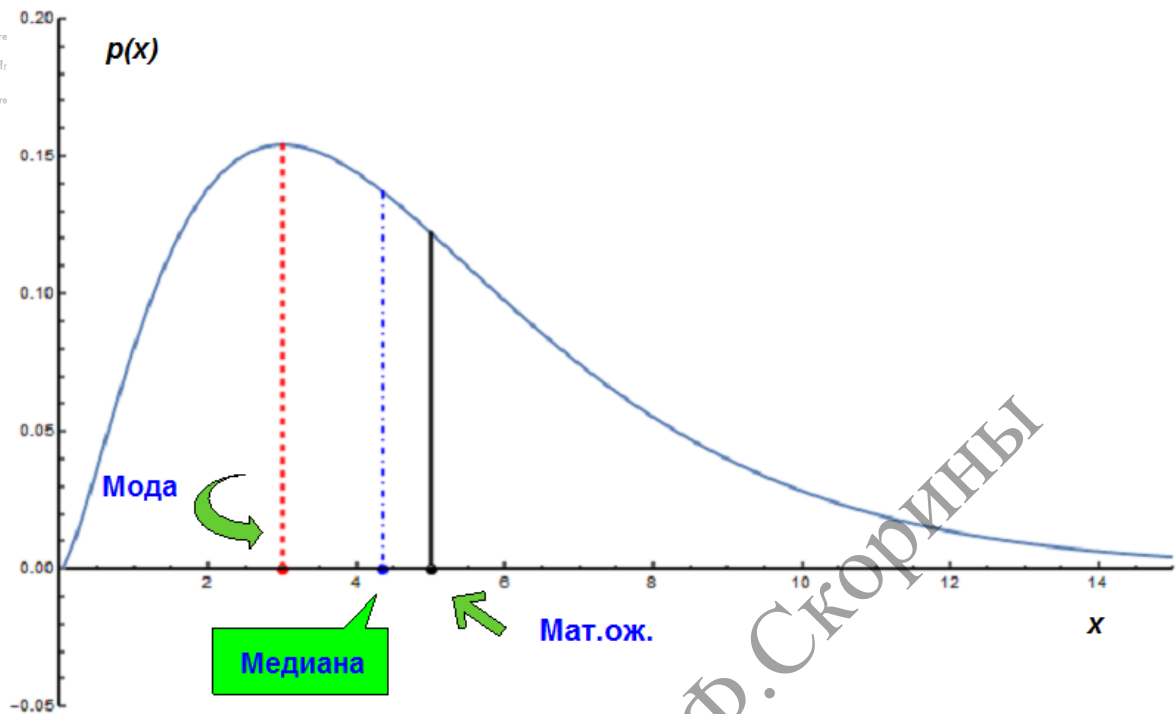


Рисунок 5– Пример распределения

Определение 2.4

Дисперсией непрерывной случайной величины с плотностью вероятности $p(x)$ называется величина $D(x)$, определяемая соотношением:

$$\nu_2 = D\{x\} = \int_{-\infty}^{\infty} (x - E\{x\})^2 p(x) dx . \quad (2.11)$$

Из определения (2.4) следует, что

$$D\{x\} = E\{(x - \mu_1)^2\} = E\{x^2\} - E^2\{x\} . \quad (2.12)$$

Свойства дисперсии (2.12):

1.

$$D\{c\} = 0 , \text{ если } c = \text{const} ; \quad (2.13)$$

2.

$$D\{cx\} = c^2 D\{x\} , \text{ если } c = \text{const} ; \quad (2.14)$$

3.

$$D\{c + x\} = D\{x\} , \text{ если } c = \text{const} . \quad (2.15)$$

Для дискретной случайной величины, с учетом (1.11) дисперсия вычисляется по формуле:

$$D\{x\} = \sum_{k=1} (x_k - E\{x\})^2 p_k. \quad (2.16)$$

Дисперсия характеризует рассеяние отдельных значений случайной величины от центра распределения и она определяет форму распределения. Дисперсия имеет размерность квадрата случайной величины и выражает как бы мощность рассеяния, поэтому для более наглядной характеристики используют среднеквадратичное отклонение σ :

$$\sigma_x = +\sqrt{D\{x\}}. \quad (2.17)$$

которое имеет размерность самой случайной величины.

Для описания относительной меры отклонения используют коэффициент вариации

$$V_x = \frac{\sigma_x}{E\{x\}}. \quad (2.18)$$

2.2 Коэффициент асимметрии. Эксцесс.

Центральные моменты 3,4 порядка дают информацию о виде кривой плотности распределения.

Определение 2.5

Третий центральный момент непрерывной случайной величины с плотностью вероятности $p(x)$, который определяется соотношением:

$$\nu_3 \equiv \int_{-\infty}^{\infty} (x - E\{x\})^3 p(x) dx, \quad (2.19)$$

характеризует асимметрию кривой $p(x)$ или скошенность распределения (например, когда один спад - крутой, а другой - пологий).

Для симметричных относительно центра распределений он равен нулю. ν_3 имеет размерность куба случайной величины, поэтому для относительной характеристики используют безразмерный **коэффициент асимметрии**:

$$\gamma_1 = \frac{\nu_3}{\nu_2^{3/2}}. \quad (2.20)$$

Четвертый центральный момент

$$\nu_4 \equiv \int_{-\infty}^{\infty} (x - E\{x\})^4 p(x) dx, \quad (2.21)$$

характеризует протяженность распределения (островершинность).

Определение 2.6

Безразмерную величину вида

$$\varepsilon = \frac{\nu_4}{\nu_2^2}$$

называют *эксцессом распределения*.

Область изменения эксцесса: $\varepsilon \in [1, \infty]$. Часто используют **коэффициент эксцесса**

$$\gamma_2 = \varepsilon - 3 \quad (2.23)$$

и **контрэксцесс** $\kappa = 1/\sqrt{\varepsilon}$.

3 Оценки параметров распределений

Определение 3.1

Функция результатов опытов, которая зависит от неизвестных статистических характеристик называют **статистикой**. Статистика зависит от случайных величин и сама является случайной величиной.

Для нахождения поведения случайной величины, полученных в результате эксперимента необходима числовая информация о моментах распределения.

Оценкой статистической характеристики $\tilde{\theta}$ называется статистика, которая принимается за неизвестное истинное значение параметра θ . Основное требование к оценке истинного значения состоит в том, чтобы большинство значений статистики сосредоточилось вблизи значений θ и вероятность больших отклонений от этого значения была мала. Также желательно, чтобы с увеличением объема выборки точность оценок также увеличилась.

Если имеется экспериментальная выборка объема n случайной величины $X = \{x_1, x_2, \dots, x_n\}$, то ее элементы можно рассматривать как n статистически независимых случайных величин. Тогда любая оценка $\tilde{\theta}$ параметра θ должна быть функций элементов выборки, т. е.

$$\tilde{\theta} = \theta(x_1, x_2, \dots, x_n) . \quad (3.1)$$

При этом закон распределения $\tilde{\theta}$ зависит от закона распределения величины X и от объема выборки.

Существуют два вида оценок параметров распределения: *точечные и интервальные*.

3.1 Точечные оценки

Под точечной оценкой параметра распределения понимают оценку одним числом. К точечным оценкам предъявляют следующие требования:

- **состоятельность**;
- **несмещенность**;
- **эффективность**;
- **устойчивость**;

надежность.

Определение 3.2

Метод оценки параметров называется **состоятельным**, если оценки, полученные с его помощью, сходятся к истинному значению с увеличением объема выборки n

Определение 3.3

Если $\tilde{\theta}$ оценка параметра θ , то оценка будет **несмещенной**, если смещение

$$E\{\tilde{\theta}\} - \theta \quad (3.2)$$

равно нулю (0), т. е. $E\{\tilde{\theta}\} = \theta$.

Определение 3.4

Оценка $\tilde{\theta}_{eff}$ является **эффективной**, если она обладает наименьшим разбросом относительно истинного значения параметра θ , т. е.

$$D\{\tilde{\theta}_{eff}\} \rightarrow \min\{\tilde{\theta}_1, \tilde{\theta}_2, \dots\} . \quad (3.3)$$

Определение 3.5

Под **устойчивостью** оценки понимают нечувствительность ее к малым отклонениям от точного распределения.

Точечная оценка для математического ожидания (центра распределения) дается выражением:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i , \quad (3.4)$$

Для дисперсии $D = \sigma^2$ такая оценка дается выражением:

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 , \quad S = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} . \quad (3.5)$$

Величины \bar{x} и S сами являются случайными величинами и, следовательно, они тоже могут иметь разброс, который характеризуется дисперсией.

Для \bar{x} :

$$S_{\bar{x}}^2 = \frac{S^2}{n}, \quad S_{\bar{x}} = \frac{S}{\sqrt{n}}. \quad (3.6)$$

Для дисперсии:

$$D[S^2] = \frac{\nu_4}{n} - \frac{\sigma^4(n-3)}{n(n-1)}, \quad (3.7)$$

где ν_4 – четвертый центральный момент, а σ – среднеквадратичное отклонение.

Несмещенными и состоятельными оценками $\tilde{\nu}_3$ и $\tilde{\nu}_4$ для третьего ν_3 и четвертого ν_4 центральных моментов соответственно являются

$$\tilde{\nu}_3 = \frac{n^2}{(n-1)(n-2)} m_3, \quad (3.8)$$

$$\begin{aligned} \tilde{\nu}_4 &= \frac{1}{(n-1)(n-2)(n-3)} \times \\ &\times [n(n^2 - 2n + 3)m_4 - 3n(2n-3)m_3^2], \end{aligned} \quad (3.9)$$

где

$$m_k = \sum_{i=1}^m (x_i - \bar{x})^k. \quad (3.10)$$

Тогда оценки (знак $\tilde{}$) для коэффициента асимметрии γ_1 и эксцесса ε (смотри соотношения (2.20) и (2.6)) определяются формулами:

$$\begin{aligned} \tilde{\gamma}_1 &= \frac{\tilde{\nu}_3}{S^3}, \\ \tilde{\varepsilon} &= \frac{\tilde{\nu}_4}{S^4}, \end{aligned} \quad (3.11)$$

где S – точечная оценка $\sqrt{D\{x\}}$.

Примечание.

В теории погрешностей, которая составляет основу обработки данных в лабораторных работах вместо обозначения S используют Δx . Эту величину называют **абсолютной погрешностью**, а также **стандартным отклонением**. Коэффициент вариации V_x называют относительной погрешностью и выражают в процентах:

$$V_x \rightarrow \epsilon_x = \frac{\Delta x}{\bar{x}} \times 100\%. \quad (3.12)$$

3.2 Интервальные оценки

Точечные оценки параметров случайных величин не позволяют судить о степени близости выборочных значений к оцениваемому параметру. Более содержательны процедуры оценивания параметров, связанные с построением интервала с известной степенью доверительности.

Для любой случайной величины X можно определить вероятности попадания в некоторый интервал $[x_{min}, x_{max}]$, если известен закон распределения вероятностей (смотри (1.15)):

$$\text{Prob}(x_{min} < x < x_{max}) = \int_{x_{min}}^{x_{max}} p(x, \theta) dx = F(x_{max}) - F(x_{min}), \quad (3.13)$$

где в плотность распределения введена величина θ , которая определяет параметры распределения.

Потребуем, чтобы вероятность попадания равнялась некоторому значению $P = 1 - \alpha$, т. е.

$$\text{Prob}(x_{min} < x < x_{max}) = P = 1 - \alpha. \quad (3.14)$$

С помощью (3.14) можно “построить” три вида интервалов с заданной вероятностью попадания p :

1. Верхний односторонний (левосторонний) интервал $]-\infty, x_{max}]$;
2. Нижний односторонний (правосторонний) интервал $]x_{min}, \infty[$;
3. Двусторонний интервал $[x_{min}, x_{max}]$.

Предполагается, что вероятность попадания в каждый их интервалов равна $P = 1 - \alpha$. Обычно величину P называют доверительной вероятностью (иногда надежностью), а величину α уровнем значимости интервала.

Интервальные оценки для моментов распределения находятся построением некоторой функции случайной величины, куда входят искомые параметры распределений θ и точечные оценки $\tilde{\theta}$. Тогда (3.14) применительно к параметрам можно записать

$$\text{Prob}(\tilde{\theta} - \delta_{min} < \theta < \tilde{\theta} + \delta_{max}) = P = 1 - \alpha, \quad (3.15)$$

который следует понимать так: вероятность того, чтобы параметр θ находится в интервале $[\tilde{\theta} - \delta_{min}, \tilde{\theta} + \delta_{max}]$ равна $P = 1 - \alpha$.

Отметим, что задача построения доверительных интервалов для параметров при произвольном объеме “экспериментальной” выборки n разработана только для нормального распределения случайной величины X . Для других распределений имеются только частные случаи.

Например, для любой выборки n из нормальной совокупности с математическим ожиданием ξ функция вида

$$t = \frac{\bar{x} - \xi}{(S/\sqrt{n})} \quad (3.16)$$

имеет t -распределение (распределение Стьюдента) с $n - 1$ степенями свободы, плотность которого определяется соотношением:

$$p_{\text{ST}}(t) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\sqrt{\pi n} \Gamma\left(\frac{n}{2}\right)} \left[1 + \frac{t^2}{n}\right]^{-\frac{n+1}{2}}. \quad (3.17)$$

Тогда можно найти такое значение $t_{n,P}$ случайной величины с распределением Стьюдента для которой выполняется

$$\text{Prob}\left(\left|\frac{\bar{x} - \xi}{S/\sqrt{n}}\right| < t_{n-1,P}\right) = P = 1 - \alpha \quad (3.18)$$

или

$$\text{Prob}\left(\bar{x} - \frac{t_{n-1,P}S}{\sqrt{n}} < \xi < \bar{x} + \frac{t_{n-1,P}S}{\sqrt{n}}\right) = P = 1 - \alpha. \quad (3.19)$$

В итоге имеем

Доверительный интервал для математического ожидания $E\{x\} = \xi$ (двусторонний), если неизвестна дисперсия распределения:

$$\left[\bar{x} - \frac{t_{n-1,1-\alpha/2}S}{\sqrt{n}} < \xi < \bar{x} + \frac{t_{n-1,1-\alpha/2}S}{\sqrt{n}}\right], \quad (3.20)$$

где $t_{n,P}$ определяется соотношением:

$$\int_{-\infty}^{t_{n,P}} p_{\text{ST}}(t) dt = P.$$

Коэффициенты $t_{n,P}$ носят название **коэффициентов Стьюдента**.

Для построения доверительного интервала для дисперсии $D\{x\} = \sigma^2$ используется, то, что функция

$$\frac{(n-1)S}{\sigma^2} \quad (3.21)$$

имеет распределение χ^2 с $n - 1$ степенями свободы.

Тогда **доверительный интервал для дисперсии** при неизвестном математическом ожидании имеет вид:

$$\left[\frac{(n-1)S^2}{\chi_{n-1, \alpha/2}^2} < D\{x\} < \frac{(n-1)S^2}{\chi_{n-1, 1-\alpha/2}^2} \right], \quad (3.22)$$

где $\chi_{n, \alpha}^2$ – квантиль распределения χ^2 :

$$\int_{\chi_{n, \alpha}^2}^{\infty} p(\chi^2) d\chi^2 = \alpha$$

с плотностью

$$p(\chi^2) = \frac{(\chi^2)^{\frac{n}{2}-1} e^{-\frac{\chi^2}{2}}}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})}, \quad (\chi^2 > 0).$$

Значения $E\{x\}$ и $D\{x\}$ лежат в интервале с доверительной вероятностью $P = 1 - \alpha$.

3.3 Доверительный интервал для случайной величины

С помощью точечных оценок, рассмотренных нами выше, результат для случайной величины x обычно записывают в виде:

$$(\bar{x} \pm \Delta x). \quad (3.23)$$

При этом согласно правилу приведения результатов:

Правило приведения результатов

Последняя значащая цифра в любом приводимом результате для \bar{x} обычно должна быть того же порядка величины (находиться в той же десятичной позиции), что и абсолютная погрешность Δx .

Однако используемые в расчетах числа должны, как правило, содержать на одну значащую цифру больше, чем это оправдано. Это уменьшит неточности, возникающие при округлении чисел. В конце расчета окончательный ответ следует округлить и избавиться от этой добавочной (и незначащей) цифры).

Придать вероятностный смысл интервалу (3.23) можно, зная лишь функцию распределения вероятности $p(x)$ (плотность вероятности). Если предположить, что наша выборка принадлежит генеральной совокупности с нормальным распределением ($E\{x\} = \bar{x}$, $D\{x\} = (\Delta x)^2$), то легко найти, что

$$\begin{aligned}
 & \text{Prob}(\bar{x} - \Delta x < x < \bar{x} + \Delta x) = \\
 &= \frac{1}{\Delta x \sqrt{2\pi}} \int_{\bar{x} - \Delta x}^{\bar{x} + \Delta x} \exp\left[-\frac{(x - \bar{x})^2}{2(\Delta x)^2}\right] dx \\
 &= \left(z = \frac{x - \bar{x}}{\Delta x}\right) = \\
 &= \frac{1}{\sqrt{2\pi}} \int_{-1}^1 \exp\left[-\frac{z^2}{2}\right] dz = \\
 &= \text{erf}\left(\frac{1}{\sqrt{2}}\right) \approx 0,682689. \tag{3.24}
 \end{aligned}$$

Поэтому в физических приложениях, часто говорят: вероятность того, что результат измерения окажется в пределах одного стандартного отклонения от истинного результата, составляет 68,3%. Можно сказать и так: вероятность, того измеряемая случайная величина лежит в интервале (3.23) составляет 68,3%.

4 Взвешенное среднее значение

В практических исследованиях часто встречается следующая ситуация: одна и та же величина измеряется в нескольких экспериментах. Это делается для наиболее надежного определения некоторой величины. При этом собирают измерения различного происхождения, выполненные разными установками (инструментами) и методами. Результаты таких измерений называют часто называют **неравноточными**.

Предположим, что у нас есть n отдельных измерений

$$x_1 \pm \Delta x_1, \quad x_2 \pm \Delta x_2, \dots, x_n \pm \Delta x_n, \quad (4.1)$$

с соответствующими погрешностями $\Delta x_1, \dots, \Delta x_n$.

Определение 4.1

Наилучшая оценка величины X , основанная на таких измерениях, равна в средневзвешенному значению \hat{x}_{sw}

$$\hat{x}_{sw} = \frac{1}{w} \sum_{i=1}^n w_i x_i, \quad (4.2)$$

$$w = \sum_{i=1}^n w_i, \quad w_i = \frac{1}{(\Delta x_i)^2}, \quad (4.3)$$

а её среднеквадратичное отклонение $\Delta \hat{x}_{sw}$

$$\Delta \hat{x}_{sw} = \frac{1}{\sqrt{w}}. \quad (4.4)$$

Результат такой обработки данных записывается в обычном виде:

$$\hat{x}_{sw} \pm \Delta \hat{x}_{sw}. \quad (4.5)$$

5 Ошибка косвенно измеряемой величины

Все, о чем мы говорили выше, относилось непосредственно к измеряемым величинам. А как определить погрешности для косвенно измеряемых величин, т. е. величин измеряемых с помощью физических законов из непосредственно измеряемых? Оказывается, погрешность косвенно измеряемых величин, связана с погрешностями непосредственно измеряемых величин.

Пусть косвенно измеряемая величина Z связана с непосредственно измеряемыми x_1, \dots, x_m (m величин) величинами функцией $Z = f(x_1, \dots, x_m)$. Если случайные величины x_i **не коррелированы друг с другом, т. е. независимы**, то абсолютная погрешность величины Z определяется соотношением

$$\sigma_Z = \Delta Z = \sqrt{\sum_{i=1}^m \left(\frac{\partial f}{\partial x_i}\right)^2 (\Delta x_i)^2}, \quad (5.1)$$

где $\Delta x_1, \dots, \Delta x_m$ среднеквадратичные отклонения непосредственно измеряемых величин.

Пример. $Z = f(x, y) = x \pm y$

$$\begin{aligned} \Delta Z^2 &= \underbrace{\left(\frac{\partial f}{\partial x}\right)^2}_{=1} (\Delta x)^2 + \underbrace{\left(\frac{\partial f}{\partial y}\right)^2}_{=1} (\Delta y)^2 = \\ &= (\Delta x)^2 + (\Delta y)^2. \end{aligned} \quad (5.2)$$

Практическое следствие этого соотношения: для создания оптимальных условий (условий с минимальной погрешностью) основные усилия должны быть направлены не на дальнейшее уточнение тех результатов измерений, которые являются наиболее точными, а на совершенствование наименее точных измерений случайных величин.

6 Графическое представление эмпирических данных

Цель обработки данных заключается в выявлении вида распределений случайных величин и оценки параметров установленного распределения.

Полученные экспериментальные данные представляют, как правило, в виде таблиц. Полученные таблицы удобно представить графически. Используя набор независимых наблюдений x_1, x_2, \dots, x_n случайной величины X , полезным первым шагом в исследовании поведения случайной величины является организация и представление их таким образом, чтобы их можно было легко интерпретировать и оценивать. Для достаточно большого количества наблюдаемых данных, полигон частот (распределения), гистограмма и кумулятивная линия является отличным графическим представлением данных, что облегчает оценку адекватности предполагаемой модели и оценку параметров распределения.

Гистограмма и полигон распределений являются графическим отображением частот, которые, в свою очередь, представляют собой оценки плотностей вероятностей $p(x)$. Кумулятивная линия - это график накопленных частот, в свою очередь оценивающих интегральную функцию распределения $F(x)$.

Как строить полигон частот

1. Построить вариационный ряд для выборки, т. е. упорядочить значения случайной величины так, чтобы выполнялось условие: $x_1 \leq x_2 \leq \dots \leq x_n$.
2. Разбить область на m интервалов (бинов).
3. Построить точки на плоскости x - ν с координатами $\{\tilde{x}_i, \nu_i\}$, ($i = 1, \dots, n$), где

$$\tilde{x}_i = \frac{x_i + x_{i+1}}{2} \quad (6.1)$$

является серединой i -того интервала, ν_i - частота попадания случайной величины X в i -тый интервал.

Это же распределение можно представить в виде гистограммы. Для построения гистограммы необходимо над каждым отрезком оси абсцисс, соответствующим интервалу значений измеряемой величины, построить прямоугольник, площадь которого пропорциональна частоте попадания в этот ин-

Viktor Andreev Viktor Andreev Viktor Andreev
Viktor Andreev Viktor Andreev Viktor Andreev
Viktor Andreev Viktor Andreev Viktor Andreev
Viktor Andreev Viktor Andreev Viktor Andreev

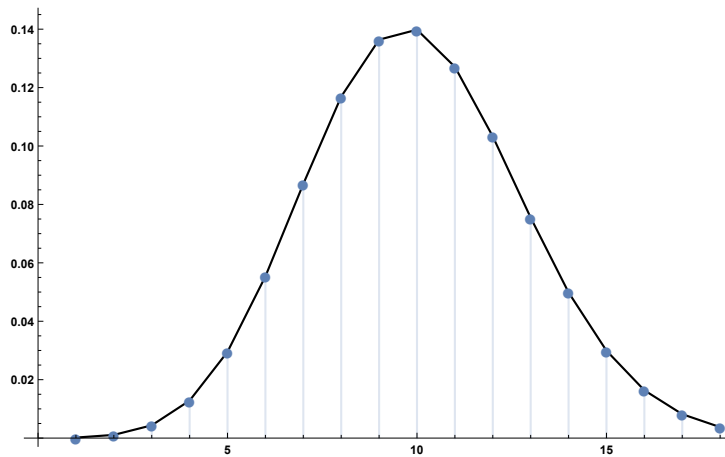


Рисунок 6– Пример полигона частот

тервал. Обычно выбирают интервалы одинаковой ширины, поэтому высота прямоугольников различна.



Рисунок 7– Пример гистограммы

Где используются обработка данных для построения гистограммы?

Для определения оценок математического ожидания, с.к.о., эксцесса не требуется какого-либо группирования данных.

Для определения медианы, сгибов, использования критерия согласия Колмогорова-Смирнова или для обнаружения промахов, экспериментальные данные необходимо расположить в порядке возрастания, т.е. построить вариационный ряд (упорядоченную выборку).

Для определения формы распределения, для использования критериев согласия Пирсона и др., для сопоставления гипотез о форме распределения и т.д. простого упорядочения выборки уже недостаточно, а выборка должна быть представлена в виде гистограммы, состоящей из m столбцов с определенной протяженностью d соответствующих им интервалов.

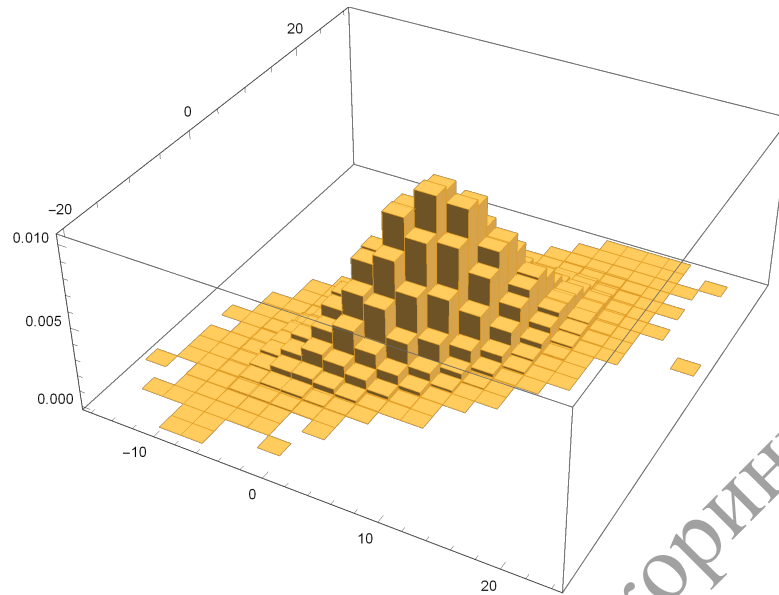


Рисунок 8– Пример 3D-гистограммы

7 Оптимальное число интервалов для получения гистограммы

Как выбрать m и d .

Оптимальное число

Оптимальное число существует!!

Оптимальное число интервалов группирования это такое число, когда ступенчатая огибающая гистограммы наиболее близка к плавной кривой распределения случайной величины.

К примеру, при группировании данных выбор большого числа интервалов ($m > n$), автоматически приведет к тому, что некоторые из них окажутся пустыми или малозначительными. Гистограмма будет отличаться от плавной кривой распределения вследствие изрезанности многими всплесками и провалами.

Тогда можно сформулировать 1-е требование к числу интервалов: Размер интервала (ячейки, бина) должен быть достаточно широким для обеспечения хороших статистических свойств будущей гистограммы (достаточно большая статистика, минимальные корреляции (связи) с соседними ячейками).

При слишком малом числе m интервалов, гистограмма отличается от

действительной кривой распределения вследствие слишком крупной ступенчатости. Из-за чего будут потеряны характерные особенности. Например, если взять $m = 1$, т.е. d равно размаху экспериментальных данных, то любое распределение сводится к равномерному, а если $m = 3$, то любое куполообразное распределение сведется к треугольному. Например, при обработке линейчатых спектров, слишком большая ячейка может привести к потере спектральной линии.

Тогда возникает **2-е требование к числу интервалов**:

Размер ячейки должен быть достаточно узким для того, чтобы прорисовывалась “тонкая структура” исследуемой величины.

Как видим, требования являются противоречивыми. Укрупнение интервалов группирования является методом “фильтрации различных случайных выбросов и провалов”, но слишком протяженные интервалы сглаживают особенности искомого закона распределения. Таким образом, задача выбора оптимального числа интервалов при построении гистограммы – это задача оптимальной фильтрации, а оптимальным числом m интервалов является максимальное возможное сглаживание случайных флуктуаций данных, которое сочетается с минимальным искажением от сглаживания самой кривой искомого распределения.

Общепринято делать интервалы **одинаковыми**. Хотя в дальнейшем увидим, что это условие необязательно. Условие равновеликости интервалов удобно с практической точки зрения.

Рекомендации по выбору m .

I группа: эвристические критерии (без доказательства).

Формула Старджеса

$$m = \log_2 n + 1 . \quad (7.1)$$

Формула Брукса и Каррузера

$$m = 5 \lg n . \quad (7.2)$$

Формула если $n > 100$

$$m = \sqrt{n} . \quad (7.3)$$

Эти три формулы являются наиболее часто встречающимися в литературе по математической статистике.

II группа: с использованием критерия χ^2 .

В ней используется рассмотрение интервалов не с равной длиной, а с **равной вероятностью** в соответствии с принимаемой моделью, т. е. предположением о законе распределения. В данном подходе неявно учитывается форма распределения.

Число интервалов с равной вероятностью, которые мы обозначили как K , отличаются от числа m с равной длиной d .

Г. Манн и А. Ваальд установили, что при $n \rightarrow \infty$ оптимальное число K равновероятных интервалов задается соотношением:

Критерий Мана-Ваальда

$$K \sim b\sqrt[5]{2} \left(\frac{n}{Z_\alpha} \right)^{2/5}, \quad (7.4)$$

где $b = 2 \div 4$.

Здесь Z_α – квантиль нормального распределения, соответствующий вероятности $P = 1 - \alpha$, α – принятый уровень значимости.

$$Z_\alpha = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{Z_\alpha} e^{-\frac{x^2}{2}} dx = 1 - \alpha.$$

Критерий Мана-Ваальда

На практике часто берут $\alpha = 0.1$, тогда

$$K \simeq 1,9n^{2/5}. \quad (7.5)$$

В итоге приходим к таким рекомендациям при использования равновероятностных бинов K :

1. найти число ячеек, используя (7.5);
2. если окажется, что n/K мало, уменьшить K чтобы выполнялось неравенство $n/K \geq 5$;
3. Сформировать K равновероятностных ячеек гистограммы на основе данных. Отметим, что если измеряемая случайная величина многомерная, то существуют различные способы формирования ячеек с одинаковым вероятностным содержанием в K ячейках.

III группа.

Поскольку для K интервалы получаются не равной длины, то это приводит к ряду неудобств при построении гистограмм, но зато при этом мы неявно закладываем при использовании χ^2 выбор K в зависимости от формы распределения.

III группа рекомендаций “устраняет” недостаток II группы возвращаясь к интервалам m с равной длиной d , но при этом и учитывает, в отличие от I группы, форму распределения (форма характеризуется эксцессом ε или контрэксцессом κ).

Примером такого подхода является соотношение:

III группа (формула И.У.Алексеевой)

$$m = \frac{4}{\kappa} \lg \frac{n}{10}, \quad \kappa = \frac{1}{\sqrt{\varepsilon}}. \quad (7.6)$$

На практике эту формулу в зависимости от n при $\kappa = const$ удобнее аппроксимировать выражением

$$m = \frac{\varepsilon + 1,5}{6} n^{2/5} \quad (7.7)$$

Трудность использования III группы состоит в том, что число интервалов часто приходится выбирать прежде, чем будут найдены оценки ε , \bar{x} , и т.д.

Эту трудность обходят следующим образом: наиболее часто встречаются распределения с ε от 1,8 до 6 (от равномерного до Лапласа, включая нормальное $\varepsilon=3$). Для этих граничных точек имеем

$$m_{min} = 0,55 n^{2/5} \quad (7.8)$$

$$m_{max} = 1,25 n^{2/5} \quad (7.9)$$

Искомое m можно выбрать близким к этому интервалу, при этом m лучше выбрать нечетным, т.к. при четном m для островершинных распределений в центре гистограммы оказывается два столбца равных по высоте и середина распределения принудительно утолщается.

“Практические” рекомендации для построения диаграмм

1. Для практического определения числа интервалов воспользоваться формулами для m_{min} (7.8) и m_{max} (7.9), или сразу формулой (7.6), выбрав при этом m нечетным.
2. Так как крайние точки могут располагаться несимметрично, то ширина d столбца гистограммы определяется по отклонению от центра ΔX_m наиболее удаленной точки:

$$d = \frac{2\Delta X_m}{m} . \quad (7.10)$$

При этом полученное значение d необходимо округлять в большую сторону, чтобы крайняя точка не оказалась за пределами крайнего столбца.

3. Величину d при этом удобно выбирать так, чтобы она делилась на 2 так, чтобы потом центральный столбец можно было бы поделить пополам для уточнения центра распределения.

8 Промахи и методы их исключения

Одним из условий правомерности статистической выборки является требование ее однородности, т.е. принадлежности всех ее членов к одной и той же генеральной совокупности.

Однако на практике это требование очень часто нарушается. И, если скажем, при обработке вручную еще можно вспомнить как (при каких условиях) были получены “подозрительные” данные, то при автоматической обработке данных необходимы методы исключения “чужих” для данной выборки результатов.

Определение 8.1

Отсчёты, резко отклоняющиеся по своим значениям от большинства других отсчетов принято называть **промахами** и исключать их из выборки.

Важно. Если серия из небольшого числа измерений содержит грубую погрешность — промах, то наличие этого промаха может сильно исказить как среднее значение измеряемой величины, так и границы доверительного интервала. Поэтому из окончательного результата необходимо исключить этот промах.

Обычно промах имеет резко отличающееся от других измерений значение. Однако это отклонение от значений других измерений не дает еще права исключить это измерение как промах, пока не проверено, не является ли это отклонение следствием статистического разброса.

Особую неприятность доставляют отсчеты, которые и не входят в компактную группу отсчетов, но и не удалены от нее на значительное расстояние. Такой отсчет называют предполагаемым промахом. В экспериментальной

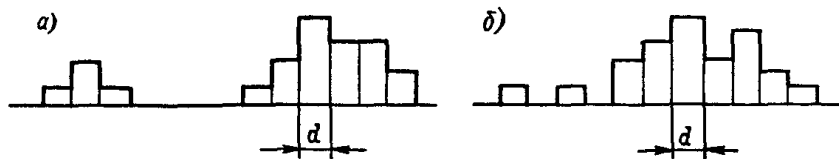


Рисунок 9– Возможные промахи

практике исследователи просто отбрасывали крайние, “слишком удаленные от центра наблюдения”. Эта процедура получила название **цензурирование выборки**.

Однако для принятия решения необходимы какие-либо формальные критерии.

Простейший метод заключается в использовании “правила 3σ ”, когда по выборке с удаленными отсчётами (предполагаемыми промахами) вычисляется оценка σ и граница $|X_{\text{group}}| = 3\sigma$, а все $|x_i| \pm 3\sigma$ отбрасываются.

Правило 3σ обосновано на неравенстве Чебышева:

$$\text{Prob}\{|x_i - \bar{x}| \geq a\} \leq \frac{\sigma^2}{a^2}, \quad (a > 0) \quad (8.1)$$

Если все $x_i \geq 0$, то

$$\text{Prob}(x \geq a) \leq \frac{\bar{x}}{a}, \quad (a > 0) \quad (8.2)$$

Если x имеет одномодальное распределение (непрерывное), то справедлива более сильная оценка

$$\text{Prob}\{|x_i - \bar{x}| \geq a\} \leq \frac{4}{9} \frac{1 + S^2}{\left(\frac{a}{\sigma} - |S_{\text{pear}}|\right)^2}, \quad (8.3)$$

где S_{pear} – пирсоновская мера асимметрии (для распределений симметричных относительно моды $S_{\text{pear}} = 0$).

$$S_{\text{pear}} = \frac{\bar{x} - \xi}{\sigma},$$

где ξ – максимальная вероятность.

$$S_{\text{pear}} = \frac{\gamma_1 (\gamma_2 + 6)}{2 (5\gamma_2 - 6\gamma_1^2 + 6)},$$

где

$$\gamma_1 = \frac{m_3}{m_2^{3/2}}, \quad \gamma_2 = \frac{m_4}{m_2^2} - 3.$$

Значения m_j определяются по формуле:

$$m_j = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^j.$$

Если положить $a = 3\sigma$, то имеем

$$P\{|x_i - \bar{x}| \geq 3\sigma\} \leq \frac{1}{9} \approx 0,11$$

или для одномодальных симметричных:

$$F\{|x_i - \bar{x}| \geq 3\sigma\} \leq \frac{41}{99} = \frac{4}{81} \approx 0,05,$$

т. е. для произвольного распределения вероятность, что $x_i \geq \bar{x}$ на 3σ составляет 11%, а для одномодальных симметричных 5%.

Чувствительность разных статистических методов к наличию аномальных наблюдений (“промахов”) в экспериментальных данных неодинакова.

Критерий Граббса

$$G > \frac{n-1}{\sqrt{n}} \sqrt{\frac{t_{\alpha/(2n), n-2}^2}{n-2 + t_{\alpha/(2n), n-2}^2}} \quad (8.4)$$

Тест Граббса определяется для гипотезы:

- H_0 : в наборе данных нет выбросов
- H_1 : В наборе данных имеется ровно один выброс.

Тест Граббса основан на предположении о нормальности выборки. То есть сначала нужно проверить, что данные могут быть разумно аппроксимированы нормальным распределением перед применением теста Граббса.

Тест Граббса обнаруживает один выброс за раз. Этот выброс исключается из набора данных, и тест повторяется до тех пор, пока не будут обнаружены выбросы. Тем не менее, множественные итерации изменяют вероятности обнаружения, и тест не должен использоваться для размеров выборок в шесть или меньше, поскольку он часто помещает большинство точек в виде выбросов.

Критерий Роснера для обнаружения нескольких выбросов

$$R_i = \max \frac{x_i - \bar{x}}{S} \quad (8.5)$$

или

$$\tau_1 = \max \left\{ \frac{x_1 - \bar{x}}{S}, \frac{x_n - \bar{x}}{S} \right\}, \text{ если } x_1 < x_2 < \dots < x_n. \quad (8.6)$$

Алгоритм критерия Роснера состоит в следующем. По начальной выборке объема n вычисляются значения \bar{x} и S и статистика τ_1 . Затем из выборки удаляется экстремальный член $x_{\min}(x_{\max})$ —в зависимости от того, какое значение более удалено от среднего. Так повторяется k раз.

Полученные значения статистик τ_{ik} ($i = 1, \dots, k$) каждый раз сравниваются с критическими значениями

$$\tau_{i,n,p}^* = \frac{n-i}{\sqrt{n-i+1}} \sqrt{\frac{t_{p,n-i-1}^2}{n+i-1+t_{p,n-i-1}^2}}, \quad p = 1 - \alpha/(2(n-i+1)) \quad (8.7)$$

для заданных n , k и вероятности p . Превышение критерием τ_{1i} критического значения $\tau_{i,n,p}^*$ для некоторого i , позволяет установить не только наличие выбросов, но и их количество (равное значению i , при котором появляется первая значимая величина критерия τ_{1i}). Вычисление последовательных статистик ведется до тех пор, пока $\tau_{1(i+1)} > \tau_{1i}$.

8.1 Другие критерии для исключения промахов

Однако, правило “ 3σ ” хорошо для нормального распределения. Действительно, при $n = 100$ появление $|x_i| \geq 3\sigma$ можно считать промахом, то для равномерного промахом можно считать уже $|x_i| = 1.8\sigma$, а для распределения Лапласа $|x_i| = 3\sigma$ есть отсчет, принадлежащий данной выборке.

На этом примере видно, что граница цензурирования зависит не только от объема выборки n , но также и от формы распределения.

Зависимость от n полуколичественно можно оценить из условия: границы цензурирования должны отсекают в среднем менее одной точки, тогда назначение границ с уровнем значимости $g = 1 - P$, где $P = \frac{n}{n+1}$ дает зависимость от n .

Для расчета количества σ для цензурирования можно использовать аппроксимационные формулы:

Формулы для расчета количества σ ()

$$\begin{aligned} t_{гр.} &= 1,2 + 3,6 (1 - 1/\sqrt{\varepsilon}) \lg \left(\frac{n}{10} \right), \\ t_{гр.} &= 1,55 + 0,8\sqrt{\varepsilon - 1} \lg \left(\frac{n}{10} \right) \end{aligned} \quad (8.8)$$

Тогда интервал цензурирования выглядит следующим образом:

$$\bar{x} \pm t_{гр.} S, \quad (8.9)$$

где S - точечная оценки среднеквадратичного отклонения (3.5).

В заключение данного раздела отметим, что все оценки \bar{x} , S и т.д. должны пересчитываться после цензурирования выборки.

9 Критерии согласия

Постановка задачи.

Одной из задач первичной обработки экспериментальных наблюдений является выбор закона распределения, который наиболее хорошо описывающего случайную величину, выборку которой наблюдают.

Рассмотрим любой эксперимент, в котором измеряется некоторая случайная величина X и мы получаем некоторую выборку объема n . Наша задача описать поведение этой случайной величины. Из теории мы знаем о большом количестве функций распределения вероятностей (нормальное, равномерное распределения, распределение Стьюдента и т.д.).

И у нас возникает мысль: А что если одно из известных теоретических (модельных) распределений подходит для описания поведения измеряемой величины? Как быть? Подойдет ли выбранное нами теоретическое распределение или не подойдет? Такая задача в математической статистике называют проверкой гипотезы.

Определение 9.1

Под статистической гипотезой понимают всякое высказывание о случайной величине (генеральной совокупности), проверяемое по выборке (по результатам наблюдений).

Соответственно, процедура сопоставления высказанной гипотезы с выборочными данными называется **проверкой гипотезы**.

Таким образом, **целью проверки гипотезы о согласии** опытного распределения с теоретическим является стремление удостовериться в том, что данная модель теоретического закона не противоречит наблюдаемым данным и использование ее не приведет к существенным ошибкам при вероятностных расчетах.

Ответ, который должен быть получен в конечном итоге выглядит следующим образом: Данная (“экспериментальная”) выборка с вероятностью p соответствует теоретическому распределению с параметрами (например, нормальному распределению с параметрами $a = 3$ и $\sigma = 1$).

Поэтому основная задача на практике состоит в определении вероятности принятия модели (гипотезы) p (или более точно вероятности, что данная выборка не соответствует данной гипотезе $\alpha = 1 - p$).

Что нужно знать о процессе проверки гипотез:

Этап 1. Выбор модели (теоретической функции распределения)

Располагая выборочными данными и руководствуясь конкретными условиями рассматриваемой задачи, формулируют гипотезу H_0 , которую называют **основной** или нулевой, и гипотезу H_1 конкурирующую с гипотезой H_0 .

Термин “конкурирующая” означает, что являются противоположными следующие два события: по выборке будет принято решение о справедливости для генеральной совокупности гипотезы H_0 ; и по выборке будет принято решение о справедливости для генеральной совокупности гипотезы H_1 . Гипотезу H_1 называют также альтернативной.

Например, если нулевая гипотеза такова: данная выборка имеет нормальное распределение с параметрами $a = 3$ и $\sigma = 1$, то альтернативная гипотеза может быть следующей: данная выборка имеет нормальное распределение, но с параметрами $a < 3$ и $\sigma = 1$ или

$$H_0 : p(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left[-\frac{(x-a)^2}{2\sigma^2} \right], \text{ где} \\ a = 3 \text{ и } \sigma = 1, \quad (9.1)$$

и

$$H_1 : p(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left[-\frac{(x-a)^2}{2\sigma^2} \right], \text{ где} \\ a < 3 \text{ и } \sigma = 1. \quad (9.2)$$

Этап 2. Расчет вероятности принятия гипотезы p или вероятности не соответствия данной модели $\alpha = 1 - p$.

Вероятность α часто называют уровнем значимости и если выборка соответствует выбранной нами модели (теоретическому распределению), то уровень значимости относительно мал ($\alpha < 0,05$) Поясним ее смысл.

Решение о том, можно ли считать высказывание H_0 справедливым для генеральной совокупности, принимается по выборочным данным, т. е. по ограниченному ряду наблюдений, следовательно, это решение может быть ошибочным.

При этом может иметь место ошибка двух родов: отвергают гипотезу H_0 , или, иначе, принимают альтернативную гипотезу H_1 , тогда как на самом деле гипотеза H_0 верна; это **ошибка первого рода**; принимают гипотезу H_0 , тогда как на самом деле высказывание H_0 неверно, т. е. верной является гипотеза H_1 это **ошибка второго рода**.

Так вот уровень значимости α – это вероятность ошибки первого рода, т. е. вероятность того, что будет принята гипотеза H_1 , если на самом деле в

генеральной совокупности верна гипотеза H_0 (или отклонение проверяемой гипотезы H_0 при её справедливости).

Вероятность ошибки второго рода часто обозначают β , т. е. вероятность того, что будет принята гипотеза H_0 , если на самом деле верна гипотеза H_1 (или не отклонении H_0 при справедливости H_1).

О мощности критерия.

При использовании критериев согласия, как правило на практике, не задают конкурирующих гипотез: рассматривается принадлежность выборки конкретному закону. А в качестве конкурирующей гипотезы — принадлежность любому другому.

Естественно, что способность критерия отличать закон, соответствующий H_0 , от других, близких к закону, соответствующему H_0 , и далёких от него, отличаются.

Определение 9.2

Мощностью критерия по отношению к конкурирующей гипотезе H_1 называется величина $1 - \beta$. Критерий тем лучше распознаёт пару конкурирующих гипотез H_0 и H_1 , чем выше его мощность.

Один из возможных вариантов проверки гипотез являются **критерии согласия**.

Критерий χ^2 (критерий Пирсона).

Рассмотрим этапы необходимые для проверки гипотез на примере критерия χ^2 .

Процедура проверки гипотез с использованием критериев типа χ^2 предусматривает группирование наблюдений.

1. Выбираем “модель”: обычно это теоретическое дифференциальное распределение вероятностей $p(x, \theta)$, где θ — один (или несколько) параметров распределения.
2. Область определения случайной величины разбивают на m непересекающихся интервалов граничными точками $x_0, x_1, \dots, x_{m-1}, x_m$, где $x_0 < x_1 < \dots < x_{m-1} < x_m$.
3. В соответствии с заданным разбиением подсчитывают число n_i выборочных значений, попавших в i -й интервал и вероятности попадания в i -й

$$\begin{aligned}
 p_i &= \int_{x_{i-1}}^{x_i} p(x, \theta) dx = \\
 &= F(x_i) - F(x_{i-1})
 \end{aligned}
 \tag{9.3}$$

соответствующие теоретическому закону с интегральной функцией распределения $F(x, \theta)$ для всех m интервалов.

При этом $\sum_{i=1}^m n_i = n$ и $\sum_{i=1}^m p_i = 1$.

4. Проводим расчет статистики критерия согласия χ^2 Пирсона с помощью соотношения

$$\chi_{\text{exp}}^2 = n \sum_{i=1}^m \frac{(n_i/n - p_i)^2}{p_i}
 \tag{9.4}$$

При проверке гипотезы при $n \rightarrow \infty$ для которой известны, как вид закона $p(x, \theta)$, так и все его параметры θ (простая гипотеза) функция χ_{exp}^2 подчиняется распределению χ_r^2 с $r = m - 1$ степенями свободы (доказано в Pearson, Karl (1900). On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling. Philosophical Magazine Series 5 50 (302): 157–175.).

5. Далее находим из уравнения величину

$$\begin{aligned}
 \int_{\chi_{\text{exp}}^2}^{\infty} p(s) ds &= \\
 \int_{\chi_{\text{exp}}^2}^{\infty} \frac{(s)^{\frac{r}{2}-1} e^{-\frac{s}{2}}}{2^{\frac{r}{2}} \Gamma(\frac{r}{2})} ds &= p, \text{ или} \\
 1 - F_r(\chi_{\text{exp}}^2) &= p,
 \end{aligned}
 \tag{9.5}$$

где $F_r(x)$ – интегральная функция вероятности распределения χ_r^2 с r степенями свободы.

6. После вычисления p получаем ответ: **Данная выборка с вероятностью p соответствует теоретическому распределению $p(x, \theta)$ с параметрами θ .**

Или: Данная выборка с вероятностью $\alpha = 1 - p$ не соответствует теоретическому распределению $p(x, \theta)$ с параметрами θ

Примечания.

- ★ На практике (например с помощью точечных оценок) удастся оценить (рассчитать) все или часть параметров распределения. Тогда статистика χ_{exp}^2 при справедливости проверяемой гипотезы подчиняется χ_r^2 -распределению с $r = m - k - 1$ степенями свободы, где k количество оцененных по выборке параметров.
- ★ Некорректное использование критериев согласия (не построен вариационный ряд для выборки, неверно выбрано число интервалов) может приводить к необоснованному принятию (чаще всего) или необоснованному отклонению проверяемой гипотезы.
- ★ Существуют и другие критерия согласия : критерий Колмогорова-Смирнова, критерий Мизеса и т.д.

Существует несколько более удобный способ понимания критерия χ^2 . Введем понятие приведенного значения $\tilde{\chi}^2$ (или $\tilde{\chi}^2$ на одну степень свободы), которое определим как

$$\tilde{\chi}^2 = \frac{\chi^2}{r} \quad (9.6)$$

Тогда, каким бы ни было число степеней свободы, наш критерий можно сформулировать следующим образом: если мы получаем значение $\tilde{\chi}^2$ порядка 1 или меньше, то у нас нет оснований сомневаться в нашем ожидаемом распределении; если мы получаем значение $\tilde{\chi}^2$ много большее, чем единица, то невероятно, чтобы наше ожидаемое распределение было верным. Т.е. если

$$\tilde{\chi}^2 \leq 1, \quad (9.7)$$

то наша выборка соответствует теоретическому распределению.

10 Алгоритм предварительной обработки экспериментальных данных

Исходя из изложенного материала можно сформировать необходимые этапы предварительной обработки наблюдений.

1. **Вычисление выборочных характеристик (точечные оценки моментов экспериментального распределения) \bar{x} , S^2 , $\tilde{\nu}_3$ (и соответственно $\tilde{\gamma}_1$), $\tilde{\nu}_4$ (и соответственно $\tilde{\varepsilon}$), а также дополнительные характеристики $S_{\bar{x}}$, и $m_{3,4}$ по формулам:**

$$\begin{aligned}\bar{x} &= \frac{1}{n} \sum_{i=1}^n x_i, \\ S^2 &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2, \\ m_k &= \sum_{i=1}^n (x_i - \bar{x})^k, \quad k = 2, 3, 4, \\ \tilde{\nu}_3 &= \frac{n^2}{(n-1)(n-2)} m_3, \\ \tilde{\gamma}_1 &= \frac{\tilde{\nu}_3}{S^3}, \\ \tilde{\nu}_4 &= \frac{(n^2 - 2n + 3) m_4 - 3n(2n-3) m_3^2}{(n-1)(n-2)(n-3)}, \\ S_{\bar{x}} &= \frac{S}{\sqrt{n}}.\end{aligned}\tag{10.1}$$

(Смотри формулы (3.4), (3.5), (3.8), (3.9)) и другие.

2. **Цензурирование выборки (отсев грубых промахов).** Последующий пересчёт точечных оценок.
3. **Построение гистограмм и кумулятивной кривой.** Визуальный анализ на соответствие теоретической плотности распределения с привлечением точечных оценок.
4. **Проверка с помощью критериев согласия на соответствие одной (или нескольких) из теоретических плотностей распределения.** Как правило, проверяют соответствие экспериментальной выборки нормальному распределению.

11 Некоторые сведения о двумерных случайных величинах

Переход от одномерных случайных величин к многомерным сопровождается введением новых понятий. Наглядные модельные представления можно построить для двух или трех переменных. Наиболее удобен случай двух переменных, и он чаще других будет использоваться в дальнейшем.

Пусть X и Y две случайные величины, имеющие набор значений x_k и y_k . Индекс k нумерует определенные значения этих величин. Для каждой из одномерных случайных величин можно ввести интегральные $F(x)$, $F(y)$ и дифференциальные функции распределения вероятностей $p(x)$ и $p(y)$. Но можно ввести и совместную функцию распределения вероятностей, описывающую поведение обеих случайных величин, как единый объект (одновременно).

Определение 11.1

Совместная функция распределения $F(x,y)$ - это вероятность, того что значения двумерной величины приписанные выборочному множеству k , удовлетворяют одновременно неравенствам $x(k) \leq x$ и $y(k) \leq y$ или

$$F(x,y) = \text{Prob}(x_k \leq x \text{ и } y_k \leq y) . \quad (11.1)$$

Легко найти, что

1.

$$F(-\infty, y) = 0 , \quad (11.2)$$

2.

$$F(x, -\infty) = 0 , \quad (11.3)$$

3.

$$F(\infty, -\infty) = 1 . \quad (11.4)$$

Совместный дифференциальный закон распределения $p(x,y)$ определяется соотношением:

$$p(x,y) = \frac{\partial^2}{\partial y \partial x} [F(x,y)] \quad (11.5)$$

Очевидно что

1.

$$p(x,y) \geq 0 , \quad (11.6)$$

$$F(x, y) = \int_{-\infty}^y \int_{-\infty}^x p(\xi, \eta) d\xi d\eta, \quad (11.7)$$

Плотность вероятности для случайных величин X и Y в отдельности, выражаются через совместную функцию распределений с помощью соотношений:

$$p(x) = \int_{-\infty}^{\infty} p(x, y) dy, \quad (11.8)$$

$$p(y) = \int_{-\infty}^{\infty} p(x, y) dx. \quad (11.9)$$

Функции (11.8) и (11.9) называют безусловными или маргинальными плотностями распределения случайных величин X и Y .

Определение 11.2

Две случайные величины X и Y являются статистически независимыми (или взаимно независимыми), если

$$p(x, y) = p(x) p(y). \quad (11.10)$$

Математическое ожидание

Если имеется функция случайных величин $G(x, y)$, то ее математическое ожидание

$$E \{G(x, y)\} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} G(x, y) p(x, y) dx dy. \quad (11.11)$$

Дисперсия

Дисперсия $D \{G(x, y)\}$ функция случайных величин $G(x, y)$ выражается следующим образом:

$$D \{G(x, y)\} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (G(x, y) - E \{G(x, y)\})^2 p(x, y) dx dy. \quad (11.12)$$

11.1 Алгебраические и центральные моменты

Алгебраические имеют соответственно следующий вид (и частном случае для (11.11), когда $G = x$ или $G = y$):

$$E \{x\} = \mu_x = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xp(x,y) dx dy , \quad (11.13)$$

$$E \{y\} = \mu_y = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} yp(x,y) dx dy . \quad (11.14)$$

My title

С учетом определений (11.8) и (11.9) для безусловных функций распределения имеем, что

$$E \{x\} = \mu_x = \int_{-\infty}^{\infty} xp(x) dx , \quad (11.15)$$

$$E \{y\} = \mu_y = \int_{-\infty}^{\infty} yp(y) dy . \quad (11.16)$$

Центральные моменты, по аналогии с одномерной случайной величиной, могут быть записаны в виде

$$D \{x\} = \sigma_x^2 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - E \{x\})^2 p(x,y) dx dy , \quad (11.17)$$

$$D \{y\} = \sigma_y^2 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (y - E \{y\})^2 p(x,y) dx dy . \quad (11.18)$$

Но кроме этих простейших моментов второго порядка (11.17) и (11.18), имеющих уже хорошо известную форму дисперсий, появляются новые возможности для образования моментов второго порядка.

Ковариация

Так можно определить дисперсию $D\{x,y\}$ посредством соотношения:

$$D\{x,y\} = \text{cov}(x,y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - E\{x\})(y - E\{y\})p(x,y) dx dy . \quad (11.19)$$

Определение 11.3

Центральный момент (11.19) называются **ковариацией** между x и y .

Если x и y статистически независимы (см. (11.10)), т. е. $p(x,y) = p(x)p(y)$, тогда ковариация обращается в нуль, Действительно, после очевидных преобразований, имеем, что

$$\begin{aligned} \text{cov}(x,y) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - E\{x\})(y - E\{y\})p(x)p(y) dx dy = \\ &= \int_{-\infty}^{\infty} p(x)(x - E\{x\}) dx \times \int_{-\infty}^{\infty} (y - E\{y\})p(y) dy = 0 . \end{aligned} \quad (11.20)$$

Во всех других случаях смешанный момент второго порядка не равен нулю. Обычно вводится безразмерная переменная - **коэффициент корреляции** $\rho(x,y)$, определяемый как

$$\rho(x,y) = \frac{\text{cov}(x,y)}{\sigma_x \sigma_y} . \quad (11.21)$$

Можно, показать, что

$$-1 \leq \rho(x,y) \leq 1 . \quad (11.22)$$

На практике бывает необходимо оценить коэффициент корреляции $\rho(x,y)$ для конкретной выборки объема n . Выборочная ковариация (несмещенная точечная оценка) $r(x,y)$ можно определить по формуле

$$r(x,y) = \frac{1}{n-1} \sum_{i=1}^n \frac{(x_i - \bar{x})(y_i - \bar{y})}{S_x S_y} p_i . \quad (11.23)$$

Здесь \bar{x} и \bar{y} - средние значения случайных переменных; p_i - частота попадания переменных x_i и y_i в i -ый интервал значений, а $S_{x,y}$ - точечные оценки среднеквадратичных отклонений $\sigma_{x,y}$.

Определение 11.4

Ковариация (11.19) имеет смешанный характер и вместе с остальными моментами образует матрицу (11.17) и (11.18), которую называют **ковариационной матрицей**, или матрицей вторых моментов, или дисперсионной матрицей, или матрицей ошибок.

Ковариационная матрица

$$\mathbf{V} = \begin{pmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{pmatrix} = \begin{pmatrix} D\{x\} & D\{x,y\} \\ D\{y,x\} & D\{y\} \end{pmatrix} \quad (11.24)$$

Поскольку $D\{x,y\} = D\{y,x\}$, то матрица корреляции симметрична, диагональные элементы которой представляют собой дисперсии.

Каждый элемент ковариационной матрицы V_{ij} можно трактовать как математическое ожидание ij -элемента произведения вектор-столбца $(\mathbf{x} - \boldsymbol{\mu})$ на вектор-строку $(\mathbf{x} - \boldsymbol{\mu})^T$, где

$$\begin{aligned} (\mathbf{x} - \boldsymbol{\mu}) &= \begin{pmatrix} x - \mu_x \\ y - \mu_y \end{pmatrix} \\ (\mathbf{x} - \boldsymbol{\mu})^T &= (x - \mu_x \quad y - \mu_y) \end{aligned} \quad (11.25)$$

Теперь в матричной записи (11.24) имеет вид

$$\mathbf{V} = E \left\{ (\mathbf{x} - \boldsymbol{\mu}) (\mathbf{x} - \boldsymbol{\mu})^T \right\} . \quad (11.26)$$

Корреляционные связи имеют важное значение в исследованиях. Удобство работы с корреляциями связано с тем, что корреляция - безразмерная величина в отличие от ковариации. Из коэффициентов корреляции можно также построить корреляционную матрицу \mathbf{R} . В двумерном случае, в следствие, того, что

$$\rho(x,x) = \rho(y,y) = 1 \quad (11.27)$$

ИМЕЕМ, ЧТО

$$\begin{aligned} \mathbf{R} &= \begin{pmatrix} \rho(x,x) & \rho(x,y) \\ \rho(x,y) & \rho(y,y) \end{pmatrix} = \\ &= \begin{pmatrix} 1 & \frac{\text{cov}(x,y)}{\sigma_x\sigma_y} \\ \frac{\text{cov}(x,y)}{\sigma_x\sigma_y} & 1 \end{pmatrix}. \end{aligned} \quad (11.28)$$

Важно, что если x и y образуют двумерную случайную величину, то дисперсия линейной функции $z = x + y$ находится с помощью соотношения:

$$D\{x + y\} = D\{x\} + D\{y\} + 2\text{cov}(x,y). \quad (11.29)$$

11.2 Примеры многомерных функций распределения вероятностей

Пусть имеем многомерную случайную величину $\{X_1, X_2, \dots, X_k\}$ размерности k . Для центральных и алгебраических моментов распределения, а также ковариационной матрицы введем следующие обозначения:

$$\begin{aligned} E\{X_i\} &= \mu_i, i = 1, \dots, k, \\ D\{X_i\} &= \sigma_i, \\ \text{cov}(X_i, X_j) &= \rho^{ij}\sigma_i\sigma_j. \end{aligned} \quad (11.30)$$

Введем дополнительные k -мерные векторы: $\mathbf{X} = \{X_1, X_2, \dots, X_k\}$ и $\boldsymbol{\mu} = \{\mu_1, \dots, \mu_k\}$ и матрицу ошибок

$$\mathbf{V}_{ij} = \begin{cases} \sigma_i^2, & \text{если } i = j \\ \rho_{ij}\sigma_i\sigma_j & \text{если } i \neq j \end{cases}. \quad (11.31)$$

Многомерное нормальное распределение

Тогда функция распределения

$$p(\mathbf{X}) = \frac{1}{(2\pi)^{k/2} |\mathbf{V}|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{X} - \boldsymbol{\mu})^T \mathbf{V}^{-1} (\mathbf{X} - \boldsymbol{\mu}) \right], \quad (11.32)$$

где $|\mathbf{V}|$ -детерминант ковариационной матрицы \mathbf{V} определяет k -мерную плотность нормального распределения.

Наиболее наглядным вариантом является двумерное нормальное распределение, которое определяется параметрами $\mu_1, \mu_2, \sigma_1, \sigma_2$ и ρ_{12} :

$$p_N(x, y) = \frac{1}{2\pi \sqrt{(1 - \rho_{12}^2) \sigma_1^2 \sigma_2^2}} \times \exp \left\{ \frac{\sigma_2^2 (x - \mu_1)^2 - 2\rho_{12}\sigma_2\sigma_1 (x - \mu_1)(y - \mu_2) + \sigma_1^2 (y - \mu_2)^2}{2(\rho_{12}^2 - 1) \sigma_1^2 \sigma_2^2} \right\} \quad (11.33)$$

Если $\mu_1 = \mu_2 = 0$ и $\sigma_1 = \sigma_2 = 1$, а $\rho_{12} = \rho$, то из (11.33) получим стандартизованное двумерное нормальное распределение:

$$p_N(x, y) = \frac{1}{2\pi \sqrt{1 - \rho^2}} \exp \left\{ \frac{x^2 - 2\rho xy + y^2}{2(\rho^2 - 1)} \right\}. \quad (11.34)$$

График стандартизованного двухмерного нормального распределения с $\rho = 0$ (отсутствуют корреляции между x и y), представлен на рисунке 11.2.

Естественно, существуют и другие теоретические функции распределения многомерных случайных величин.

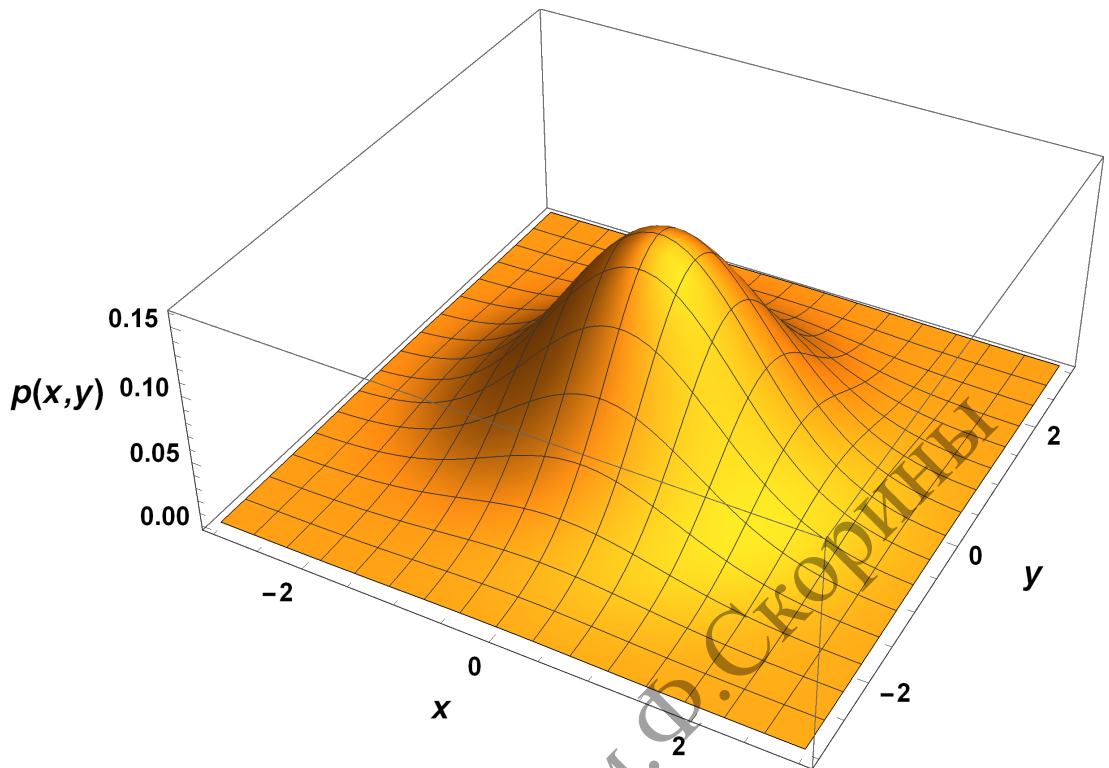


Рисунок 10– Двумерное нормальное распределение с $\mu_1 = \mu_2 = 0$ и $\sigma_1 = \sigma_2 = 1$, а $\rho = 0$

12 Корреляционный и регрессионный анализ

Статистические связи между случайными величинами можно изучать методами корреляционного и регрессионного анализа. Основной целью регрессионного анализа является установление формы и изучение зависимости между переменными.

Целью корреляционного анализа является определить степень взаимосвязи между исследуемыми случайными величинами. В случае, если исследуется связь двух переменных, корреляционный анализ называют **парным**; если число переменных более двух — **множественным**.

Статистическая зависимость между переменными, при которой каждому значению одной или нескольких случайных величин соответствует определенное среднее значение другой, называется **корреляционной**.

Корреляционный и регрессионный анализ решает две основные задачи:

1. Определение уравнения (или системы уравнений) связи, т.е. в установлении математической формы, в которой выражается данная связь. Это очень важно, так как от правильного выбора формы связи зависит конечный результат изучения взаимосвязи между признаками.

Viktor Andreev Viktor Andreev Viktor Andreev
Viktor Andreev Viktor Andreev Viktor Andreev
Viktor Andreev Viktor Andreev Viktor Andreev
Viktor Andreev Viktor Andreev Viktor Andreev

2. Вычисление степени взаимосвязи, т.е. меры связи между признаками с целью установить степень влияния данных случайных величин (факторов) на конечную случайную величину (результат). Эта задача решается путем определения параметров корреляционного уравнения различными математическими методами.

Далее проводятся оценка и анализ полученных результатов при помощи различных показателей регрессионно-корреляционного метода (коэффициентов детерминации, линейной и множественной корреляции и т.д.), а также проверка существенности связи между изучаемыми признаками.

В результате использование корреляционного и регрессионного анализа можно решить такие задачи, как

- Взаимосвязь. Есть ли взаимосвязь между случайными величинами, и насколько велико их влияние на конечную величину.
- Прогнозирование. Если известно поведение одного параметра, то можно предсказать поведение другого параметра, коррелирующего с первым в области, где нет экспериментальной информации.
- Обнаружение неизвестных причинных связей.

Более часто используемым показателем степени тесноты корреляционной связи является линейный коэффициент корреляции. При расчете этого показателя учитываются не только отклонения индивидуальных значений случайной величины от среднего значения, но и сама величина этих отклонений.

В самом общем случае, при большом числе наблюдений одно и то же значение случайной величины X может встретиться n_x раз, одно и то же значение случайной величины Y может встретиться n_y раз, а одна и та же пара чисел $\{x, y\}$ может наблюдаться n_{xy} раз. Поэтому данные наблюдений группируют, т.е. подсчитывают частоты n_x , n_y и n_{xy} . Все сгруппированные данные записывают в виде таблицы, которую называют корреляционной.

Таблица 1– Пример корреляционной таблицы

$y \backslash x$	1	2	3	4	5	n_y
0	–	–	–	6	4	10
1	–	–	1	4	6	11
2	–	5	9	5	–	19
3	3	7	–	–	–	10
n_x	3	12	10	15	10	50

В первой строке указаны наблюдаемые значения $\{1, 2, 3, 4, 5\}$ случайной величины X , а в первом столбце таблицы – наблюдаемые значения $\{0, 1, 2, 3\}$ случайной величины Y .

На пересечении строк и столбцов находятся частоты n_{xy} наблюдаемых пар значений случайных величин X и Y . Например, частота 7 указывает, что пара чисел $\{2, 3\}$ наблюдалась 7 раз. Все частоты помещены в прямоугольнике. В последнем столбце записаны суммы частот строк. В последней строке записаны суммы частот столбцов. Общее число наблюдений $n = 50$.

Если имеем k значений случайной величины X и m значений случайной величины Y в корреляционной таблице с числом попаданий n_i , n_j и n_{ij} для x_i , y_j и пар $\{x_i, y_j\}$ соответственно, то оценку коэффициент корреляции можно в данном случае можно найти из выражения:

$$r_{xy} = \frac{n \sum_{i=1}^k \sum_{j=1}^m x_i y_j n_{ij} - \left(\sum_{i=1}^k x_i n_i \right) \left(\sum_{j=1}^m y_j n_j \right)}{\sqrt{n \sum_{i=1}^k x_i^2 n_i - \left(\sum_{i=1}^k x_i n_i \right)^2} \sqrt{n \sum_{j=1}^m y_j^2 n_j - \left(\sum_{j=1}^m y_j n_j \right)^2}}. \quad (12.1)$$

В упрощенном варианте, когда имеется n значений пар $\{x_i, y_i\}$ двумерной случайной величины, выражение коэффициент корреляции (12.1) упрощается и приобретает вид:

$$r_{xy} = \frac{n \sum_{i=1}^n x_i y_i - \left(\sum_{i=1}^n x_i \right) \left(\sum_{j=1}^n y_j \right)}{\sqrt{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2} \sqrt{n \sum_{j=1}^n y_j^2 - \left(\sum_{j=1}^n y_j \right)^2}}. \quad (12.2)$$

Последний вариант часто встречается в физических исследованиях.

Иногда одному значению x_i соответствует m значений $\{\tilde{y}_{i1}, \tilde{y}_{i2}, \dots, \tilde{y}_{i,m}\}$ величины Y . Тогда, после нахождения среднего значения и среднеквадратичного отклонения по известным формулам:

$$y_i = \frac{1}{m} \sum_{j=1}^m \tilde{y}_{ij},$$

$$s_{y_i} = \sqrt{\frac{1}{m-1} \sum_{j=1}^m (y_i - \tilde{y}_{ij})^2}, \quad (12.3)$$

наборам пар $\{x_i, y_i\}$ с дополнительным набором ошибок $\{s_{y_1}, \dots, s_{y_n}\}$.

Для определения степени связи используют обычно таблицу Чеглока (смотри таблицу 2).

Таблица 2– Пример корреляционной таблицы

Значение $ r_{xy} $	Степень связи
0,1 – 0,3	слабая
0,3 – 0,5	умеренная
0,5 – 0,7	заметная
0,7 – 0,9	высокая
0,9 – 0,99	весьма высокая связь

Репозиторий ГГУ им. Ф. Скорины

13 Линейный регрессионный анализ

Корреляционный анализ позволяет установить степень взаимосвязи двух и более случайных величин. Однако наряду с этим желательно иметь модель этой связи, которая дала бы возможность предсказывать значения одной случайной величины по конкретным значениям другой. Методы решения подобных задач составляют раздела математической статистики “регрессионный анализ”.

Для упрощения изложения рассмотрим случай двух случайных величин x и y . Линейная связь между двумя случайными величинами означает, что прогноз значения y по данному значению x имеет вид:

$$\hat{y} = A + Bx \quad (13.1)$$

Если данные связаны идеальной линейной зависимостью $|r_{xy}| = 1$, то предсказанное значение \tilde{y}_i будет соответствовать наблюдаемому значению y_i при любом данном x_i . На практике идеальная зависимость отсутствует (за счет случайных разбросов данных, за счет нелинейных эффектов).

Однако, если предположить наличие линейной связи, то можно подобрать A и B , которые дадут возможность предсказывать ожидаемое значение y_i для любого данного x_i (при этом предсказанное \tilde{y}_i не совпадает с наблюдаемыми y_i , однако оно будет равно среднему значению всех таких наблюдаемых значений).

Наиболее общепринятая процедура определения коэффициентов A и B состоит в таком выборе, который минимизирует сумму квадратов отклонений наблюдаемых значений от предсказанного значения y_i .

Этот метод называют методом наименьших квадратов (МНК). Он разработан в 1795-1805 гг. Лежандром и Гауссом.

13.1 Метод наименьших квадратов

Рассчитывается отклонение $y_i - \hat{y}_i = \Delta_i$, а затем находятся коэффициенты так, что

$$Q = \sum_{i=1}^n \Delta_i^2 \rightarrow \min \quad (13.2)$$

Для этого надо взять частные производные функции Q (13.2) по параметрам A и B и приравнять их к нулю.

В итоге имеем, что

$$\begin{cases} \frac{\partial Q}{\partial A} = 0, \\ \frac{\partial Q}{\partial B} = 0 \end{cases} . \quad (13.3)$$

Из этой системы получаем оценки пар A и B .

Проделав цепочку преобразований

$$\frac{\partial Q}{\partial A} = -2 \sum_{i=1}^n (y_i - A - Bx_i) = 0 . \quad (13.4)$$

$$An + b \sum_{i=1}^n x_i = \sum_{i=1}^n y_i ,$$

$$\frac{\partial Q}{\partial B} = -2 \sum_{i=1}^n (y_i - A - Bx_i) x_i = 0 ,$$

$$- \sum_{i=1}^n y_i x_i + A \sum_{i=1}^n x_i + B \sum_{i=1}^n x_i^2 = 0 ,$$

$$An + B \left(\sum_{i=1}^n x_i \right) = \left(\sum_{i=1}^n y_i \right) ,$$

$$\left(\sum_{i=1}^n x_i \right) A + B \left(\sum_{i=1}^n x_i^2 \right) = \left(\sum_{i=1}^n x_i y_i \right)$$

в итоге имеем:

$$\begin{cases} A = \frac{\sum_{i=1}^n x_i^2 \sum_{i=1}^n y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2} , \\ B = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2} . \end{cases} \quad (13.5)$$

Эти оценки можно использовать для модели $\hat{y} = A + Bx$, которую называют **прямой линейной регрессии y на x** .

13.2 Свойства метода наименьших квадратов

Из первого уравнения имеем:

$$A + B \sum_{i=1}^n \frac{x_i}{n} = \sum_{i=1}^n \frac{y_i}{n} = A + B\bar{x} = \bar{y}, \quad (13.6)$$

т.е. кривая регрессии проходит через центр тяжести экспериментальных точек.

Если теперь до проведения МНК все исходные данные отцентрировать, т.е. \bar{x} и \bar{y} перенести в начало координат $\bar{x}=\bar{y}=0$, то первое уравнение превратится в тождество, т.е. система уравнений сокращается, что особенно важно, когда имеет место ситуация с большим числом коэффициентов.

Уравнение прямой можно записать:

$$\hat{y} = A + Bx = \bar{y} + B(x - \bar{x}). \quad (13.7)$$

Вторая особенность МНК состоит в том, что полученные этим методом оценки необратимы, т.е. если имеется модель регрессии y на x : $y = A + Bx$, то регрессию x на y нельзя обратить:

$$x = \frac{(y - A)}{B} = \frac{1}{b}y - A. \quad (13.8)$$

Коэффициент наклона при y обозначим как B' . Тогда имеем, что

$$x = B'y - A. \quad (13.9)$$

Можно найти по аналогии с вышеизложенным, что

$$B' = \frac{\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}}{\sum_{i=1}^n y_i^2 - n\bar{y}^2}. \quad (13.10)$$

Коэффициент (13.10) связан с B через коэффициент корреляции r_{xy} соотношением

$$\sqrt{B'B} = r_{xy}. \quad (13.11)$$

13.3 Точность оценок A и B .

Пусть $E\{A\} = a$, где a – “истинное” значение. Тогда доверительные интервалы для коэффициентов линейной регрессии A и B запишутся в виде:

$$\begin{aligned} a - A &= \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)^{\frac{1}{2}} S_{y/x} t_{n-2, \alpha/2} , \\ b - B &= \left(\frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)^{\frac{1}{2}} S_{y/x} t_{n-2, \alpha/2} , \end{aligned} \quad (13.12)$$

где

$$S_{y/x} = \left[\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2} \right]^{\frac{1}{2}} = \left[\left(\frac{n-1}{n-2} \right) S_y^2 (1 - r_{xy}^2) \right]^{\frac{1}{2}} , \quad (13.13)$$

а $t_{n-2, \alpha/2}$ - квантили распределения Стьюдента.

13.4 Редукция некоторых задач нелинейного регрессионного анализа к линейному

Как сводить некоторые задачи регрессионного анализа к линейному регрессионному анализу.

Например, модель типа

$$\frac{1}{y} = a_1 + b_1 x$$

путем замены переменных $\frac{1}{y} = y_1$ имеем $y_1 = a_1 + b_1 x$.

Аналогично

$$y = a_1 + \frac{b_1}{x} \rightarrow x ,$$

$$\frac{y}{x} = a_1 + b_1 x ,$$

$$y = ax^b \lg y = \lg a + b \lg x ,$$

$$\lg y \rightarrow y \lg x \rightarrow x ,$$

$$y = ab^x \lg y = \lg a + x \lg b .$$

Такие преобразования, строго говоря, допустимы, если исходные величины измерены точно. Например: $y = Ax^B$ измерен с Δy : $y = Ax^B + \Delta y$.

Тогда линеаризированная форма будет отличаться от исходного за счет неясности преобразования Δy . Однако в первичном анализе Δy можно пренебречь, если это необоснованно, то вводим переменную y' путем

$$y' = y - \Delta y = Ax^B .$$

Метод МНК достаточно чувствителен к неоднородностям статистики. Наличие неоднородной статистики приводит иногда к абсурдным результатам, поэтому окончательное решение об МНК должно приниматься по однородной, очищенной статистике.

Репозиторий ГГУ им. Ф. Скорины

14 Элементы нелинейного регрессионного анализа

Пусть одному значению x_i из выборки объема n соответствует m значений $\{\tilde{y}_{i1}, \tilde{y}_{i2}, \dots, \tilde{y}_{i,k}\}$ величины Y . Тогда, после нахождения среднего значения и среднеквадратичного отклонения, получим набор пар $\{x_i, y_i\}$ с дополнительным набором ошибок $\{\sigma_{y_1}, \dots, \sigma_{y_n}\}$. Такая ситуация возникает, например, при построении гистограмм.

Цель регрессионного анализа найти уравнение взаимосвязи случайных величин Y и X (в данном случае, парная регрессионная модель) вида: $y = \phi(x; \theta)$. В силу воздействия неучтенных случайных факторов отдельные значения y_i будут в большей или меньшей мере отклоняться от функции регрессии $\phi(x; \theta)$. В этом случае уравнение взаимосвязи двух переменных (парная регрессионная модель) может быть представлено в виде:

$$y = \phi(x; \theta) + \epsilon, \quad (14.1)$$

где $\epsilon = \{\epsilon_1, \dots, \epsilon_n\}$ случайная величина, характеризующая отклонение от функции регрессии. Эту переменную часто называют возмущением; величина $\theta = \{\theta_1, \theta_2, \dots, \theta_k\}$ определяет набор параметров регрессионной модели, которые необходимо определить.

Основные положения регрессионного анализа:

Основные положения регрессионного анализа

1. Одним из основных требований регрессионного анализа является, условие равенства нулю математического ожидания “возмущения” ϵ :

$$E \{\epsilon_i\} = 0. \quad (14.2)$$

2. Также предполагают, что случайные величины ϵ_i подчиняются одномерному нормальному распределению, если все y_i и y_j не коррелируют с другом или многомерному нормальному распределению, если такая корреляция имеет место.

3. Дисперсия “возмущения” ϵ_i (или зависимой переменной y для любого i) задается соотношением:

$$D \{\epsilon_i\} = \sigma_i^2. \quad (14.3)$$

Наиболее простой вариант (14.3) состоит в требовании того, чтобы все дисперсии отклонений были одинаковые.

В случае постоянства дисперсии и отсутствия корреляции, оценки параметров регрессии полученные, например, с помощью метода наименьших квадратов, обладают важными свойствами, а именно:

- несмещенность;
- состоятельность;
- эффективность.

В случае нелинейной зависимости, свойства постоянства дисперсии и отсутствия корреляции могут не выполняться, тогда полученные оценки параметров регрессии не будут обладать указанными характеристиками. Или связь между переменными x и y линейна, но на исследуемый показатель воздействует фактор, не включенный в модель.

Для решения этой задачи в случае линейной регрессии использовался метод наименьших квадратов. Рассмотрим некоторые методы поиска оптимальной оценки $\hat{\theta}$ в случае нелинейной регрессии. Естественно, эти методы могут быть использованы и в линейном варианте функциональной зависимости.

14.1 Метод максимального правдоподобия

Пусть имеется выборка из n независимых результатов наблюдений $x = \{x_1, x_2, \dots, x_n\}$ с плотностью вероятности $p(x, \theta)$ (θ - набор параметров).

Определение 14.1

Совместная функция вероятности

$$L(x_1, x_2, \dots, x_n, \theta) = p(x_1, \theta) p(x_2, \theta) \cdots p(x_n, \theta) = \prod_{i=1}^n p(x_i, \theta) \quad (14.4)$$

называют функцией максимального правдоподобия.

Совместная функция вероятности $L(x_1, x_2, \dots, x_n, \theta)$ может быть представлена в виде (14.4) в силу независимости результатов каждого измерения. Две функции правдоподобия являются равными, если одна есть произведение второй на некоторую скалярную величину.

$L(x_1, x_2, \dots, x_n, \theta)$ рассматривают как функцию θ при фиксированных, полученных в измерениях x . Чем больше значение $L(x_1, x_2, \dots, x_n, \theta)$, тем

более вероятна или правдоподобна, выборка значений $\{x_1, x_2, \dots, x_n\}$ при заданном значении θ . Поэтому $L(x_1, x_2, \dots, x_n, \theta)$ и называют функцией правдоподобия.

Метод максимального правдоподобия или метод наибольшего правдоподобия (ММП, ML, MLE — англ. maximum likelihood estimation) в математической статистике — это метод оценивания неизвестного параметра путём максимизации функции правдоподобия. Основан на предположении о том, что вся информация о статистической выборке содержится в функции правдоподобия.

Метод максимального правдоподобия состоит в поиске максимума функции $L(x_1, x_2, \dots, x_n, \theta)$, при этом $\{x_1, x_2, \dots, x_n\}$ считаются постоянными, а параметр θ — переменным. Далее находят такое значение θ , при котором функция $L(x_1, x_2, \dots, x_n, \theta)$ становится наиболее правдоподобной, т.е. **принимает максимальное значение**.

Для эффективного использования ММП требуются достаточно большие объемы выборок, точное знание анализируемого закона распределения, достаточно устойчивые распределения и не очень большое число неизвестных параметров.

Доказано, что вторая производная от функции правдоподобия $L(x_1, x_2, \dots, x_n, \theta)$ меньше нуля и, таким образом, равенство нулю первой производной дает действительно максимальное значение $L(x_1, x_2, \dots, x_n, \theta)$. Максимум определяют по стандартной методике, однако вместо самой функции для удобства вычислений берут логарифм функции и ищут его максимум. Поскольку максимумы функции правдоподобия и логарифмической функции совпадают, то для удобства вычислений берут логарифм функции и ищут его максимум.

Введем обозначение

$$\ell(x_1, x_2, \dots, x_n, \theta) = \ln L(x_1, x_2, \dots, x_n, \theta) = \sum_{i=1}^n \ln p(x_i, \theta) . \quad (14.5)$$

Система уравнений для получения оценок параметров $\theta = \{\theta_1, \theta_2, \dots, \theta_k\}$

МОЖНО ЗАПИСАТЬ В ВИДЕ

$$\left\{ \begin{array}{l} \frac{\partial \ell(x_1, x_2, \dots, x_n, \boldsymbol{\theta})}{\partial \theta_1} = 0, \\ \frac{\partial \ell(x_1, x_2, \dots, x_n, \boldsymbol{\theta})}{\partial \theta_2} = 0, \\ \dots \\ \frac{\partial \ell(x_1, x_2, \dots, x_n, \boldsymbol{\theta})}{\partial \theta_k} = 0, \end{array} \right. \quad (14.6)$$

Решение системы уравнений (14.6) дает набор оптимальных значений параметров $\tilde{\boldsymbol{\theta}} = \{\tilde{\theta}_1, \tilde{\theta}_2, \dots, \tilde{\theta}_k\}$.

Наибольшие трудности возникают не при оценке параметров с помощью решения системы (14.6), а при выявлении погрешностей, с которыми эти оценки сделаны. Возможные (неявные и естественные) связи между параметрами $\boldsymbol{\theta} = \{\theta_1, \theta_2, \dots, \theta_k\}$ приводят к необходимости учитывать корреляции (ковариации) между оценками различных параметров $\tilde{\boldsymbol{\theta}} = \{\tilde{\theta}_1, \tilde{\theta}_2, \dots, \tilde{\theta}_k\}$.

Введем обозначения для сокращения записи соотношений:

$$\begin{aligned} \ell(x_1, x_2, \dots, x_n, \boldsymbol{\theta}) &= \ell(\boldsymbol{\theta}), \\ L(x_1, x_2, \dots, x_n, \tilde{\boldsymbol{\theta}}) &= L_{max}. \end{aligned} \quad (14.7)$$

Для оценки погрешностей параметров $\tilde{\boldsymbol{\theta}} = \{\tilde{\theta}_1, \tilde{\theta}_2, \dots, \tilde{\theta}_k\}$ разложим логарифмическую функцию правдоподобия в ряд Тейлора в окрестностях точек $\{\tilde{\theta}_1, \tilde{\theta}_2, \dots, \tilde{\theta}_k\}$:

$$\begin{aligned} \ell(\boldsymbol{\theta}) &= \ell(\tilde{\boldsymbol{\theta}}) + \sum_{i=1}^k \left(\frac{\partial \ell}{\partial \theta_i} \right)_{\tilde{\boldsymbol{\theta}}} (\theta_i - \tilde{\theta}_i) + \\ &+ \frac{1}{2} \sum_{p=1}^k \sum_{m=1}^k \left(\frac{\partial^2 \ell}{\partial \theta_p \partial \theta_m} \right)_{\tilde{\boldsymbol{\theta}}} (\theta_i - \tilde{\theta}_i) (\theta_m - \tilde{\theta}_m) + \dots \end{aligned} \quad (14.8)$$

Величины $(\theta_i - \tilde{\theta}_i)$ трактуются как дисперсии оценок максимального правдоподобия $\tilde{\theta}_i$.

Из определения оценок $\tilde{\boldsymbol{\theta}}$ следует, что второй член в разложении равен нулю для всех k . Пренебрегая членами разложения выше квадратичных, логарифмическую функцию правдоподобия можно записать в матричном виде:

$$\ell(\boldsymbol{\theta}) = \ln L_{max} + \frac{1}{2} (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}^T) \mathbf{A} (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}) , \quad (14.9)$$

$$\mathbf{A} = \begin{pmatrix} \frac{\partial^2 \ell}{\partial \theta_1^2} & \frac{\partial^2 \ell}{\partial \theta_1 \theta_2} & \cdots & \frac{\partial^2 \ell}{\partial \theta_1 \theta_k} \\ \frac{\partial^2 \ell}{\partial \theta_1 \theta_2} & \frac{\partial^2 \ell}{\partial \theta_2^2} & \cdots & \frac{\partial^2 \ell}{\partial \theta_2 \theta_k} \\ \cdots & \cdots & \cdots & \cdots \\ \frac{\partial^2 \ell}{\partial \theta_1 \theta_k} & \frac{\partial^2 \ell}{\partial \theta_2 \theta_k} & \cdots & \frac{\partial^2 \ell}{\partial \theta_k^2} \end{pmatrix}_{\boldsymbol{\theta} = \tilde{\boldsymbol{\theta}}} . \quad (14.10)$$

В асимптотическом пределе $n \rightarrow \infty$, элементы матрицы практически не зависят от конкретной выборки $\{x_1, x_2, \dots, x_n\}$, и их можно заменить математическими ожиданиями:

$$\mathbf{B} = \mathbf{E} \{ \mathbf{A} \} = \begin{pmatrix} \mathbf{E} \left\{ \frac{\partial^2 \ell}{\partial \theta_1^2} \right\} & \mathbf{E} \left\{ \frac{\partial^2 \ell}{\partial \theta_1 \theta_2} \right\} & \cdots & \mathbf{E} \left\{ \frac{\partial^2 \ell}{\partial \theta_1 \theta_k} \right\} \\ \mathbf{E} \left\{ \frac{\partial^2 \ell}{\partial \theta_1 \theta_2} \right\} & \mathbf{E} \left\{ \frac{\partial^2 \ell}{\partial \theta_2^2} \right\} & \cdots & \mathbf{E} \left\{ \frac{\partial^2 \ell}{\partial \theta_2 \theta_k} \right\} \\ \cdots & \cdots & \cdots & \cdots \\ \mathbf{E} \left\{ \frac{\partial^2 \ell}{\partial \theta_1 \theta_k} \right\} & \mathbf{E} \left\{ \frac{\partial^2 \ell}{\partial \theta_2 \theta_k} \right\} & \cdots & \mathbf{E} \left\{ \frac{\partial^2 \ell}{\partial \theta_k^2} \right\} \end{pmatrix}_{\boldsymbol{\theta} = \tilde{\boldsymbol{\theta}}} . \quad (14.11)$$

После потенцирования (14.9) следует, что функция правдоподобия имеет вид

$$L(\boldsymbol{\theta}) = L_{max} \exp \left\{ \frac{1}{2} (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}^T) \mathbf{B} (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}) \right\} . \quad (14.12)$$

Сопоставление с известными теоретическими распределениями приводит к выводу, что (14.12) представляет собой k -мерное нормальное распределение со средними оценками $\{\tilde{\theta}_1, \tilde{\theta}_2, \dots, \tilde{\theta}_k\}$ и ковариационной матрицей

$$\mathbf{V} = -\mathbf{B}^{-1} . \quad (14.13)$$

Диагональные элементы ковариационной матрицы (14.13) - являются дисперсиями $c_{11} = \sigma_{\tilde{\theta}_1}, \dots, c_{kk} = \sigma_{\tilde{\theta}_k}$ оценок максимального правдоподобия $\{\tilde{\theta}_1, \tilde{\theta}_2, \dots, \tilde{\theta}_k\}$, а недиагональные элементы представляют собой ковариации между всевозможными парами оценок $c_{ij} = \text{cov}(\tilde{\theta}_i, \tilde{\theta}_j)$.

Коэффициент корреляции между оценками определяется формулой

$$\rho(\tilde{\theta}_i, \tilde{\theta}_j) = \frac{\text{cov}(\tilde{\theta}_i, \tilde{\theta}_j)}{\sigma_{\tilde{\theta}_i} \sigma_{\tilde{\theta}_j}}, \quad i \neq j. \quad (14.14)$$

Поскольку $L(\boldsymbol{\theta})$ в окрестности $\tilde{\boldsymbol{\theta}}$ представляет собой k -мерное нормальное распределение, то квадратичная форма

$$Q(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}) = (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})^T \mathbf{V}^{-1} (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}) \quad (14.15)$$

распределена по закону χ^2 с r степенями свободы. Поэтому можно определить вероятностное утверждение:

$$\text{Prob}[Q(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}) \leq Q_\alpha] = \alpha, \quad (14.16)$$

где квантиль распределения χ^2 с r степенями свободы Q_α уровня α задается соотношением

$$\int_0^{Q_\alpha} \frac{2^{-r/2}}{\Gamma(r/2)} \exp[-x/2] x^{r/2-1} dx = 1 - \alpha. \quad (14.17)$$

Гиперэллипсоид, определяемый уравнением

$$Q(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}) = Q_\alpha \quad (14.18)$$

будет задавать доверительную область в пространстве переменных $\tilde{\boldsymbol{\theta}}$ с вероятностным содержанием (доверительной вероятностью) $P = 1 - \alpha$.

Для приближенного построения гиперэллипсоида (14.18), который определяет возможные значения оптимальных параметров с доверительной вероятностью P , можно использовать уравнение вида (смотри (14.9)):

$$\ln L(\boldsymbol{\theta}) \approx \ln L_{max} - \Delta(\ln L(\boldsymbol{\theta})) \quad (14.19)$$

Таким образом, как следует из (14.17) и (14.18), значение $\Delta(\ln L(\boldsymbol{\theta}))$ определяется задаваемым уровнем достоверности (С.Л.) и числом параметров, входящих в набор $\boldsymbol{\theta}$.

С помощью (14.19) для доверительной вероятности $P = 95\%$ можно найти для числа параметров 1, 2 и 3, числовое значение $\Delta(\ln L(\boldsymbol{\theta})) \approx 3,84, 5,99$

и 7,82 соответственно. Аналогичные значения для стандартного отклонения σ $P = 68,27\%$ приводят к значениям $\Delta(\ln L(\theta)) \approx 1,00, 2,30$ и $3,53$.

В этом случае, можно построить контурные графики функции

$$\ln L(\theta) - \ln L_{max} = \Delta(\ln L(\theta)) \quad (14.20)$$

для различных пар параметров и найти среднеквадратичные отклонения $\sigma_{\tilde{\theta}_1}, \dots, \sigma_{\tilde{\theta}_k}$ оценок максимального правдоподобия $\{\tilde{\theta}_1, \tilde{\theta}_2, \dots, \tilde{\theta}_k\}$.

14.2 Метод наименьших квадратов

Метод наименьших квадратов также используется для получения “наилучших” (оптимальных) значений модельных параметров при описании экспериментальных значений.

Метод наименьших квадратов совпадает с методом максимального правдоподобия в следующем частном случае. Рассмотрим набор из n независимых измерений y_i в известных точках x_i .

Из y_i из $\mathbf{y} = \{y_1 \pm \sigma_1, \dots, y_n \pm \sigma_n\}$ подчиняются подчиняются нормальному распределению с математическим ожиданием $\phi(x_i, \theta)$ и известной дисперсией σ_i^2 . Таким образом предполагается, что описание значений y_i с помощью модельной регрессионной кривой $\phi(x_i, \theta)$ должно быть оптимальным. Цель состоит в том, чтобы построить оценки для неизвестных параметров θ .

Для получения оптимальных значений $\tilde{\theta}$ набора модельных параметров θ используют функцию

$$\chi^2(\theta) = -2 \ln L(\theta) + \text{const} = \sum_{i=1}^n \left[\frac{y_i - \phi(x_i, \theta)}{\sigma_i} \right]^2. \quad (14.21)$$

Соотношение (14.21) можно получить из определения функции максимального правдоподобия (14.1), используя плотность нормального распределения

$$p(y_i) = \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp \left[-\frac{(y_i - \phi(x_i, \theta))^2}{2\sigma_i^2} \right]. \quad (14.22)$$

Исходя из требования минимального значения функции $\chi^2(\theta)$ при $\theta = \tilde{\theta}$, т. е.

$$\chi^2(\tilde{\theta}) = \chi_{min}^2 \quad (14.23)$$

получаем систему уравнений

$$\begin{cases} \frac{\partial \chi^2(\boldsymbol{\theta})}{\partial \theta_1} = 0, \\ \frac{\partial \chi^2(\boldsymbol{\theta})}{\partial \theta_2} = 0, \\ \dots \\ \frac{\partial \chi^2(\boldsymbol{\theta})}{\partial \theta_k} = 0, \end{cases} \quad (14.24)$$

для определения “оптимальных” значений $\tilde{\boldsymbol{\theta}}$. Процедура аналогична линейному регрессионному анализу, с той лишь разницей, что поиск ‘оптимальных’ значений $\tilde{\boldsymbol{\theta}}$, как правило, находится численно, а не аналитически.

Значение χ_{min}^2 является мерой уровня согласия между экспериментальными данными и подогнанной кривой (fit):

$$\chi_{min}^2 = \sum_{i=1}^n \left[\frac{y_i - \phi(x_i, \tilde{\boldsymbol{\theta}})}{\sigma_i} \right]^2. \quad (14.25)$$

Поэтому χ_{min}^2 можно использовать как статистику соответствия предполагаемую функциональную форму $\phi(x_i, \tilde{\boldsymbol{\theta}})$. Известно, что (14.25) имеет распределение χ^2 с $n_d = n - k$ степенями свободы.

В этом случае, в соответствии с разделом проверки гипотез, если $\chi_{min}^2/n_d \leq 1$, то совпадение будет “хорошим”. Или можно найти вероятность соответствия модели $\phi(x_i, \tilde{\boldsymbol{\theta}})$ экспериментальным данным P (p -value) из соотношения:

$$P = \frac{1}{2^{n_d} \int_{\chi_{min}^2}^{\infty} \Gamma\{n_d/2\}} t^{n_d/2-1} \exp(-t/2). \quad (14.26)$$

Следующим шагом решения данной задачи состоит в нахождении доверительных интервалов (ошибок $\sigma_{\tilde{\theta}_i}$) для оптимальных значений $\tilde{\boldsymbol{\theta}} = \{\tilde{\theta}_1, \dots, \tilde{\theta}_k\}$. Эта задача решается так же как и для метода максимального правдоподобия, если предположить, что оценки $\tilde{\theta}_i, i = 1, \dots, k$ подчиняются k -мерному нормальному распределению (14.12).

Отличием, состоит в том, что приближенного построения гиперэллипсоида (14.18), который определяет возможные значения оптимальных парамет-

ров с вероятностным содержанием С.Л., используется уравнение вида:

$$\chi^2(\boldsymbol{\theta}) = \chi_{min}^2 + \Delta\chi_{crit}^2. \quad (14.27)$$

Значение $\Delta\chi_{crit}^2$ определяется задаваемым уровнем достоверности p и числом параметров, входящих в набор $\boldsymbol{\theta}$ и совпадает с $\Delta(\ln L(\boldsymbol{\theta}))$.

Замечание.

Если y_i не являются независимыми, и имеют ковариационную матрицу $\mathbf{V}_{ij} = \text{cov}(y_i, y_j)$, то наилучшие оценки параметров $\tilde{\boldsymbol{\theta}} = \{\tilde{\theta}_1, \dots, \tilde{\theta}_k\}$ путем поиска минимума функционала

$$\chi^2(\boldsymbol{\theta}) = \sum_{i=1}^n \sum_{j=1}^n (y_i - \phi(x_i, \boldsymbol{\theta})) (\mathbf{V})_{ij}^{-1} (y_j - \phi(x_j, \boldsymbol{\theta})). \quad (14.28)$$

Репозиторий ГГУ им. Ф.Скорины

15 Сравнение моделей

При поиске зависимости между измеряемыми переменными по выборочным данным наиболее важной задачей является поиск модели регрессии. Построение эмпирического уравнения регрессии - начальный этап анализа. При этом, на первом этапе не всегда удастся получить “оптимальную” модель регрессии. Тогда для описания зависимости между случайными величинами может использоваться несколько моделей (уравнений регрессии). Поэтому возникает задача об оценке “качества” полученных моделей.

Качество модели оценивается по следующим направлениям:

1. Содержательная оценка качества модели.
2. Статистическая оценка качества модели.

Содержательный анализ подразумевает анализ физического (экономического и т.д.) смысла модели. Результат данной работы должен быть ответы на вопросы вида:

- ✓ *Являются ли те факторы, которые используются в модели, значимыми для описания данного физического явления или процесса?*
- ✓ *Какой физический смысл параметров регрессионного уравнения (модели)?*

Статистическая оценка качества модели включает следующие этапы:

- ✓ Проверка статистической значимости параметров уравнения регрессии при помощи различных критериев.
- ✓ Проверку общего качества уравнения регрессии.
- ✓ Проверку свойств данных, выполнение которых предполагалось при оценивании уравнения.

15.1 Проверка статистической значимости параметров уравнения регрессии

После того, как найдено уравнение регрессии (например, линейной регрессии (13.1)), проводится оценка значимости как уравнения в целом, так и отдель-

ных его параметров.

Оценка значимости уравнения регрессии в целом дается с помощью F -критерия Фишера. При этом выдвигается нулевая гипотеза о том, что коэффициент регрессии B равен нулю, т. е. $H_0 : B = 0$. Следовательно, фактор случайной величины x не оказывает влияния на результат y .

Перед расчетом критерия проводятся анализ дисперсии. Общая сумма квадратов отклонений (СКО) y от среднего значения \bar{y} раскладывается на две части - “объясненную” и “необъясненную”:

$$\underbrace{\sum_{i=1}^n (y_i - \bar{y}_i)^2}_{\text{Общая СКО}} = \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y}_i)^2}_{\text{Факторная СКО}} + \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i)^2}_{\text{Остаточная СКО}}. \quad (15.1)$$

Проиллюстрируем соотношение (15.1) на примере двух крайних вариантов. Если общая СКО в точности равна остаточной, то случайный фактор x не оказывает влияния на результат, вся дисперсия y обусловлена воздействием прочих факторов. При этом прямая для линейной регрессии параллельна оси Ox .

Когда общая СКО равна факторной, то никакие другие (прочие) факторы не влияют на результат. Величина y связана с x с вероятностью 100%, и остаточная СКО равна нулю.

Однако на практике в правой части (15.1) присутствуют оба слагаемых. Пригодность линии регрессии для прогноза зависит от того, какая часть общей вариации y приходится на объясненную вариацию. Если объясненная СКО будет больше остаточной СКО, то уравнение регрессии статистически значимо и фактор x оказывает существенное воздействие на результат y . Это равносильно тому, что коэффициент корреляции будет приближаться к единице.

Для использования критерия Фишера для оценки значимости регрессии необходимо ввести понятие число степеней свободы для различных величин.

Определение 15.1

Число степеней свободы (*df-degrees of freedom*) - это минимально необходимое число значений зависимой переменной, которых достаточно для получения искомой характеристики выборки и которые могут свободно изменяться с учетом того, что для этой выборки известны все другие величины, используемые для расчета искомой характеристики. Или другими словами-это **число независимо варьируемых значений признака**.

Уравнение для определения F -статистики в случае многомерной регрессии имеет вид:

$$F_{\Phi} = \left(\frac{\sum_{i=1}^n (\hat{y}_i - \bar{y}_i)^2}{m} \right) \left\{ 1 / \frac{\left(\sum_{i=1}^n (y_i - \hat{y})^2 \right)}{(n - m - 1)} \right\}, \quad (15.2)$$

где n - объем выборки; m - число независимых переменных в факторной части СКО.

Поясним на примере линейной регрессии. Факторную СКО в этом случае можно выразить так:

$$\begin{aligned} \sum_{i=1}^n (\hat{y}_i - \bar{y}_i)^2 &= \sum_{i=1}^n ([A + Bx_i] - [A + B\bar{x}])^2 = \\ &= \sum_{i=1}^n (Bx_i - B\bar{x})^2 = B^2 \sum_{i=1}^n (x_i - \bar{x})^2. \end{aligned} \quad (15.3)$$

СКО (15.3) зависит только от одного параметра B . Следовательно, факторная СКО в случае линейной регрессии имеет одну степень свободы, т. е. $m = 1$.

Можно и рассчитать число m и другим способом. Для того вспомним выражение (13.7) для регрессионной прямой

$$\hat{y} = A + Bx = \bar{y} + B(x - \bar{x}). \quad (15.4)$$

Как видно из (15.4) прямая регрессии определяется только одним параметром.

Разделив каждую СКО на свое число степеней свободы, получим среднее квадратичные отклонения (или дисперсии на одну степень свободы):

$$\begin{aligned}
 D_{\text{общ.}} &= \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y}_i)^2, \\
 D_{\text{факт.}} &= \frac{1}{m} \sum_{i=1}^n (\hat{y}_i - \bar{y}_i)^2, \\
 D_{\text{ост.}} &= \frac{1}{n-m-1} \sum_{i=1}^n (y_i - \hat{y}_i)^2.
 \end{aligned} \tag{15.5}$$

С учетом (15.5) F -критерий (15.2) можно записать в виде:

$$F_{\Phi} = \frac{D_{\text{факт.}}}{D_{\text{ост.}}}. \tag{15.6}$$

В теории вероятности доказано, что для выборки из генеральной совокупности у которой отсутствует связь между зависимой y и независимой x (или x_1, x_2, \dots, x_k) переменной имеет распределение Фишера:

$$\begin{aligned}
 p_{F_R}(x, f_1, f_2) &= \frac{f_1^{f_1/2} f_2^{f_2/2} x^{f_1/2-1} (f_1 x + f_2)^{-\frac{1}{2}(f_1+f_2)}}{B\left(\frac{f_1}{2}, \frac{f_2}{2}\right)}, \quad x \geq 0, \\
 F_R(x, f_1, f_2) &= I_{\frac{f_1 x}{f_1 x + f_2}}\left(\frac{f_1}{2}, \frac{f_2}{2}\right).
 \end{aligned} \tag{15.7}$$

где $f_{1,2}$ - параметры распределения; функция $I_n(a,b)$ - регуляризованная неполная бэ́та-функция.

Благодаря этому для осуществления статистической проверки значимости уравнения регрессии формулируется нулевая гипотеза об отсутствии связи между переменными (все коэффициенты при переменных равны нулю) и выбирается уровень значимости α .

Уровень значимости α - это вероятность совершить ошибку первого рода, т. е. отвергнуть в результате проверки верную нулевую гипотезу. В рассматриваемом случае совершить ошибку первого рода означает признать по выборке наличие связи между переменными в генеральной совокупности, **когда на самом деле ее там нет.**

Вычисленное значение F_{Φ} признается достоверным (отличным от единицы), если оно больше табличного квантиля распределения Фишера F_{α}

$$\int_{F_{\alpha}}^{\infty} p_{F_R}(x, m, n - m - 1) dx = \alpha, \quad (15.8)$$

т.е.

$$F_{\Phi} \geq F_{\alpha}. \quad (15.9)$$

В этом случае нулевая гипотеза об отсутствии связи между переменными отклоняется и делается вывод о существенности превышения $D_{\text{факт.}}$ над $D_{\text{ост.}}$. Или другими словами делается вывод о существенной статистической зависимости между зависимой y и независимыми величинами $x = x_1, \dots, x_k$.

Если $F_{\Phi} \leq F_{\alpha}$, то вероятность нулевой гипотезы выше заданного уровня α (например, $0,05 \div 0,10$), и эта гипотеза не может быть отклонена без серьезного риска сделать неправильный вывод о наличии связи между y и x . Уравнение регрессии считается статистически незначимым, гипотеза H_0 не отклоняется.

С помощью компьютерных вычислений, практически всегда можно найти α путем решения уравнения:

$$\int_{F_{\Phi}}^{\infty} p_{F_R}(x, m, n - m - 1) dx = 1 - F_R(F_{\Phi}, m, n - m - 1) = \alpha. \quad (15.10)$$

В если α относительно велико ($\alpha > 0,1$), то о наличии корреляции говорит сложно (или с большими оговорками).

Как и в случае парной регрессии, статистическая значимость коэффициентов множественной линейной регрессии с m объясняющими переменными проверяется на основе t -критерия Стьюдента. Величина стандартной ошибки совместно с t -распределением Стьюдента при $n - m$ степенях свободы применяется для проверки существенности коэффициента регрессии и для расчета его доверительных интервалов. Для этого для каждого параметра регрессии a_i , полученного в результате вычислений рассчитывается величина

$$t_{\text{exp}} = \frac{a_i}{\sigma_{a_i}}, \quad (15.11)$$

т.е. величина параметров регрессии сравнивается с его стандартной ошибкой.

Значение (15.11) сравнивается с табличным значением при определенном уровне значимости α и числе степеней свободы. По сути проверяется нулевая

гипотеза в виде $H_0: a_i = 0$. Если $t_{\text{exp}} > t_{n-m, \alpha/2}$, то гипотеза $H_0: a_i = 0$ должна быть отклонена, а статистическая связь y с x считается установленной. В случае $t_{\text{exp}} > t_{n-m, \alpha/2}$ нулевая гипотеза не может быть отклонена, и влияние x на y признается несущественным.

15.2 Проверка общего качества уравнения регрессии

Оценить общее качество уравнения регрессии означает установить соответствует ли математическая модель, выражающая зависимость между переменными экспериментальным данным и достаточно ли включенных в модель переменных, объясняющих поведение предсказанной величины (y). Оценить общее качества модели = оценить надежность модели = оценить достоверность уравнения регрессии.

15.3 Коэффициент детерминации

Для оценки качества различных моделей используют коэффициент детерминации.

Определение 15.2

Коэффициент детерминации (R^2 — R-квадрат) — это доля дисперсии зависимой переменной, объясняемая рассматриваемой моделью зависимости, то есть объясняющими переменными.

Более точно — это единица минус доля необъяснённой дисперсии (дисперсии случайной ошибки модели, или условной по факторам дисперсии зависимой переменной) в дисперсии зависимой переменной. Его рассматривают как универсальную меру зависимости одной случайной величины от множества других. В частном случае линейной зависимости R^2 является квадратом так называемого множественного коэффициента корреляции между зависимой переменной и объясняющими переменными. В частности, для модели парной линейной регрессии коэффициент детерминации равен квадрату обычного коэффициента корреляции между y и x .

Истинный коэффициент детерминации модели зависимости случайной величины y от факторов x определяется следующим образом:

$$R^2 = 1 - \frac{D\{y|x\}}{D\{y\}} = 1 - \frac{\sigma^2}{\sigma_y^2}, \quad (15.12)$$

где $D\{y|x\} = \sigma^2$ — условная (по факторам x) дисперсия зависимой переменной (дисперсия случайной ошибки модели).

В данном определении используются истинные параметры, характеризующие распределение случайных величин. Если использовать выборочную оценку значений соответствующих дисперсий, то получим формулу для выборочного коэффициента детерминации (который обычно и подразумевается под коэффициентом детерминации):

Коэффициент детерминации

$$R^2 = 1 - \frac{\hat{\sigma}^2}{\hat{\sigma}_y^2} = 1 - \frac{SS_{res}/n}{SS_{tot}/n} = 1 - \frac{SS_{res}}{SS_{tot}}, \quad (15.13)$$

где

$$SS_{res} = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (15.14)$$

сумма квадратов остатков регрессии, y_i, \hat{y}_i — фактические (“экспериментальные”) и расчётные значения объясняемой переменной, а величина

$$SS_{tot} = \sum_{i=1}^n (y_i - \bar{y})^2 = n\hat{\sigma}_y^2 \quad (15.15)$$

так называемая **общая сумма квадратов**, а

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i.$$

В случае линейной регрессии с константой $SS_{tot} = SS_{reg} + SS_{res}$, где

$$SS_{reg} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \quad (15.16)$$

объяснённая сумма квадратов, поэтому получаем более простое определение в этом случае — коэффициент детерминации — это доля объяснённой суммы квадратов в общей:

$$R^2 = \frac{SS_{reg}}{SS_{tot}} \quad (15.17)$$

Необходимо подчеркнуть, что формула (15.17) справедлива только для модели с константой, в общем случае необходимо использовать предыдущую формулу.

Интерпретация

1. Коэффициент детерминации для модели с константой принимает значения от 0 до 1. Чем ближе значение коэффициента к 1, тем сильнее зависимость. При оценке регрессионных моделей это интерпретируется как соответствие модели данным. Для приемлемых моделей предполагается, что коэффициент детерминации должен быть хотя бы не меньше 50% (в этом случае коэффициент множественной корреляции превышает по модулю 70%). Модели с коэффициентом детерминации выше 80% можно признать достаточно хорошими (коэффициент корреляции превышает 90%). Значение коэффициента детерминации 1 означает функциональную зависимость между переменными.

Интерпретация

2. При отсутствии статистической связи между объясняемой переменной и факторами, статистика nR^2 для линейной регрессии имеет

асимптотическое распределение $\chi^2(k - 1)$, где $k - 1$ — количество факторов модели. В случае линейной регрессии с нормально распределёнными случайными ошибками статистика $F = \frac{R^2/(k-1)}{(1-R^2)/(n-k)}$ имеет точное (для выборок любого объёма) распределение Фишера $F(k - 1, n - k)$ (так называемый F-тест). Информация о распределении этих величин позволяет проверить статистическую значимость регрессионной модели исходя из значения коэффициента детерминации. Фактически в этих тестах проверяется гипотеза о равенстве истинного коэффициента детерминации нулю.

3. В общем случае коэффициент детерминации может быть и отрицательным, это говорит о крайней неадекватности модели: простое среднее приближает лучше.

15.4 Недостаток R^2 и альтернативные показатели

Основная проблема применения (выборочного) R^2 заключается в том, что его значение увеличивается (не уменьшается) от добавления в модель новых переменных, даже если эти переменные никакого отношения к объясняемой переменной не имеют! Поэтому сравнение моделей с разным количеством факторов с помощью коэффициента детерминации, вообще говоря, некорректно. Для этих целей можно использовать альтернативные показатели.

Скорректированный (adjusted) R^2

Для того, чтобы была возможность сравнивать модели с разным числом факторов так, чтобы число факторов не влияло на статистику R^2 обычно используется скорректированный коэффициент детерминации, в котором используются несмещённые оценки дисперсий:

$$R_{adj}^2 = 1 - \frac{s^2}{s_y^2} = 1 - \frac{SS_{res}/(n - k)}{SS_{tot}/(n - 1)} = 1 - (1 - R^2) \frac{(n - 1)}{(n - k)} \leq R^2 \quad (15.18)$$

который даёт “штраф” за дополнительно включённые факторы, где n — количество наблюдений, а k — количество параметров.

Данный показатель всегда меньше единицы, но теоретически может быть и меньше нуля (только при очень маленьком значении обычного коэффициента детерминации и большом количестве факторов). Поэтому теряется интерпретация показателя как “доли”. Тем не менее, применение показателя в сравнении вполне обоснованно.

Для моделей с одинаковой зависимой переменной и одинаковым объёмом выборки сравнение моделей с помощью скорректированного коэффициента детерминации эквивалентно их сравнению с помощью остаточной дисперсии $s^2 = SS_{res}/(n - k)$ или стандартной ошибки модели s . Разница только в том, что последние критерии чем меньше, тем лучше.

Для оценки качества различных моделей при описании данных также используют информационные критерии.

Определение 15.3

Информационный критерий — применяемая в статистике мера относительного качества статистических моделей, учитывающая степень “подгонки” модели под данные с корректировкой (“штрафом”) на используемое количество оцениваемых параметров. То есть критерии основаны на некоем компромиссе между точностью и сложностью модели. Критерии различаются тем, как они обеспечивают этот баланс.

Информационный характер критериев связан с концепцией информационной энтропии и расстоянием Кульбака-Лейблера, на основе которой был разработан исторически первый критерий — критерий Акаике (AIC), предложенный в 1974 году Хиротсугу Акаике (H. Akaike, “A new look at the statistical model identification” IEEE Transactions on Automatic Control, vol. 19, № . 6, pp. 716-723, 1974.)

Информационные критерии используются исключительно для сравнения моделей между собой, без содержательной интерпретации значений этих критериев. Они не позволяют тестировать модели в смысле проверки статистических гипотез. Обычно чем меньше значения критериев, тем выше относительное качество модели.

Информационный критерий Акаике (AIC)

Предложен Хиротугу Акаике в 1971 году, описан и исследован им же в 1973, 1974, 1983 годах. Первоначально аббревиатура AIC, предложенная автором, расшифровывалась как an information criterion (“некий информационный критерий”), однако последующие авторы называли его Akaike information criterion. Исходная расчетная формула критерия имеет вид:

$$AIC = 2k - 2L \quad (15.19)$$

где L - значение логарифмической функции правдоподобия построенной модели, k -количество использованных (оцененных) параметров.

Многие современные авторы, а также во многих программных продуктах применяется несколько иная формула, предполагающая деление на объем выборки n , по которой строилась модель:

$$AIC = \frac{2k - 2L}{n} . \quad (15.20)$$

Данный подход позволяет сравнивать модели, оцененные по выборках разного объема.

Чем меньше значение критерия, тем лучше модель.

Многие другие критерии являются модификациями AIC.

Байесовский информационный критерий (BIC) или критерий Шварца (SC)

Байесовский информационный критерий (Bayesian information criterion — BIC) предложен Шварцем в 1978 году, поэтому часто он называется также критерием Шварца (Schwarz criterion — SC). Он разработан исходя из байесовского подхода и является наиболее часто используемой модификацией AIC:

$$BIC = SC = k \ln n - 2L . \quad (15.21)$$

Как видно из формулы, данный критерий налагает большой штраф на увеличение количества параметров по сравнению с AIC, так как $\ln n$ больше 2 уже при количестве 8 наблюдений.

Прочие информационные критерии

Состоятельный критерий Акаике (Consistent AIC – CAIC) предложенный в 1987 году Боздоганом:

$$\text{CAIC} = (1 + \ln n)k - 2L . \quad (15.22)$$

Данный критерий асимптотически эквивалентен BIC. Тот же автор в 1994 году предложил модификации, увеличивающие коэффициент при количестве параметров (вместо 2-3 или 4 для AIC₃ и AIC₄).

Скорректированный критерий Акаике (Corrected AIC – AIC_c), который рекомендуется применять на малых выборках (предложен в 1978 году Sugiura):

$$\text{AIC}_c = \text{AIC} + \frac{2k(k+1)}{n-k-1} . \quad (15.23)$$

Данный критерий, наряду с AIC и BIC выдается в результатах оценки моделей в Wolfram Mathematica при использовании оператора NonlinearModelFit.

Критерий Ханна-Куинна (Hannan-Quinn, HQ) предложен авторами в 1979 году

$$\text{HQ} = 2k \ln \ln n/n - 2L/n . \quad (15.24)$$

Имеются также модификации AIC, использующие более сложные штрафные функции, зависящие от различных характеристик.

Замечание

Высокие значения коэффициента детерминации, вообще говоря, не свидетельствуют о наличии причинно-следственной зависимости между переменными (также как и в случае обычного коэффициента корреляции). Например, если объясняемая переменная и факторы, на самом деле не связанные с объясняемой переменной, имеют возрастающую динамику, то коэффициент детерминации будет достаточно высок. **Поэтому логическая и смысловая адекватность модели имеют первостепенную важность.** Кроме того, необходимо использовать критерии для всестороннего анализа качества модели.

Репозиторий ГГУ им. Ф. Скорины

16 Рекомендуемая литература

Рекомендуемая литература

1. Лавренчик, В.Н. Постановка физического эксперимента и статистическая обработка его результатов/ В.Н. Лавренчик.. – М.: Энергоатомиздат, 1986. - 272 с.
2. Бендат Дж., Пирсол А. Прикладной анализ случайных данных/ Дж.Бендат, А.Пирсол. – М.: Мир,1989. -504 с.
3. Новицкий, П.В. Оценка погрешностей результатов измерений/ П.В.Новицкий, И.А..Зограф - Л.: Энергоатомиздат, 1991. - 304 с
4. Тейлор Дж. Введение в теорию ошибок/ Дж.Тейлор. - М.: Мир, 1985. - 45 с.
5. Гмурман, В. Е. Теория вероятностей и математическая статистика: Учеб. пособие для вузов/В. Е. Гмурман. - 9-е изд., стер. - М.: Высш. шк., 2003. - 479 с.

Дополнительная

Рекомендуемая литература

6. Дьяконов, В.П. Mathematica 4: учебный курс/ В.П. Дьяконов. - СПб.: Питер, 2001 . - 654 с.
7. Воробьев Е.М. Введение в систему Mathematica./ Е.М.Воробьев. - М.: Финансы и статистика, 1998. - 345 с.

8. Львовский, Б.Н. Статистические методы построения эмпирических формул: Учеб. пособие для втузов/ Б.Н.Львовский. – М.: Высш. шк., 1988 - 239 с.
9. Кобзарь А. И. Прикладная математическая статистика. Для инженеров и научных работников/ А. И. Кобзарь. – М.: Физматлит, 2006. – 816 с.
10. Боровков Л. Л. Математическая статистика/ Л. Л.Боровков-Учебник.- М.: Наука. Главная редакция физико-математической литературы, 1984. – 472 с.

11. Кассандрова, О. Н. Обработка результатов наблюдений/ О. Н. Кассандрова, В. В. Лебедев - М.:Наука, Главная редакция физ.-мат. литературы, 1970 г. – 104 с.

12. Ивченко, Г. И. Математическая статистика/ Г.И.Ивченко, Ю. И. Медведев Учеб. пособие для вузов. - М.: Высш. шк., 1984. – 248 с.

Репозиторий ГГУ им. Ф. Скорины

Учреждение образования
“Гомельский государственный университет
имени Франциска Скорины”
Кафедра теоретической физики

Статистические методы обработки данных

Специальность 1-31 04 08 Компьютерная физика

Лекции: 30 часов
Практические занятия: нет
Лабораторные занятия: 34 часа

Подготовил:
Андреев
Виктор Васильевич
доктор физ.-мат.наук, доцент

Гомель, 2018

Содержание I

- 1 Основные понятия теории вероятности
 - Частота попадания случайной величины
 - Интегральный закон распределения вероятности.
 - Дифференциальный закон распределения вероятности
 - Квантили распределений
 - Примеры законов распределений
- 2 Моменты распределения.
 - Математическое ожидание. Дисперсия.
 - Коэффициент асимметрии. Эксцесс.
- 3 Оценки параметров распределений
 - Точечные оценки
 - Интервальные оценки
 - Доверительный интервал для случайной величины
- 4 Взвешенное среднее значение
- 5 Ошибка косвенного измеряемой величины
- 6 Графическое представление эмпирических данных
- 7 Оптимальное число интервалов для получения гистограммы
- 8 Промахи и методы их исключения

Содержание II

- 9 Промахи и методы их исключения
 - Другие критерии для исключения промахов
- 10 Критерии согласия
- 11 Алгоритм предварительной обработки экспериментальных данных
- 12 Некоторые сведения о двумерных случайных величинах
 - Алгебраические и центральные моменты
 - Примеры многомерных функций распределения вероятностей
- 13 Корреляционный и регрессионный анализ
- 14 Линейный регрессионный анализ
 - Метод наименьших квадратов
 - Свойства метода наименьших квадратов
 - Точность оценок A и B .
 - Редукция некоторых задач к линейному анализу
- 15 Элементы нелинейного регрессионного анализа
 - Метод максимального правдоподобия
 - Метод наименьших квадратов
- 16 Сравнение моделей
 - Проверка статистической значимости параметров уравнения регрессии

Содержание III

- Проверка общего качества уравнения регрессии
- Коэффициент детерминации
- Недостаток R^2 и альтернативные показатели

17 Рекомендуемая литература

1. Основные понятия теории вероятности

Для исследования свойств объекта проводят измерения, позволяющие количественно дать характеристики свойств этого объекта. На практике производится ограниченное число измерений n при одинаковых условиях. В результате получаем множество событий (значений) исследуемой величины $X: \{x_1, x_2, \dots, x_n\}$, которое представляет собой выборку объема n . Таким образом, возможные исходы некоторого эксперимента (измерения) представляют собой множество точек, которые можно сочетать разными способами. **Такие сочетания называют событиями.**

1. 1. Частота попадания случайной величины

В том случае, (т.е. когда одна и та же характеристика объекта принимает различные значения) говорят о том, что величина X является **случайной величиной**. Случайная величина - это действительное число (или набор действительных чисел), которое заключено между $-\infty$ и $+\infty$, которое сопоставляется каждой возможной точке из числа значений этой характеристики. При соответствующих условиях для каждое событие можно характеризовать **частотой появления** ν_n .

Определение 1.1

Частотой появления ν_n события A называется отношение числа m появления данного события к общему числу проведенных одинаковых испытаний, в каждом из которых могло появиться или не появиться данное событие:

$$\nu_n = \frac{m}{n}. \quad (1.1)$$

Определение 1.2

Если число испытаний n велико, то, как правило, частоты появления данного события A в различных сериях измерений отличаются мало друг от друга. Это утверждение записывают следующим образом:

$$\lim_{n \rightarrow \infty} \nu_n = p. \quad (1.2)$$

Число p называются вероятностью (англ.-probability) случайного события A .

Отметим, что существуют такие события у которых частота появления может сильно отличаться от вероятности, даже при большом числе испытаний.

Как видим из вышеприведенного примера, для описания случайного поведения величины необходима совокупность, содержащая неограниченное число значений измеряемой величины ($n = \infty$). Такая выборка называется генеральной совокупностью. Генеральная совокупность часто используется как важное абстрактное понятие, необходимое в теоретических расчетах, связанных с исследованием поведения физической величины как случайной величины.

Различают два основных типа случайных величин: дискретные случайные величины и непрерывные случайные величины. Если величина X имеет **конечное** число (счетное множество) из последовательности возможных значений $\{x_1, x_2, \dots, x_k, \dots\}$ то такая величина называется *дискретной случайной величиной*. Если случайная величина может принимать любое значение из интервала возможных значений, то такая величина называется *непрерывной случайной величиной*.

Для характеристики частоты появления различных значений случайной величины X теория вероятностей предлагает пользоваться указанием **закона распределения вероятностей** различных значений этой величины.

При этом различают два вида описания законов распределения:

1. **интегральный закон распределения вероятности;**
2. **дифференциальный закон распределения вероятности.**

Определение 1.3

Интегральным законом, или функцией распределения вероятностей $F(x)$ случайной величины X , называют функцию, значения которой представляют вероятность того, что значения x_k случайной величины X меньше некоторого значения x (x – некоторое произвольное число).

Данное утверждение символически записывается в виде

$$F(x) = \text{Prob}(x_k < x), \quad (1.3)$$

где $\text{Prob}(x_k < x)$ и представляет собой вероятность события в вышеприведенном определении.

Очевидно, что

$$F(a) \leq F(b), \quad (1.4)$$

при $a \leq b$ (неубывающая функция)

$$F(-\infty) = 0, \quad F(\infty) = 1. \quad (1.5)$$

Для дискретной одномерной случайной величины X закон распределения удобно представить в виде таблицы

$$X = \begin{pmatrix} x_1, x_2, \dots, x_n \\ p_1, p_2, \dots, p_n \end{pmatrix} \quad (1.6)$$

где x_1, x_2, \dots, x_n – значения случайной величины X , а p_1, p_2, \dots, p_n – вероятности появления этих значений.

Другими словами $\text{Prob}(X = x_i) = p_i$. В этом случае интегральный закон вероятности в соответствии с законом распределения вероятности имеет вид:

$$F(x) = \text{Prob}(x_k < x) = \sum_{i=1}^k p_i. \quad (1.7)$$

1. 3. Дифференциальный закон распределения вероятности

Для случайной величины с непрерывной и дифференцируемой функцией распределения $F(x)$ можно найти дифференциальный закон распределения вероятностей:

$$p(x) = \frac{dF(x)}{dx}. \quad (1.8)$$

$p(x)$ называют **красой** плотности распределения вероятностей (или просто **плотность вероятности**.)

Свойства $p(x)$:Свойства $p(x)$

1. $p(x) \geq 0$
2. Из (1.8) следует, что

$$F(x) = \int_{-\infty}^x p(\xi) d\xi . \quad (1.9)$$

3. Условие нормировки:

$$\int_{-\infty}^{\infty} p(x) dx = 1 . \quad (1.10)$$

Для дискретной случайной величины из определения (1.6) следует, что

$$p(x) = \sum_{i=1}^n \delta(x - x_i) p_i, \quad (1.11)$$

где одномерная δ -функция Дирака определена соотношениями:

$$\delta(x) = \begin{cases} +\infty, & \text{если } x = 0 \\ 0, & \text{если } x \neq 0 \end{cases},$$
$$\int_{-\infty}^{\infty} \delta(x) dx = 1. \quad (1.12)$$

1. 4. Квантили распределений

Определение 1.4

Квантилем случайной величины x вероятности q (или уровня значимости $\alpha = 1 - q$), называется величина x_q , для которой имеем:

$$F(x_q) = \text{Prob}(x < x_q) = \int_{-\infty}^{x_q} p(x) dx = q \quad (1.13)$$

или

$$\text{Prob}(x > x_q) = \int_{x_q}^{\infty} p(x) dx = \alpha = 1 - q. \quad (1.14)$$

Отметим, что с помощью законов распределений вероятности можно найти вероятность нахождения случайной величины в некотором интервале $[a, b]$:

$$\text{Prob}(a < x < b) = \int_a^b p(x) dx = F(b) - F(a). \quad (1.15)$$

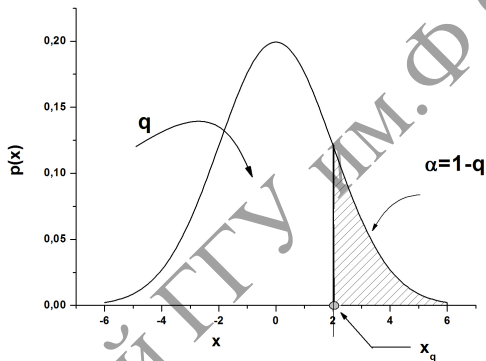


Рисунок 1: Иллюстрация квантиля распределения $p(x)$

1. 5. Примеры законов распределений

Нормальное распределение

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(x-a)^2}{2\sigma^2}\right], \quad (1.16)$$

a и σ - параметры нормального распределения.

Рисунок 2: Нормальное распределение с $a = 0$ и $\sigma = 1$

Распределение χ^2 с n степенями свободы

$$p(\chi^2) = \frac{(\chi^2)^{(\frac{n}{2}-1)} e^{-\frac{\chi^2}{2}}}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})}, \quad (\chi^2 > 0), \quad (1.17)$$

n – параметры χ^2 распределения.

Рисунок 3: Распределение χ^2 с $n = 3, 4, 5, 6$ степенями свободы

Равномерное распределение

$$p(x) = \frac{1}{b-a}, (b > a), \quad (1.18)$$

a, b - параметры распределения.

Рисунок 4: Равномерное распределение с $a = 0$ и $b = 1, 2, 4$

Дискретные распределения

Дискретные распределения:**Распределение Пуассона:**

$$p(x) = \frac{e^{-\mu} \mu^k}{k!} \quad k \geq 0, \quad (1.19)$$

 μ - параметр распределения, а переменная $k = 0, 1, \dots$

2. Моменты распределения

Моменты k -того порядка для непрерывной случайной величины записываются в виде:

$$\mu_k = \int_{-\infty}^{\infty} x^k p(x) dx, \quad (2.1)$$

где μ_k – алгебраический момент k -того порядка.

$$\nu_k = \int_{-\infty}^{\infty} (x - \mu_1)^k p(x) dx, \quad (2.2)$$

где ν_k – центральный момент k -того порядка.

2. 1. Математическое ожидание. Дисперсия

Первый алгебраический момент называют **математическим ожиданием**:

Определение 2.1

Математическим ожиданием непрерывной случайной величины с плотностью вероятности $p(x)$ называется величина $E\{x\}$, определяемая соотношением:

$$\mu_1 = E\{x\} = \int_{-\infty}^{\infty} x p(x) dx . \quad (2.3)$$

Для функции случайной величины $g(x)$ соотношение (2.3) обобщается следующим образом:

$$E\{g(x)\} = \int_{-\infty}^{\infty} g(x)p(x)dx . \quad (2.4)$$

Иногда, для сокращения записи, используем следующее обозначение для математического ожидания:

$$E\{g(x)\} \equiv \langle g(x) \rangle . \quad (2.5)$$

Используя (1.11), (2.3) и свойство δ -функции

$$\int_{-\infty}^{\infty} f(x)\delta(x-a)dx = f(a) \quad (2.6)$$

для дискретной случайной величины получаем, что математическое ожидание вычисляется по формуле:

$$E\{x\} = \sum_{i=1}^n p_i x_i . \quad (2.7)$$

Свойства математического ожидания (2.4):

1

$$E\{c\} = c, \text{ если } c = \text{const}; \quad (2.8)$$

2

$$E\left\{\sum_{k=1}^n c_k g_k(x)\right\} = \sum_{k=1}^n c_k E\{g_k(x)\}, \text{ если } c_k = \text{const}; \quad (2.9)$$

3

$$E\left\{\sum_{k=1}^n c_k \xi_k\right\} = \sum_{k=1}^n c_k E\{\xi_k\}, \text{ если } \xi_k \text{ случайные величины}. \quad (2.10)$$

Математическое ожидание характеризует центр распределения $p(x)$. Однако следует отметить, что не для всех распределений существует математическое ожидание.

Наиболее общей характеристикой центра распределения следует считать медиану.

Определение 2.2

Медиана это такое значение случайной величины x_m для которой вероятности появления различных значений случайной величины X $p_1 = \text{Prob}(X < x_m)$ и $p_2 = \text{Prob}(X > x_m)$ равны между собой, т. е. $p_1 = p_2 = 0,5$.

Другими словами можно сказать, что медиана это квантиль распределения вероятности $q = 0,5$.

Также центр распределения может характеризоваться модой распределения.

Определение 2.3

Мода распределения это такое значение случайной величины x_{mod} для которой плотность вероятности появления значения случайной величины X максимальна, т. е. $p(x_{mod}) = \max p(x)$.

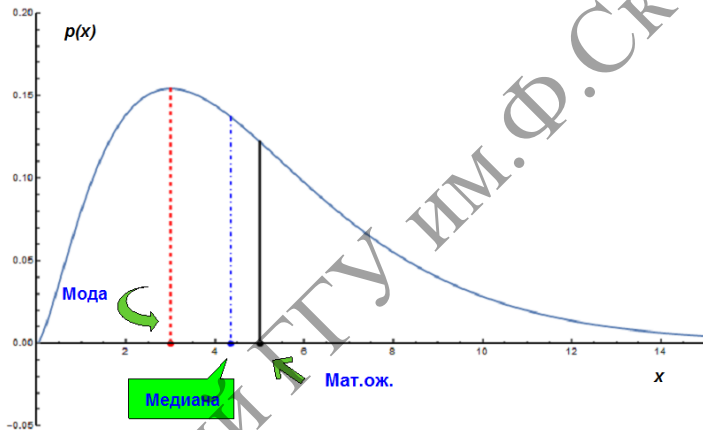


Рисунок 5: Пример распределения

Определение 2.4

Дисперсией непрерывной случайной величины с плотностью вероятности $p(x)$ называется величина $D(x)$, определяемая соотношением:

$$\nu_2 \equiv D\{x\} = \int_{-\infty}^{\infty} (x - E\{x\})^2 p(x) dx . \quad (2.11)$$

Из определения (2.4) следует, что

$$D\{x\} = E\{(x - \mu_1)^2\} = E\{x^2\} - E^2\{x\} . \quad (2.12)$$

Свойства дисперсии

Свойства дисперсии (2.12):

1

$$D\{c\} = 0, \text{ если } c = \text{const}; \quad (2.13)$$

2

$$D\{cx\} = c^2 D\{x\}, \text{ если } c = \text{const}; \quad (2.14)$$

3

$$D\{c+x\} = D\{x\}, \text{ если } c = \text{const}. \quad (2.15)$$

Для дискретной случайной величины, с учетом (1.11) дисперсия вычисляется по формуле:

$$D\{x\} = \sum_{k=1} (x_k - E\{x\})^2 p_k. \quad (2.16)$$

Дисперсия характеризует рассеяние отдельных значений случайной величины от центра распределения и она определяет форму распределения. Дисперсия имеет размерность квадрата случайной величины и выражает как бы мощность рассеяния, поэтому для более наглядной характеристики используют среднеквадратичное отклонение σ :

$$\sigma_x = +\sqrt{D\{x\}}. \quad (2.17)$$

которое имеет размерность самой случайной величины.

Для описания относительной меры отклонения используют коэффициент вариации

$$V_x = \frac{\sigma_x}{E\{x\}}. \quad (2.18)$$

2. 2. Коэффициент асимметрии. Эксцесс

Центральные моменты 3,4 порядка дают информацию о виде кривой плотности распределения.

Определение 2.5

Третий центральный момент непрерывной случайной величины с плотностью вероятности $p(x)$, который определяется соотношением:

$$\nu_3 \equiv \int_{-\infty}^{\infty} (x - E\{x\})^3 p(x) dx, \quad (2.19)$$

характеризует асимметрию кривой $p(x)$ или скошенность распределения (например, когда один спад - крутой, а другой - пологий).

Для симметричных относительно центра распределений он равен нулю. ν_3 имеет размерность куба случайной величины, поэтому для относительной характеристики используют безразмерный коэффициент асимметрии:

$$\gamma_1 = \frac{\nu_3}{\nu_2^{3/2}}. \quad (2.20)$$

Четвертый центральный момент

$$\nu_4 \equiv \int_{-\infty}^{\infty} (x - E\{x\})^4 p(x) dx, \quad (2.21)$$

характеризует протяженность распределения (островершинность).

Определение 2.6

Безразмерную величину вида

$$\varepsilon = \frac{\nu_4}{\nu_2^2}$$

называют **эксцессом** распределения.

Область изменения эксцесса: $\varepsilon \in [1, \infty]$. Часто используют **коэффициент эксцесса**

$$\gamma_2 = \varepsilon - 3 \quad (2.23)$$

и **контрэксцесс** $\kappa = 1/\sqrt{\varepsilon}$.

3. Оценки параметров распределений

Определение 3.1

Функция результатов опытов, которая зависит от неизвестных статистических характеристик называют **статистикой**. Статистика зависит от случайных величин и сама является случайной величиной.

Для нахождения поведения случайной величины, полученных в результате эксперимента необходима числовая информация о моментах распределения.

Оценкой статистической характеристики $\tilde{\theta}$ называется статистика, которая принимается за неизвестное истинное значение параметра θ . Основное требование к оценке истинного значения состоит в том, чтобы большинство значений статистики сосредоточилось вблизи значений θ и вероятность больших отклонений от этого значения была мала. Также желательно, чтобы с увеличением объема выборки точность оценок также увеличилась.

Если имеется экспериментальная выборка объема n случайной величины $X = \{x_1, x_2, \dots, x_n\}$, то ее элементы можно рассматривать как n статистически независимых случайных величин. Тогда любая оценка $\tilde{\theta}$ параметра θ должна быть функцией элементов выборки, т. е.

$$\tilde{\theta} = \theta(x_1, x_2, \dots, x_n) . \quad (3.1)$$

При этом закон распределения $\tilde{\theta}$ зависит от закона распределения величины X и от объема выборки.

Существуют два вида оценок параметров распределения: *точечные* и *интервальные*.

3. 1. Точечные оценки

Требования к оценкам

Под точечной оценкой параметра распределения понимают оценку одним числом. К точечным оценкам предъявляют следующие требования:

- **состоятельность**;
- **несмещенность**;
- **эффективность**;
- **устойчивость**;
- **надежность**.

Определение 3.2

Метод оценки параметров называется **состоятельным**, если оценки, полученные с его помощью, сходятся к истинному значению с увеличением объема выборки n

Определение 3.3

Если $\tilde{\theta}$ оценка параметра θ , то оценка будет **несмещенной**, если смещение

$$E\{\tilde{\theta}\} - \theta \quad (3.2)$$

равно нулю (0), т. е. $E\{\tilde{\theta}\} = \theta$.

Определение 3.4

Оценка $\tilde{\theta}_{eff}$ является эффективной, если она обладает наименьшим разбросом относительно истинного значения параметра θ , т. е.

$$D\{\tilde{\theta}_{eff}\} \rightarrow \min\{\tilde{\theta}_1, \tilde{\theta}_2, \dots\}. \quad (3.3)$$

Определение 3.5

Под **устойчивостью** оценки понимают нечувствительность ее к малым отклонениям от точного распределения.

Точечная оценка для математического ожидания (центра распределения) дается выражением:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad (3.4)$$

Для дисперсии $D = \sigma^2$ такая оценка дается выражением:

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2, \quad S = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}. \quad (3.5)$$

Величины \bar{x} и S сами являются случайными величинами и, следовательно, они тоже могут иметь разброс, который характеризуется дисперсией.

Для \bar{x} :

$$S_{\bar{x}}^2 = \frac{S^2}{n}, \quad S_{\bar{x}} = \frac{S}{\sqrt{n}}. \quad (3.6)$$

Для дисперсии:

$$D[S^2] = \frac{\nu_4}{n} - \frac{\sigma^4(n-3)}{n(n-1)}, \quad (3.7)$$

где ν_4 – четвертый центральный момент, а σ – среднеквадратичное отклонение.

Несмещенными и состоятельными оценками $\tilde{\nu}_3$ и $\tilde{\nu}_4$ для третьего ν_3 и четвертого ν_4 центрального моментов соответственно являются

$$\tilde{\nu}_3 = \frac{n^2}{(n-1)(n-2)} m_3, \quad (3.8)$$

$$\begin{aligned} \tilde{\nu}_4 &= \frac{1}{(n-1)(n-2)(n-3)} \times \\ &\times [n(n^2 - 2n + 3) m_4 - 3n(2n-3) m_3^2], \end{aligned} \quad (3.9)$$

где

$$m_k = \sum_{i=1}^m (x_i - \bar{x})^k. \quad (3.10)$$

Тогда оценки (знак $\tilde{\cdot}$) для коэффициента асимметрии γ_1 и эксцесса ε (смотри соотношения (2.20) и (2.6)) определяются формулами:

$$\begin{aligned}\tilde{\gamma}_1 &= \frac{\tilde{\nu}_3}{S^3}, \\ \tilde{\varepsilon} &= \frac{\tilde{\nu}_4}{S^4},\end{aligned}\tag{3.11}$$

где S – точечная оценка $\sqrt{D\{x\}}$.

Примечание

Примечание.

В теории погрешностей, которая составляет основу обработки данных в лабораторных работах вместо обозначения S используют Δx . Эту величину называют **абсолютной погрешностью**, а также **стандартным отклонением**. Коэффициент вариации V_x называют относительной погрешностью и выражают в процентах:

$$V_x \rightarrow \epsilon_x = \frac{\Delta x}{\bar{x}} \times 100\% . \quad (3.12)$$

3. 2. Интервальные оценки

Точечные оценки параметров случайных величин не позволяют судить о степени близости выборочных значений к оцениваемому параметру. Более содержательны процедуры оценивания параметров, связанные с построением интервала с известной степенью доверительности.

Для любой случайной величины X можно определить вероятности попадания в некоторый интервал $[x_{min}, x_{max}]$, если известен закон распределения вероятностей (смотри (1.15)):

$$\text{Prob}(x_{min} < x < x_{max}) = \int_{x_{min}}^{x_{max}} p(x, \theta) dx = F(x_{max}) - F(x_{min}), \quad (3.13)$$

где в плотность распределения введена величина θ , которая определяет параметры распределения.

Потребуем, чтобы вероятность попадания равнялась некоторому значению $P = 1 - \alpha$, т. е.

$$\text{Prob}(x_{min} < x < x_{max}) = P = 1 - \alpha. \quad (3.14)$$

С помощью (3.14) можно “построить” три вида интервалов с заданной вероятностью попадания p :

- 1 Верхний односторонний (левосторонний) интервал $]-\infty, x_{max}]$;
- 2 Нижний односторонний (правосторонний) интервал $[x_{min}, \infty[$;
- 3 Двусторонний интервал $[x_{min}, x_{max}]$.

Предполагается, что вероятность попадания в каждый их интервалов равна $P = 1 - \alpha$. Обычно величину P называют доверительной вероятностью (иногда надежностью), а величину α уровнем значимости интервала.

Интервальные оценки для моментов распределения находятся построением некоторой функции случайной величины, куда входят искомые параметры распределений θ и точечные оценки $\tilde{\theta}$. Тогда (3.14) применительно к параметрам можно записать

$$\text{Prob} \left(\tilde{\theta} - \delta_{min} < \theta < \tilde{\theta} + \delta_{max} \right) = P = 1 - \alpha, \quad (3.15)$$

который следует понимать так: вероятность того, чтобы параметр θ находится в интервале $\left[\tilde{\theta} - \delta_{min}, \tilde{\theta} + \delta_{max} \right]$ равна $P = 1 - \alpha$.

Отметим, что задача построения доверительных интервалов для параметров при произвольном объеме “экспериментальной” выборки n разработана только для нормального распределения случайной величины X . Для других распределений имеются только частные случаи.

Например, для любой выборки n из нормальной совокупности с математическим ожиданием ξ функция вида

$$t = \frac{\bar{x} - \xi}{(S/\sqrt{n})} \quad (3.16)$$

имеет t -распределение (распределение Стьюдента) с $n - 1$ степенями свободы, плотность которого определяется соотношением:

$$p_{ST}(t) = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{\pi} n \Gamma(\frac{n}{2})} \left[1 + \frac{t^2}{n} \right]^{-\frac{n+1}{2}}. \quad (3.17)$$

Тогда можно найти такое значение $t_{n,P}$ случайной величины с распределением Стьюдента для которой выполняется

$$\text{Prob} \left(\left| \frac{\bar{x} - \xi}{S/\sqrt{n}} \right| < t_{n-1,P} \right) = P = 1 - \alpha \quad (3.18)$$

или

$$\text{Prob} \left(\bar{x} - \frac{t_{n-1,P} S}{\sqrt{n}} < \xi < \bar{x} + \frac{t_{n-1,P} S}{\sqrt{n}} \right) = P = 1 - \alpha . \quad (3.19)$$

В итоге имеем

Доверительный интервал для математического ожидания
 $E\{x\} = \xi$ (двусторонний), если неизвестна дисперсия
 распределения:

$$\left[\bar{x} - \frac{t_{n-1, 1-\alpha/2} S}{\sqrt{n}} < \xi < \bar{x} + \frac{t_{n-1, 1-\alpha/2} S}{\sqrt{n}} \right], \quad (3.20)$$

где $t_{n,P}$ определяется соотношением:

$$\int_{-\infty}^{t_{n,P}} p_{ST}(t) dt = P.$$

Коэффициенты $t_{n,P}$ носят название **коэффициентов Стьюдента**.

Для построения доверительного интервала для дисперсии $D\{x\} = \sigma^2$ используется, то, что функция

$$\frac{(n-1) S^2}{\sigma^2} \quad (3.21)$$

имеет распределение χ^2 с $n - 1$ степенями свободы.

Тогда **доверительный интервал** для дисперсии при неизвестном математическом ожидании имеет вид:

$$\left[\frac{(n-1) S^2}{\chi_{n-1, \alpha/2}^2} < D\{x\} < \frac{(n-1) S^2}{\chi_{n-1, 1-\alpha/2}^2} \right], \quad (3.22)$$

где $\chi_{n, \alpha}^2$ – квантиль распределения χ^2 :

$$\int_{\chi_{n, \alpha}^2}^{\infty} p(\chi^2) d\chi^2 = \alpha$$

с плотностью

$$p(\chi^2) = \frac{(\chi^2)^{(\frac{n}{2}-1)} e^{-\frac{\chi^2}{2}}}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})}, \quad (\chi^2 > 0).$$

Значения $E\{x}$ и $D\{x\}$ лежат в интервале с доверительной вероятностью $P = 1 - \alpha$.

3. 3. Доверительный интервал для случайной величины

С помощью точечных оценок, рассмотренных нами выше, результат для случайной величины x обычно записывают в виде:

$$(\bar{x} \pm \Delta x) . \quad (3.23)$$

При этом согласно правилу приведения результатов:

Правило приведения результатов

Последняя значащая цифра в любом приводимом результате для \bar{x} обычно должна быть того же порядка величины (находиться в той же десятичной позиции), что и абсолютная погрешность Δx .

Однако используемые в расчетах числа должны, как правило, содержать на одну значащую цифру больше, чем это оправдано. Это уменьшит неточности, возникающие при округлении чисел. В конце расчета окончательный ответ следует округлить и избавиться от этой добавочной (и незначащей) цифры).

Придать вероятностный смысл интервалу (3.23) можно, зная лишь функцию распределения вероятности $p(x)$ (плотность вероятности). Если предположить, что наша выборка принадлежит генеральной совокупности с нормальным распределением ($E\{x\} = \bar{x}$, $D\{x\} = (\Delta x)^2$), то легко найти, что

$$\begin{aligned}
 & \text{Prob}(\bar{x} - \Delta x < x < \bar{x} + \Delta x) = \\
 & = \frac{1}{\Delta x \sqrt{2\pi}} \int_{\bar{x} - \Delta x}^{\bar{x} + \Delta x} \exp\left[-\frac{(x - \bar{x})^2}{2(\Delta x)^2}\right] dx \\
 & = \left(z = \frac{x - \bar{x}}{\Delta x}\right) = \frac{1}{\sqrt{2\pi}} \int_{-1}^1 \exp\left[-\frac{z^2}{2}\right] dz = \\
 & = \text{erf}\left(\frac{1}{\sqrt{2}}\right) \approx 0,682689. \tag{3.24}
 \end{aligned}$$

Поэтому в физических приложениях, часто говорят: вероятность того, что результат измерения окажется в пределах одного стандартного отклонения от истинного результата, составляет 68,3%. Можно сказать и так: вероятность, того измеряемая случайная величина лежит в интервале (3.23) составляет 68,3%.

4. Взвешенное среднее значение

В практических исследованиях часто встречается следующая ситуация: одна и та же величина измеряется в нескольких экспериментах. Это делается для наиболее надежного определения некоторой величины. При этом собирают измерения различного происхождения, выполненные разными установками (инструментами) и методами. Результаты таких измерений называют часто называют **неравноточными**.

Предположим, что у нас есть n отдельных измерений

$$x_1 \pm \Delta x_1, \quad x_2 \pm \Delta x_2, \quad \dots, \quad x_n \pm \Delta x_n, \quad (4.1)$$

с соответствующими погрешностями $\Delta x_1, \dots, \Delta x_n$.

Определение 4.1

Наилучшая оценка величины X , основанная на таких измерениях, равна в средневзвешенному значению \hat{x}_{sw}

$$\hat{x}_{sw} = \frac{1}{w} \sum_{i=1}^n w_i x_i, \quad (4.2)$$

$$w = \sum_{i=1}^n w_i, \quad w_i = \frac{1}{(\Delta x_i)^2}, \quad (4.3)$$

а её среднеквадратичное отклонение $\Delta \hat{x}_{sw}$

$$\Delta \hat{x}_{sw} = \frac{1}{\sqrt{w}}. \quad (4.4)$$

Результат такой обработки данных записывается в обычном виде:

$$\hat{x}_{sw} \pm \Delta \hat{x}_{sw} . \quad (4.5)$$

5. Ошибка косвенного измеряемой величины

Все, о чем мы говорили выше, относилось непосредственно к измеряемым величинам. А как определить погрешности для косвенно измеряемых величин, т. е. величин измеряемых с помощью физических законов из непосредственно измеряемых? Оказывается, погрешность косвенно измеряемых величин, связана с погрешностями непосредственно измеряемых величин.

Пусть косвенно измеряемая величина Z связана с непосредственно измеряемыми x_1, \dots, x_m (m величин) величинами функцией $Z = f(x_1, \dots, x_m)$. Если случайные величины x_i **не коррелированы друг с другом, т. е. независимы**, то абсолютная погрешность величины Z определяется соотношением

$$\sigma_Z = \Delta Z = \sqrt{\sum_{i=1}^m \left(\frac{\partial f}{\partial x_i}\right)^2 (\Delta x_i)^2}, \quad (5.1)$$

где $\Delta x_1, \dots, \Delta x_m$ среднеквадратичные отклонения непосредственно измеряемых величин.

Пример. $Z = f(x, y) = x \pm y$

$$\begin{aligned} \Delta Z^2 &= \underbrace{\left(\frac{\partial f}{\partial x}\right)^2}_{=1} (\Delta x)^2 + \underbrace{\left(\frac{\partial f}{\partial y}\right)^2}_{=1} (\Delta y)^2 = \\ &= (\Delta x)^2 + (\Delta y)^2. \end{aligned} \quad (5.2)$$

Практическое следствие этого соотношения: для создания оптимальных условий (условий с минимальной погрешностью) основные усилия должны быть направлены не на дальнейшее уточнение тех результатов измерений, которые являются наиболее точными, а на совершенствование наименее точных измерений случайных величин.

6. Графическое представление эмпирических данных

Цель обработки данных заключается в выявлении вида распределений случайных величин и оценки параметров установленного распределения.

Полученные экспериментальные данные представляют, как правило, в виде таблиц. Полученные таблицы удобно представить графически. Используя набор независимых наблюдений x_1, x_2, \dots, x_n случайной величины X , полезным первым шагом в исследовании поведения случайной величины является организация и представление их таким образом, чтобы их можно было легко интерпретировать и оценивать. Для достаточно большого количества наблюдаемых данных, полигон частот (распределения), гистограмма и кумулятивная линия является отличным графическим представлением данных, что облегчает оценку адекватности предполагаемой модели и оценку параметров распределения.

Гистограмма и полигон распределений являются графическим отображением частот, которые, в свою очередь, представляют собой оценки плотностей вероятностей $p(x)$. Кумулятивная линия - это график накопленных частот, в свою очередь оценивающих интегральную функцию распределения $F(x)$.

Как строить полигон частот

- 1 Построить вариационный ряд для выборки, т. е. упорядочить значения случайной величины так, чтобы выполнялось условие:
 $x_1 \leq x_2 \leq \dots \leq x_n$.
- 2 Разбить область на m интервалов (бинов).
- 3 Построить точки на плоскости x - ν с координатами $\{\tilde{x}_i, \nu_i\}$, ($i = 1, \dots, n$), где

$$\tilde{x}_i = \frac{x_i + x_{i+1}}{2} \quad (6.1)$$

является серединой i -того интервала, ν_i - частота попадания случайной величины X в i -тый интервал.

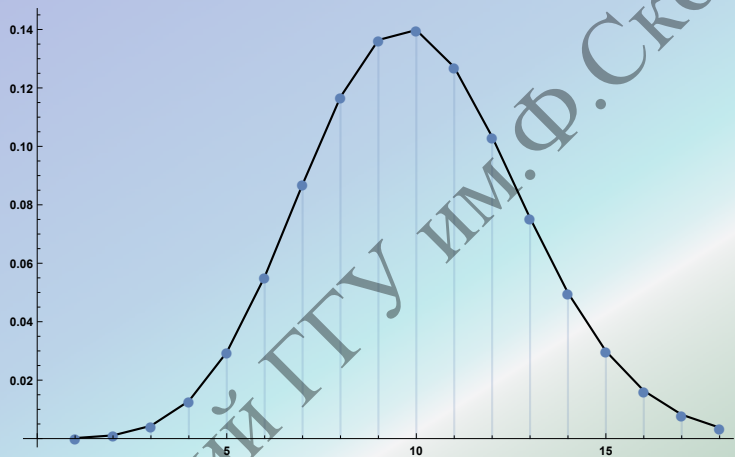


Рисунок 6: Пример полигона частот

Это же распределение можно представить в виде гистограммы. Для построения гистограммы необходимо над каждым отрезком оси абсцисс, соответствующим интервалу значений измеряемой величины, построить прямоугольник, площадь которого пропорциональна частоте попадания в этот интервал. Обычно выбирают интервалы одинаковой ширины, поэтому высота прямоугольников различна.

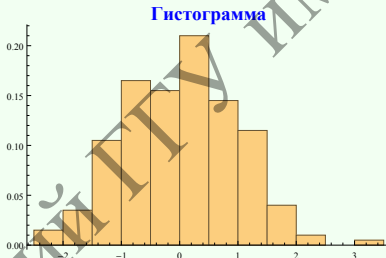


Рисунок 7: Пример гистограммы

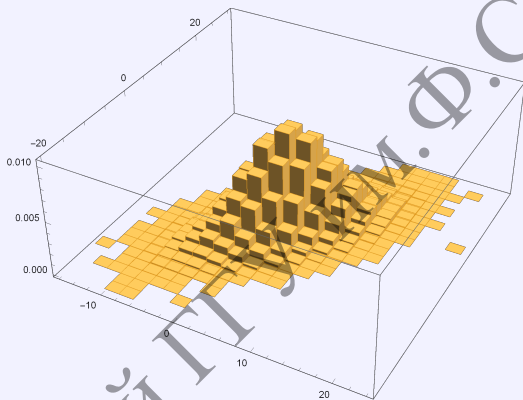


Рисунок 8: Пример 3D-гистограммы

Где используются обработка данных для построения гистограммы?

Для определения оценок математического ожидания, с.к.о., эксцесса не требуется какого-либо группирования данных.

Для определения медианы, сгибов, использования критерия согласия Колмогорова-Смирнова или для обнаружения промахов, экспериментальные данные необходимо расположить в порядке возрастания, т.е. построить вариационный ряд (упорядоченную выборку).

Для определения формы распределения, для использования критериев согласия Пирсона и др., для сопоставления гипотез о форме распределения и т.д. простого упорядочения выборки уже недостаточно, а выборка должна быть представлена в виде гистограммы, состоящей из m столбцов с определенной протяженностью d соответствующих им интервалов.

7. Оптимальное число интервалов для получения гистограммы

Как выбрать m и d .

Оптимальное число

Оптимальное число существует!!

Оптимальное число интервалов группирования это такое число, когда ступенчатая огибающая гистограммы наиболее близка к плавной кривой распределения случайной величины.

К примеру, при группировании данных выбор большого числа интервалов ($m > n$), автоматически приведет к тому, что некоторые из них окажутся пустыми или малозначительными. Гистограмма будет отличаться от плавной кривой распределения вследствие изрезанности многими всплесками и провалами.

Тогда можно сформулировать 1-е требование к числу интервалов:
Размер интервала (ячейки, бина) должен быть достаточно широким для обеспечения хороших статистических свойств будущей гистограммы (достаточно большая статистика, минимальные корреляции (связи) с соседними ячейками).

При слишком малом числе m интервалов, гистограмма отличается от действительной кривой распределения вследствие слишком крупной ступенчатости. Из-за чего будут потеряны характерные особенности. Например, если взять $m = 1$, т.е. d равно размаху экспериментальных данных, то любое распределение сводится к равномерному, а если $m = 3$, то любое куполообразное распределение сведется к треугольному. Например, при обработке линейчатых спектров, слишком большая ячейка может привести к потере спектральной линии.

Тогда возникает **2-е требование к числу интервалов**:

Размер ячейки должен быть достаточно узким для того, чтобы прорисовывалась “тонкая структура” исследуемой величины.

Как видим, требования являются противоречивыми. Укрупнение интервалов группирования является методом “фильтрации различных случайных выбросов и провалов”, но слишком протяженные интервалы сглаживают особенности искомого закона распределения. Таким образом, задача выбора оптимального числа интервалов при построении гистограммы – это задача оптимальной фильтрации, а оптимальным числом m интервалов является максимальное возможное сглаживание случайных флуктуаций данных, которое сочетается с минимальным искажением от сглаживания самой кривой искомого распределения.

Общепринято делать интервалы одинаковыми. Хотя в дальнейшем увидим, что это условие необязательно. Условие равновеликости интервалов удобно с практической точки зрения.

Рекомендации по выбору m .

I группа: эвристические критерии (без доказательства).

Формула Старджеса

$$m = \log_2 n + 1 . \quad (7.1)$$

Формула Брукса и Каррузера

$$m = 5 \lg n . \quad (7.2)$$

Формула если $n > 100$

$$m = \sqrt{n}. \quad (7.3)$$

Эти три формулы являются наиболее часто встречающимися в литературе по математической статистике.

II группа: с использованием критерия χ^2 .

В ней используется рассмотрение интервалов не с равной длиной, а с **равной вероятностью** в соответствии с принимаемой моделью, т. е. предположением о законе распределения. В данном подходе неявно учитывается форма распределения.

Число интервалов с равной вероятностью, которые мы обозначили как K , отличаются от числа m с равной длиной d .

Г. Манн и А. Ваальд установили, что при $n \rightarrow \infty$ оптимальное число K равновероятных интервалов задается соотношением:

Критерий Мана-Ваальда

$$K \sim b\sqrt[5]{2} \left(\frac{n}{Z_\alpha} \right)^{2/5}, \quad (7.4)$$

где $b = 2 \div 4$.

Здесь Z_α – квантиль нормального распределения, соответствующий вероятности $P = 1 - \alpha$, α – принятый уровень значимости.

$$Z_\alpha = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{Z_\alpha} e^{-\frac{x^2}{2}} dx = 1 - \alpha.$$

Критерий Мана-Ваальда

На практике часто берут $\alpha = 0.1$, тогда

$$K \simeq 1,9 n^{2/5} . \quad (7.5)$$

В итоге приходим к таким рекомендациям при использования равновероятностных бинов K :

- 1 найти число ячеек, используя (7.5);
- 2 если окажется, что n/K мало, уменьшить K чтобы выполнялось неравенство $n/K \geq 5$;
- 3 Сформировать K равновероятностных ячеек гистограммы на основе данных. Отметим, что если измеряемая случайная величина многомерная, то существуют различные способы формирования ячеек с одинаковым вероятностным содержанием в K ячейках.

III группа.

Поскольку для K интервалы получаются не равной длины, то это приводит к ряду неудобств при построении гистограмм, но зато при этом мы неявно закладываем при использовании χ^2 выбор K в зависимости от формы распределения.

III группа рекомендаций “устраняет” недостаток II группы возвращаясь к интервалам m с равной длиной d , но при этом и учитывает, в отличие от I группы, форму распределения (форма характеризуется эксцессом ε или контрэксцессом κ).

Примером такого подхода является соотношение:

III группа (формула И.У.Алексеевой)

$$m = \frac{4}{\kappa} \lg \frac{n}{10}, \quad \kappa = \frac{1}{\sqrt{\varepsilon}}. \quad (7.6)$$

На практике эту формулу в зависимости от n при $\kappa = const$ удобнее аппроксимировать выражением

$$m = \frac{\varepsilon + 1,5}{6} n^{2/5} \quad (7.7)$$

Трудность использования III группы состоит в том, что число интервалов часто приходится выбирать прежде, чем будут найдены оценки ε , \bar{x} , и т.д.

Эту трудность обходят следующим образом: наиболее часто встречаются распределения с ε от 1,8 до 6 (от равномерного до Лапласа, включая нормальное $\varepsilon=3$). Для этих граничных точек имеем

$$m_{min} = 0,55 n^{2/5} \quad (7.8)$$

$$m_{max} = 1,25 n^{2/5} \quad (7.9)$$

Искомое m можно выбрать близким к этому интервалу, при этом m лучше выбрать нечетным, т.к. при четном m для островершинных распределений в центре гистограммы оказывается два столбца равных по высоте и середина распределения принудительно утолщается.

“Практические” рекомендации для построения диаграмм

- 1 Для практического определения числа интервалов воспользоваться формулами для m_{min} (7.8) и m_{max} (7.9), или сразу формулой (7.6), выбрав при этом m нечетным.
- 2 Так как крайние точки могут располагаться несимметрично, то ширина d столбца гистограммы определяется по отклонению от центра ΔX_m наиболее удаленной точки:

$$d = \frac{2\Delta X_m}{m} . \quad (7.10)$$

При этом полученное значение d необходимо округлять в большую сторону, чтобы крайняя точка не оказалась за пределами крайнего столбца.

- 3 Величину d при этом удобно выбирать так, чтобы она делилась на 2 так, чтобы потом центральный столбец можно было бы поделить пополам для уточнения центра распределения.

8. Оптимальное число интервалов для получения гистограммы

Одним из условий правомерности статистической выборки является требование ее однородности, т.е. принадлежности всех ее членов к одной и той же генеральной совокупности.

Однако на практике это требование очень часто нарушается. И, если скажем, при обработке вручную еще можно вспомнить как (при каких условиях) были получены “подозрительные” данные, то при автоматической обработке данных необходимы методы исключения “чужих” для данной выборки результатов.

Определение 9.1

Отсчёты, резко отклоняющиеся по своим значениям от большинства других отсчетов принято называть **промахами** и исключать их из выборки.

Важно. Если серия из небольшого числа измерений содержит грубую погрешность — промах, то наличие этого промаха может сильно исказить как среднее значение измеряемой величины, так и границы доверительного интервала. Поэтому из окончательного результата необходимо исключить этот промах.

Обычно промах имеет резко отличающееся от других измерений значение. Однако это отклонение от значений других измерений не дает еще права исключить это измерение как промах, пока не проверено, не является ли это отклонение следствием статистического разброса.

Особую неприятность доставляют отсчеты, которые и не входят в компактную группу отсчетов, но и не удалены от нее на значительное расстояние. Такой отсчет называют предполагаемым промахом.

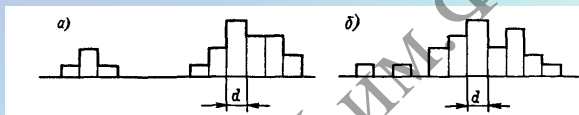


Рисунок 9: Возможные промахи

В экспериментальной практике исследователи просто отбрасывали крайние, “слишком удаленные от центра наблюдения”. Эта процедура получила название **цензурирование выборки**.

Однако для принятия решения необходимы какие-либо формальные критерии.

Простейший метод заключается в использовании “правила 3σ ”, когда по выборке с удаленными отсчётами (предполагаемыми промахами) вычисляется оценка σ и граница $|X_{\text{group}}| = 3\sigma$, а все $|x_i| \pm 3\sigma$ отбрасываются.

Правило 3σ обосновано на неравенстве Чебышева:

$$\text{Prob}\{|x_i - \bar{x}| \geq a\} \leq \frac{\sigma^2}{a^2}, \quad (a > 0) \quad (9.1)$$

Если все $x_i \geq 0$, то

$$\text{Prob}(x \geq a) \leq \frac{\bar{x}}{a}, \quad (a > 0) \quad (9.2)$$

Если x имеет одномодальное распределение (непрерывное), то справедлива более сильная оценка

$$\text{Prob}\{|x_i - \bar{x}| \geq a\} \leq \frac{4}{9} \frac{1 + S^2}{\left(\frac{a}{\sigma} - |S_{pear}|\right)^2}, \quad (9.3)$$

где S_{pear} – пирсоновская мера асимметрии (для распределений симметричных относительно моды $S_{pear} = 0$).

$$S_{pear} = \frac{\bar{x} - \xi}{\sigma},$$

где ξ – максимальная вероятность.

$$S_{pear} = \frac{\gamma_1(\gamma_2 + 6)}{2(5\gamma_2 - 6\gamma_1^2 + 6)},$$

где

$$\gamma_1 = \frac{m_3}{m_2^{3/2}}, \quad \gamma_2 = \frac{m_4}{m_2^2} - 3.$$

Значения m_j определяются по формуле:

$$m_j = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^j.$$

Если положить $a = 3\sigma$, то имеем

$$P\{|x_i - \bar{x}| \geq 3\sigma\} \leq \frac{1}{9} \approx 0,11$$

или для одномодальных симметричных:

$$F\{|x_i - \bar{x}| \geq 3\sigma\} \leq \frac{4}{9} \frac{1}{9} = \frac{4}{81} \approx 0,05,$$

т. е. для произвольного распределения вероятность, что $x_i \geq \bar{x}$ на 3σ составляет 11%, а для одномодальных симметричных 5%.

Чувствительность разных статистических методов к наличию аномальных наблюдений (“промахов”) в экспериментальных данных неодинакова.

Критерий Граббса

$$G > \frac{n-1}{\sqrt{n}} \sqrt{\frac{t_{\alpha/(2n), n-2}^2}{n-2 + t_{\alpha/(2n), n-2}^2}} \quad (9.4)$$

Тест Граббса определяется для гипотезы:

- H_0 : в наборе данных нет выбросов
- H_1 : В наборе данных имеется ровно один выброс.

Тест Граббса основан на предположении о нормальности выборки. То есть сначала нужно проверить, что данные могут быть разумно аппроксимированы нормальным распределением перед применением теста Граббса.

Тест Граббса обнаруживает один выброс за раз. Этот выброс исключается из набора данных, и тест повторяется до тех пор, пока не будут обнаружены выбросы. Тем не менее, множественные итерации изменяют вероятности обнаружения, и тест не должен использоваться для размеров выборок в шесть или меньше, поскольку он часто помещает большинство точек в виде выбросов.

Критерий Роснера для обнаружения нескольких выбросов

$$R_i = \max \frac{x_i - \bar{x}}{S} \quad (9.5)$$

или

$$\tau_1 = \max \left\{ \frac{x_1 - \bar{x}}{S}, \frac{x_n - \bar{x}}{S} \right\}, \text{ если } x_1 < x_2 < \dots < x_n. \quad (9.6)$$

Алгоритм критерия Роснера состоит в следующем. По начальной выборке объема n вычисляются значения \bar{x} и S и статистика τ_1 . Затем из выборки удаляется экстремальный член $x_{\min}(x_{\max})$ —в зависимости от того, какое значение более удалено от среднего. Так повторяется k раз.

Полученные значения статистик τ_{ik} ($i = 1, \dots, k$) каждый раз сравниваются с критическими значениями

$$\tau_{i,n,p}^* = \frac{n-i}{\sqrt{n-i+1}} \sqrt{\frac{t_{p,n-i-1}^2}{n+i-1+t_{p,n-i-1}^2}}, \quad p = 1 - \alpha/(2(n-i+1)) \quad (9.7)$$

для заданных n , k и вероятности p . Превышение критерием τ_{1i} критического значения $\tau_{i,n,p}^*$ для некоторого i , позволяет установить не только наличие выбросов, но и их количество (равное значению i , при котором появляется первая значимая величина критерия τ_{1i}). Вычисление последовательных статистик ведется до тех пор, пока $\tau_{1(i+1)} > \tau_{1i}$.

9. 1. Другие критерии для исключения промахов

Однако, правило “ 3σ ” хорошо для нормального распределения. Действительно, при $n = 100$ появление $|x_i| \geq 3\sigma$ можно считать промахом, то для равномерного промахом можно считать уже $|x_i| = 1.8\sigma$, а для распределения Лапласа $|x_i| = 3\sigma$ есть отсчет, принадлежащий данной выборке.

На этом примере видно, что граница цензурирования зависит не только от объема выборки n , но также и от формы распределения.

Зависимость от n полуколичественно можно оценить из условия: границы цензурирования должны отсекают в среднем менее одной точки, тогда назначение границ с уровнем значимости $g = 1 - P$, где $P = \frac{n}{n+1}$ дает зависимость от n .

Для расчета количества σ для цензурирования можно использовать аппроксимационные формулы:

Формулы для расчета количества σ ()

$$\begin{aligned}t_{гр.} &= 1,2 + 3,6 \left(1 - 1/\sqrt{\varepsilon}\right) \lg \left(\frac{n}{10}\right), \\t_{гр.} &= 1,55 + 0,8\sqrt{\varepsilon - 1} \lg \left(\frac{n}{10}\right)\end{aligned}\quad (9.8)$$

Тогда интервал цензурирования выглядит следующим образом:

$$\bar{x} \pm t_{гр.} S, \quad (9.9)$$

где S - точечная оценки среднеквадратичного отклонения (3.5).

В заключение данного раздела отметим, что все оценки \bar{x} , S и т.д. должны пересчитываться после цензурирования выборки.

10. Критерии согласия

Постановка задачи.

Одной из задач первичной обработки экспериментальных наблюдений является выбор закона распределения, который наиболее хорошо описывающего случайную величину, выборку которой наблюдают.

Рассмотрим любой эксперимент, в котором измеряется некоторая случайная величина X и мы получаем некоторую выборку объема n . Наша задача описать поведение этой случайной величины. Из теории мы знаем о большом количестве функций распределения вероятностей (нормальное, равномерное распределения, распределение Стьюдента и т.д.).

И у нас возникает мысль: А что если одно из известных теоретических (модельных) распределений подходит для описания поведения измеряемой величины? Как быть? Подойдет ли выбранное нами теоретическое распределение или не подойдет? Такая задача в математической статистике называют проверкой гипотезы.

Определение 10.1

Под статистической гипотезой понимают всякое высказывание о случайной величине (генеральной совокупности), проверяемое по выборке (по результатам наблюдений).

Соответственно, процедура сопоставления высказанной гипотезы с выборочными данными называется **проверкой гипотезы**.

Таким образом, **целью проверки гипотезы о согласии** опытного распределения с теоретическим является стремление удостовериться в том, что данная модель теоретического закона не противоречит наблюдаемым данным и использование ее не приведет к существенным ошибкам при вероятностных расчетах.

Ответ, который должен быть получен в конечном итоге выглядит следующим образом: Данная (“экспериментальная”) выборка с вероятностью p соответствует теоретическому распределению с параметрами (например, нормальному распределению с параметрами $a = 3$ и $\sigma = 1$).

Поэтому основная задача на практике состоит в определении вероятности принятия модели (гипотезы) p (или более точно вероятности, что данная выборка не соответствует данной гипотезе $\alpha = 1 - p$).

Что нужно знать о процессе проверки гипотез:

Этап 1. Выбор модели (теоретической функции распределения)

Располагая выборочными данными и руководствуясь конкретными условиями рассматриваемой задачи, формулируют гипотезу H_0 , которую называют **основной** или нулевой, и гипотезу H_1 конкурирующую с гипотезой H_0 .

Термин “конкурирующая” означает, что являются противоположными следующие два события: по выборке будет принято решение о справедливости для генеральной совокупности гипотезы H_0 ; и по выборке будет принято решение о справедливости для генеральной совокупности гипотезы H_1 . Гипотезу H_1 называют также альтернативной.

Например, если нулевая гипотеза такова: данная выборка имеет нормальное распределение с параметрами $a = 3$ и $\sigma = 1$, то альтернативная гипотеза может быть следующей: данная выборка имеет нормальное распределение, но с параметрами $a < 3$ и $\sigma = 1$ или

$$H_0: p(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(x-a)^2}{2\sigma^2}\right], \text{ где}$$
$$a = 3 \text{ и } \sigma = 1, \quad (10.1)$$

и

$$H_1: p(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(x-a)^2}{2\sigma^2}\right], \text{ где}$$
$$a < 3 \text{ и } \sigma = 1. \quad (10.2)$$

Этап 2. Расчет вероятности принятия гипотезы p или вероятности не соответствия данной модели $\alpha = 1 - p$.

Вероятность α часто называют уровнем значимости и если выборка соответствует выбранной нами модели (теоретическому распределению), то уровень значимости относительно мал ($\alpha < 0,05$)
Поясним ее смысл.

Решение о том, можно ли считать высказывание H_0 справедливым для генеральной совокупности, принимается по выборочным данным, т. е. по ограниченному ряду наблюдений, следовательно, это решение может быть ошибочным.

При этом может иметь место ошибка двух родов: отвергают гипотезу H_0 , или, иначе, принимают альтернативную гипотезу H_1 , тогда как на самом деле гипотеза H_0 верна; это **ошибка первого рода**; принимают гипотезу H_0 , тогда как на самом деле высказывание H_0 неверно, т. е. верной является гипотеза H_1 это **ошибка второго рода**.

Так вот уровень значимости α – это вероятность ошибки первого рода, т. е. вероятность того, что будет принята гипотеза H_1 , если на самом деле в генеральной совокупности верна гипотеза H_0 (или отклонение проверяемой гипотезы H_0 при её справедливости).

Вероятность ошибки второго рода часто обозначают β , т. е. вероятность того, что будет принята гипотеза H_0 , если на самом деле верна гипотеза H_1 (или не отклонении H_0 при справедливости H_1).

О мощности критерия.

При использовании критериев согласия, как правило на практике, не задают конкурирующих гипотез: рассматривается принадлежность выборки конкретному закону. А в качестве конкурирующей гипотезы — принадлежность любому другому.

Естественно, что способность критерия отличать закон, соответствующий H_0 , от других, близких к закону, соответствующему H_0 , и далёких от него, отличаются.

Определение 10.2

Мощностью критерия по отношению к конкурирующей гипотезе H_1 называется величина $1 - \beta$. Критерий тем лучше распознаёт пару конкурирующих гипотез H_0 и H_1 , чем выше его мощность.

Один из возможных вариантов проверки гипотез являются **критерии согласия**.

Критерий χ^2 (критерий Пирсона).

Рассмотрим этапы необходимые для проверки гипотез на примере критерия χ^2 .

Процедура проверки гипотез с использованием критериев типа χ^2 предусматривает группирование наблюдений.

1. Выбираем “модель”: обычно это теоретическое дифференциальное распределение вероятностей $p(x, \theta)$, где θ – один (или несколько) параметров распределения.
2. Область определения случайной величины разбивают на m непересекающихся интервалов граничными точками $x_0, x_1, \dots, x_{m-1}, x_m$, где $x_0 < x_1 < \dots < x_{m-1} < x_m$.

3. В соответствии с заданным разбиением подсчитывают число n_i выборочных значений, попавших в i -й интервал и вероятности попадания в i -й интервал

$$p_i = \int_{x_{i-1}}^{x_i} p(x, \theta) dx = F(x_i) - F(x_{i-1}) \quad (10.3)$$

соответствующие теоретическому закону с интегральной функцией распределения $F(x, \theta)$ для всех m интервалов. При этом $\sum_{i=1}^m n_i = n$ и $\sum_{i=1}^m p_i = 1$.

4. Проводим расчет статистики критерия согласия χ^2 Пирсона с помощью соотношения

$$\chi_{\text{exp}}^2 = n \sum_{i=1}^m \frac{(n_i/n - p_i)^2}{p_i} \quad (10.4)$$

При проверке гипотезы при $n \rightarrow \infty$ для которой известны, как вид закона $p(x, \theta)$, так и все его параметры θ (простая гипотеза) функция χ_{exp}^2 подчиняется распределению χ_r^2 с $r = m - 1$ степенями свободы (доказано в Pearson, Karl (1900). On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling. Philosophical Magazine Series 5 50 (302): 157—175.).

5. Далее находим из уравнения величину

$$\int_{\chi_{\text{exp}}^2}^{\infty} p(s) ds = \int_{\chi_{\text{exp}}^2}^{\infty} \frac{(s)^{\left(\frac{r}{2}-1\right)} e^{-\frac{s}{2}}}{2^{\frac{r}{2}} \Gamma\left(\frac{r}{2}\right)} ds = 1 - F_r(\chi_{\text{exp}}^2) = p, \quad (10.5)$$

где $F_r(x)$ – интегральная функция вероятности распределения χ_r^2 с r степенями свободы.

6. После вычисления p получаем ответ: **Данная выборка с вероятностью p соответствует теоретическому распределению $p(x, \theta)$ с параметрами θ .**

Или: **Данная выборка с вероятностью $\alpha = 1 - p$ не соответствует теоретическому распределению $p(x, \theta)$ с параметрами θ**

Примечания.

- ★ На практике (например с помощью точечных оценок) удается оценить(рассчитать) все или часть параметров распределения. Тогда статистика χ_{exp}^2 при справедливости проверяемой гипотезы подчиняется χ_r^2 -распределению с $r = m - k - 1$ степенями свободы, где k количество оцененных по выборке параметров.
- ★ Некорректное использование критериев согласия (не построен вариационный ряд для выборки, неверно выбрано число интервалов) может приводить к необоснованному принятию (чаще всего) или необоснованному отклонению проверяемой гипотезы.
- ★ Существуют и другие критерия согласия.

$\tilde{\chi}^2$ of d.f.

Существует несколько более удобный способ понимания критерия χ^2 . Введем понятие приведенного значения $\tilde{\chi}^2$ (или $\tilde{\chi}^2$ на одну степень свободы), которое определим как

$$\tilde{\chi}^2 = \frac{\chi^2}{r} \quad (10.6)$$

Тогда, каким бы ни было число степеней свободы, наш критерий можно сформулировать следующим образом: если мы получаем значение $\tilde{\chi}^2$ порядка 1 или меньше, то у нас нет оснований сомневаться в нашем ожидаемом распределении; если мы получаем значение $\tilde{\chi}^2$ много большее, чем единица, то невероятно, чтобы наше ожидаемое распределение было верным. Т.е. если

$$\tilde{\chi}^2 \leq 1, \quad (10.7)$$

то наша выборка соответствует теоретическому распределению.

11. Алгоритм предварительной обработки экспериментальных данных

Исходя из изложенного материала можно сформировать необходимые этапы предварительной обработки наблюдений.

- 1. Вычисление выборочных характеристик (точечные оценки моментов экспериментального распределения) \bar{x} , S^2 , $\tilde{\nu}_3$** (и соответственно $\tilde{\gamma}_1$), $\tilde{\nu}_4$ (и соответственно $\tilde{\varepsilon}$), а также дополнительные характеристики $S_{\bar{x}}$, и $m_{3,4}$ по формулам:

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

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 ,$$

$$m_k = \sum_{i=1}^n (x_i - \bar{x})^k , k = 2, 3, 4 ,$$

$$\tilde{\nu}_3 = \frac{n^2}{(n-1)(n-2)} m_3 ,$$

$$\tilde{\gamma}_1 = \frac{\tilde{\nu}_3}{S^3} ,$$

$$\tilde{\nu}_4 = \frac{(n^2 - 2n + 3) m_4 - 3n(2n-3) m_3^2}{(n-1)(n-2)(n-3)} ,$$

$$S_{\bar{x}} = \frac{S}{\sqrt{n}} .$$

(11.1)

(Смотри формулы (3.4), (3.5), (3.8),(3.9)) и другие.

2. Цензурирование выборки (отсев грубых промахов). Последующий пересчёт точечных оценок.
3. Построение гистограмм и кумулятивной кривой. Визуальный анализ на соответствие теоретической плотности распределения с привлечением точечных оценок.
4. Проверка с помощью критериев согласия на соответствие одной (или нескольких) из теоретических плотностей распределения. Как правило, проверяют соответствие экспериментальной выборки нормальному распределению.

12. Некоторые сведения о двумерных случайных величинах

Переход от одномерных случайных величин к многомерным сопровождается введением новых понятий. Наглядные модельные представления можно построить для двух или трех переменных. Наиболее удобен случай двух переменных, и он чаще других будет использоваться в дальнейшем.

Пусть X и Y две случайные величины, имеющие набор значений x_k и y_k . Индекс k нумерует определенные значения этих величин. Для каждой из одномерных случайных величин можно ввести интегральные $F(x)$, $F(y)$ и дифференциальные функции распределения вероятностей $p(x)$ и $p(y)$. Но можно ввести и совместную функцию распределения вероятностей, описывающие поведение обеих случайных величин, как единый объект (одновременно).

Определение 12.1

Совместная функция распределения $F(x, y)$ - это вероятность, того что значения двумерной величины приписанные выборочному множеству k , удовлетворяют одновременно неравенствам $x(k) \leq x$ и $y(k) \leq y$ или

$$F(x, y) = \text{Prob}(x_k \leq x \text{ и } y_k \leq y) . \quad (12.1)$$

Легко найти, что

1.

$$F(-\infty, y) = 0 , \quad (12.2)$$

2.

$$F(x, -\infty) = 0 , \quad (12.3)$$

3.

$$F(\infty, -\infty) = 1 . \quad (12.4)$$

Совместный дифференциальный закон распределения $p(x, y)$ определяется соотношением:

$$p(x, y) = \frac{\partial^2}{\partial y \partial x} [F(x, y)] \quad (12.5)$$

Очевидно что

1

$$p(x, y) \geq 0, \quad (12.6)$$

2

$$F(x, y) = \int_{-\infty}^y \int_{-\infty}^x p(\xi, \eta) d\xi d\eta, \quad (12.7)$$

Плотность вероятности для случайных величин X и Y в отдельности, выражаются через совместную функцию распределений с помощью соотношений:

$$p(x) = \int_{-\infty}^{\infty} p(x, y) dy, \quad (12.8)$$

$$p(y) = \int_{-\infty}^{\infty} p(x, y) dx. \quad (12.9)$$

Функции (12.8) и (12.9) называют безусловными или маргинальными плотностями распределения случайных величин X и Y .

Определение 12.2

Две случайные величины X и Y являются статистически независимыми (или взаимно независимыми), если

$$p(x, y) = p(x)p(y). \quad (12.10)$$

Математическое ожидание

Если имеется функция случайных величин $G(x, y)$, то ее математическое ожидание

$$E\{G(x, y)\} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} G(x, y)p(x, y) dx dy. \quad (12.11)$$

Дисперсия

Дисперсия $D\{G(x, y)\}$ функция случайных величин $G(x, y)$ выражается следующим образом:

$$D\{G(x, y)\} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (G(x, y) - E\{G(x, y)\})^2 p(x, y) dx dy . \quad (12.12)$$

12. 1. Алгебраические и центральные моменты

Алгебраические имеют соответственно следующий вид (и частном случае для (12.11), когда $G = x$ или $G = y$):

$$E \{x\} = \mu_x = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x p(x, y) dx dy, \quad (12.13)$$

$$E \{y\} = \mu_y = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y p(x, y) dx dy. \quad (12.14)$$

С учетом определений (12.8) и (12.9) для безусловных функций распределения имеем, что

$$E\{x\} = \mu_x = \int_{-\infty}^{\infty} xp(x) dx, \quad (12.15)$$

$$E\{y\} = \mu_y = \int_{-\infty}^{\infty} yp(y) dy. \quad (12.16)$$

Центральные моменты, по аналогии с одномерной случайной величиной, могут быть записаны в виде

$$D\{x\} = \sigma_x^2 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - E\{x\})^2 p(x, y) dx dy, \quad (12.17)$$

$$D\{y\} = \sigma_y^2 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (y - E\{y\})^2 p(x, y) dx dy. \quad (12.18)$$

Но кроме этих простейших моментов второго порядка (12.17) и (12.18), имеющих уже хорошо известную форму дисперсий, появляются новые возможности для образования моментов второго порядка.

Ковариация

Так можно определить дисперсию $D\{x, y\}$ посредством соотношения:

$$D\{x, y\} = \text{cov}(x, y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - E\{x\})(y - E\{y\}) p(x, y) dx dy. \quad (12.19)$$

Определение 12.3

Центральный момент (12.19) называются **ковариацией** между x и y .

Если x и y статистически независимы (см. (12.10)), т. е. $p(x, y) = p(x)p(y)$, тогда ковариация обращается в нуль. Действительно, после очевидных преобразований, имеем, что

$$\begin{aligned} \operatorname{cov}(x, y) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - E\{x\})(y - E\{y\}) p(x) p(y) dx dy = \\ &= \int_{-\infty}^{\infty} p(x)(x - E\{x\}) dx \times \int_{-\infty}^{\infty} (y - E\{y\}) p(y) dy = 0 \end{aligned} \quad (12.20)$$

Во всех других случаях смешанный момент второго порядка не равен нулю. Обычно вводится безразмерная переменная – **коэффициент корреляции** $\rho(x, y)$, определяемый как

$$\rho(x, y) = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y} \quad (12.21)$$

Можно, показать, что

$$-1 \leq \rho(x, y) \leq 1. \quad (12.22)$$

На практике бывает необходимо оценить коэффициент корреляции $\rho(x, y)$ для конкретной выборки объема n . Выборочная ковариация (несмещенная точечная оценка) $r(x, y)$ можно определить по формуле

$$r(x, y) = \frac{1}{n-1} \sum_{i=1}^n \frac{(x_i - \bar{x})(y_i - \bar{y})}{S_x S_y} p_i. \quad (12.23)$$

Здесь \bar{x} и \bar{y} - средние значения случайных переменных; p_i - частота попадания переменных x_i и y_i в i -ый интервал значений, а $S_{x,y}$ - точечные оценки среднеквадратичных отклонений $\sigma_{x,y}$.

Определение 12.4

Ковариация (12.19) имеет смешанный характер и вместе с остальными моментами образует матрицу (12.17) и (12.18), которую называют **ковариационной матрицей**, или матрицей вторых моментов, или дисперсионной матрицей, или матрицей ошибок.

Ковариационная матрица

$$\mathbf{V} = \begin{pmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{pmatrix} = \begin{pmatrix} D\{x\} & D\{x, y\} \\ D\{y, x\} & D\{y\} \end{pmatrix}. \quad (12.24)$$

Поскольку $D\{x, y\} = D\{y, x\}$, то матрица корреляции симметрична, диагональные элементы которой представляют собой дисперсии.

Каждый элемент ковариационной матрицы V_{ij} можно трактовать как математическое ожидание ij -элемента произведения вектор-столбца $(\mathbf{x} - \boldsymbol{\mu})$ на вектор-строку $(\mathbf{x} - \boldsymbol{\mu})^T$, где

$$\begin{aligned}(\mathbf{x} - \boldsymbol{\mu}) &= \begin{pmatrix} x - \mu_x \\ y - \mu_y \end{pmatrix} \\ (\mathbf{x} - \boldsymbol{\mu})^T &= (x - \mu_x \quad y - \mu_y) \end{aligned} \quad (12.25)$$

Теперь в матричной записи (12.24) имеет вид

$$V = E \left\{ (\mathbf{x} - \boldsymbol{\mu}) (\mathbf{x} - \boldsymbol{\mu})^T \right\}. \quad (12.26)$$

Корреляционные связи имеют важное значение в исследованиях. Удобство работы с корреляциями связано с тем, что корреляция - безразмерная величина в отличие от ковариации. Из коэффициентов корреляции можно также построить корреляционную матрицу \mathbf{R} . В двумерном случае, в следствие, того, что

$$\rho(x, x) = \rho(y, y) = 1 \quad (12.27)$$

имеем, что

$$\begin{aligned} \mathbf{R} &= \begin{pmatrix} \rho(x, x) & \rho(x, y) \\ \rho(x, y) & \rho(y, y) \end{pmatrix} = \\ &= \begin{pmatrix} 1 & \frac{\text{cov}(x, y)}{\sigma_x \sigma_y} \\ \frac{\text{cov}(x, y)}{\sigma_x \sigma_y} & 1 \end{pmatrix}. \end{aligned} \quad (12.28)$$

Важно, что если x и y образуют двумерную случайную величину, то дисперсия линейной функции $z = x + y$ находится с помощью соотношения:

$$D\{x + y\} = D\{x\} + D\{y\} + 2\text{cov}(x, y) . \quad (12.29)$$

12. 2. Примеры многомерных функций распределения вероятностей

Пусть имеем многомерную случайную величину $\{X_1, X_2, \dots, X_k\}$ размерности k . Для центральных и алгебраических моментов распределения, а также ковариационной матрицы введем следующие обозначения:

$$\begin{aligned} E\{X_i\} &= \mu_i, i = 1, \dots, k, \\ D\{X_j\} &= \sigma_j^2, \\ \text{cov}(X_i, X_j) &= \rho^{ij} \sigma_i \sigma_j. \end{aligned} \quad (12.30)$$

Введем дополнительные k -мерные векторы: $\mathbf{X} = \{X_1, X_2, \dots, X_k\}$ и $\boldsymbol{\mu} = \{\mu_1, \dots, \mu_k\}$ и матрицу ошибок

$$\mathbf{V}_{ij} = \begin{cases} \sigma_i^2, & \text{если } i = j \\ \rho_{ij} \sigma_i \sigma_j & \text{если } i \neq j \end{cases}. \quad (12.31)$$

Многомерное нормальное распределение

Тогда функция распределения

$$p(\mathbf{X}) = \frac{1}{(2\pi)^{k/2} |\mathbf{V}|^{1/2}} \exp\left[-\frac{1}{2} (\mathbf{X} - \boldsymbol{\mu})^T \mathbf{V}^{-1} (\mathbf{X} - \boldsymbol{\mu})\right] \quad (12.32)$$

где $|\mathbf{V}|$ -детерминант ковариационной матрицы \mathbf{V} определяет k -мерную плотность нормального распределения.

Наиболее наглядным вариантом является двумерное нормальное распределение, которое определяется параметрами $\mu_1, \mu_2, \sigma_1, \sigma_2$ и ρ_{12} :

$$p_N(x, y) = \frac{1}{2\pi\sqrt{(1-\rho_{12}^2)\sigma_1^2\sigma_2^2}} \times \exp\left\{\frac{\sigma_2^2(x-\mu_1)^2 - 2\rho_{12}\sigma_2\sigma_1(x-\mu_1)(y-\mu_2) + \sigma_1^2(y-\mu_2)^2}{2(\rho_{12}^2-1)\sigma_1^2\sigma_2^2}\right\} \quad (12.33)$$

Если $\mu_1 = \mu_2 = 0$ и $\sigma_1 = \sigma_2 = 1$, а $\rho_{12} = \rho$, то из (12.33) получим стандартизованное двумерное нормальное распределение:

$$p_N(x, y) = \frac{1}{2\pi\sqrt{1-\rho^2}} \exp\left\{\frac{x^2 - 2\rho xy + y^2}{2(\rho^2 - 1)}\right\}. \quad (12.34)$$

График стандартизованного двумерного нормального распределения с $\rho = 0$ (отсутствуют корреляции между x и y), представлен на рисунке 129.

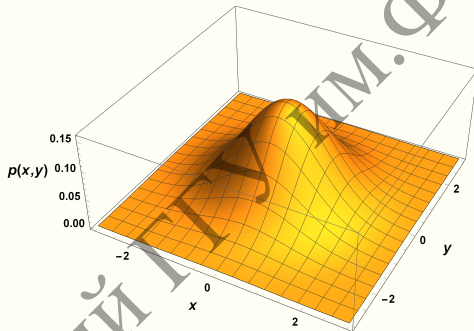


Рисунок 10: Двумерное нормальное распределение с $\mu_1 = \mu_2 = 0$ и $\sigma_1 = \sigma_2 = 1$, а $\rho = 0$

Естественно, существуют и другие теоретические функции распределения многомерных случайных величин.

13. Корреляционный и регрессионный анализ

Статистические связи между случайными величинами можно изучать методами корреляционного и регрессионного анализа. Основной целью регрессионного анализа является установление формы и изучение зависимости между переменными.

Целью корреляционного анализа является определить степень взаимосвязи между исследуемыми случайными величинами. В случае, если исследуется связь двух переменных, корреляционный анализ называют парным; если число переменных более двух — множественным.

Статистическая зависимость между переменными, при которой каждому значению одной или нескольких случайных величин соответствует определенное среднее значение другой, называется **корреляционной**.

Корреляционный и регрессионный анализ решает две основные задачи:

1. Определение уравнения (или системы уравнений) связи, т.е. в установлении математической формы, в которой выражается данная связь. Это очень важно, так как от правильного выбора формы связи зависит конечный результат изучения взаимосвязи между признаками.
2. Вычисление степени взаимосвязи, т.е. меры связи между признаками с целью установить степень влияния данных случайных величин (факторов) на конечную случайную величину (результат). Эта задача решается путем определения параметров корреляционного уравнения различными математическими методами.

Далее проводятся оценка и анализ полученных результатов при помощи различных показателей регрессионно-корреляционного метода (коэффициентов детерминации, линейной и множественной корреляции и т.д.), а также проверка существенности связи между изучаемыми признаками.

В результате использование корреляционного и регрессионного анализа можно решить такие задачи, как

- Взаимосвязь. Есть ли взаимосвязь между случайными величинами, и насколько велико их влияние на конечную величину.
- Прогнозирование. Если известно поведение одного параметра, то можно предсказать поведение другого параметра, коррелирующего с первым в области, где нет экспериментальной информации.
- Обнаружение неизвестных причинных связей.

Более часто используемым показателем степени тесноты корреляционной связи является линейный коэффициент корреляции. При расчете этого показателя учитываются не только отклонения индивидуальных значений случайной величины от среднего значения, но и сама величина этих отклонений.

В самом общем случае, при большом числе наблюдений одно и то же значение случайной величины X может встретиться n_x раз, одно и то же значение случайной величины Y может встретиться n_y раз, а одна и та же пара чисел $\{x, y\}$ может наблюдаться n_{xy} раз. Поэтому данные наблюдений группируют, т.е. подсчитывают частоты n_x , n_y и n_{xy} . Все сгруппированные данные записывают в виде таблицы, которую называют корреляционной.

Таблица 1: Пример корреляционной таблицы

$y \backslash x$	1	2	3	4	5	n_y
0	–	–	–	6	4	10
1	–	–	1	4	6	11
2	–	5	9	5	–	19
3	3	7	–	–	–	10
n_x	3	12	10	15	10	50

В первой строке указаны наблюдаемые значения $\{1, 2, 3, 4, 5\}$ случайной величины X , а в первом столбце таблицы – наблюдаемые значения $\{0, 1, 2, 3\}$ случайной величины Y .

На пересечении строк и столбцов находятся частоты n_{xy} наблюдаемых пар значений случайных величин X и Y . Например, частота 7 указывает, что пара чисел $\{2, 3\}$ наблюдалась 7 раз. Все частоты помещены в прямоугольнике.

В последнем столбце записаны суммы частот строк. В последней строке записаны суммы частот столбцов. Общее число наблюдений $n = 50$.

Если имеем k значений случайной величины X и m значений случайной величины Y в корреляционной таблице с числом попаданий n_i , n_j и n_{ij} для x_i , y_j и пар $\{x_i, y_j\}$ соответственно, то оценку коэффициент корреляции можно в данном случае можно найти из выражения:

$$r_{xy} = \frac{n \sum_{i=1}^k \sum_{j=1}^m x_i y_j n_{ij} - \left(\sum_{i=1}^k x_i n_i \right) \left(\sum_{j=1}^m y_j n_j \right)}{\sqrt{n \sum_{i=1}^k x_i^2 n_i - \left(\sum_{i=1}^k x_i n_i \right)^2} \sqrt{n \sum_{j=1}^m y_j^2 n_j - \left(\sum_{j=1}^m y_j n_j \right)^2}}. \quad (13.1)$$

В упрощенном варианте, когда имеется n значений пар $\{x_i, y_i\}$ двумерной случайной величины, выражение коэффициент корреляции (13.1) упрощается и приобретает вид:

$$r_{xy} = \frac{n \sum_{i=1}^n x_i y_i - \left(\sum_{i=1}^n x_i \right) \left(\sum_{j=1}^n y_j \right)}{\sqrt{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2} \sqrt{n \sum_{j=1}^n y_j^2 - \left(\sum_{j=1}^n y_j \right)^2}}. \quad (13.2)$$

Последний вариант часто встречается в физических исследованиях.

Иногда одному значению x_i соответствует m значений $\{\tilde{y}_{i1}, \tilde{y}_{i2}, \dots, \tilde{y}_{ik}\}$ величины Y . Тогда, после нахождения среднего значения и среднеквадратичного отклонения по известным формулам:

$$y_i = \frac{1}{m} \sum_{j=1}^m \tilde{y}_{ij},$$

$$s_{y_i} = \sqrt{\frac{1}{m-1} \sum_{j=1}^m (y_i - \tilde{y}_{ij})^2}, \quad (13.3)$$

наборам пар $\{x_i, y_i\}$ с дополнительным набором ошибок $\{s_{y_1}, \dots, s_{y_n}\}$.

Таблица Чеглока

Для определения степени связи используют обычно таблицу Чеглока (смотри таблицу 2).

Таблица 2: Пример корреляционной таблицы

Значение $ r_{xy} $	Степень связи
0,1 – 0,3	слабая
0,3 – 0,5	умеренная
0,5 – 0,7	заметная
0,7 – 0,9	высокая
0,9 – 0,99	весьма высокая связь

14. Линейный регрессионный анализ

Корреляционный анализ позволяет установить степень взаимосвязи двух и более случайных величин. Однако наряду с этим желательно иметь модель этой связи, которая дала бы возможность предсказывать значения одной случайной величины по конкретным значениям другой. Методы решения подобных задач составляют раздела математической статистики “регрессионный анализ”.

Для упрощения изложения рассмотрим случай двух случайных величин x и y . Линейная связь между двумя случайными величинами означает, что прогноз значения y по данному значению x имеет вид:

$$\hat{y} = A + Bx \quad (14.1)$$

Если данные связаны идеальной линейной зависимостью $|r_{xy}| = 1$, то предсказанное значение \tilde{y}_i будет соответствовать наблюдаемому значению y_i при любом данном x_i . На практике идеальная зависимость отсутствует (за счет случайных разбросов данных, за счет нелинейных эффектов).

Однако, если предположить наличие линейной связи, то можно подобрать A и B , которые дадут возможность предсказывать ожидаемое значение y_i для любого данного x_i (при этом предсказанное \tilde{y}_i не совпадает с наблюдаемыми y_i , однако оно будет равно среднему значению всех таких наблюдаемых значений).

Наиболее общепринятая процедура определения коэффициентов A и B состоит в таком выборе, который минимизирует сумму квадратов отклонений наблюдаемых значений от предсказанного значения y_i . Этот метод называют методом наименьших квадратов (МНК). Он разработан в 1795-1805 гг. Лежандром и Гауссом.

14. 1. Метод наименьших квадратов

Рассчитывается отклонение $y_i - \hat{y}_i = \Delta_i$, а затем находятся коэффициенты так, что

$$Q = \sum_{i=1}^n \Delta_i^2 \rightarrow \min \quad (14.2)$$

Для этого надо взять частные производные функции Q (14.2) по параметрам A и B и приравнять их к нулю.

В итоге имеем, что

$$\begin{cases} \frac{\partial Q}{\partial A} = 0, \\ \frac{\partial Q}{\partial B} = 0 \end{cases} \quad (14.3)$$

Из этой системы получаем оценки пар A и B .

Проделав цепочку преобразований

$$\frac{\partial Q}{\partial A} = -2 \sum_{i=1}^n (y_i - A - Bx_i) = 0. \quad (14.4)$$

$$An + b \sum_{i=1}^n x_i = \sum_{i=1}^n y_i,$$

$$\frac{\partial Q}{\partial B} = -2 \sum_{i=1}^n (y_i - A - Bx_i) x_i = 0,$$

$$- \sum_{i=1}^n y_i x_i + A \sum_{i=1}^n x_i + B \sum_{i=1}^n x_i^2 = 0,$$

$$An + B \left(\sum_{i=1}^n x_i \right) = \left(\sum_{i=1}^n y_i \right),$$

$$\left(\sum_{i=1}^n x_i \right) A + B \left(\sum_{i=1}^n x_i^2 \right) = \left(\sum_{i=1}^n x_i y_i \right)$$

В итоге имеем:

$$\left\{ \begin{array}{l} A = \frac{\sum_{i=1}^n x_i^2 \sum_{i=1}^n y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2}, \\ B = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2} \end{array} \right. \quad (14.5)$$

Эти оценки можно использовать для модели $\hat{y} = A + Bx$, которую называют прямой линейной регрессии y на x .

14. 2. Свойства метода наименьших квадратов

Из первого уравнения имеем:

$$A + B \sum_{i=1}^n \frac{x_i}{n} = \sum_{i=1}^n \frac{y_i}{n} = A + B\bar{x} = \bar{y}, \quad (14.6)$$

т.е. кривая регрессии проходит через центр тяжести экспериментальных точек.

Если теперь до проведения МНК все исходные данные отцентрировать, т.е. \bar{x} и \bar{y} перенести в начало координат $\bar{x}=\bar{y}=0$, то первое уравнение превратится в тождество, т.е. система уравнений сокращается, что особенно важно, когда имеет место ситуация с большим числом коэффициентов.

Уравнение прямой можно записать:

$$\hat{y} = A + Bx = \bar{y} + B(x - \bar{x}) . \quad (14.7)$$

Вторая особенность МНК состоит в том, что полученные этим методом оценки необратимы, т.е. если имеется модель регрессии y на x : $y = A + Bx$, то регрессию x на y нельзя обратить:

$$x = \frac{(y - A)}{B} = \frac{1}{b}y - A . \quad (14.8)$$

Коэффициент наклона при y обозначим как B' . Тогда имеем, что

$$x = B'y - A. \quad (14.9)$$

Можно найти по аналогии с вышеизложенным, что

$$B' = \frac{\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}}{\sum_{i=1}^n y_i^2 - n\bar{y}^2}. \quad (14.10)$$

Коэффициент (14.10) связан с B через коэффициент корреляции r_{xy} соотношением

$$\sqrt{B'B} = r_{xy}. \quad (14.11)$$

14. 3. Точность оценок A и B Ошибки A и B

Пусть $E\{A\} = a$, где a – “истинное” значение. Тогда доверительные интервалы для коэффициентов линейной регрессии A и B запишутся в виде:

$$a - A = \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)^{\frac{1}{2}} S_{y/x} t_{n-2, \alpha/2},$$

$$b - B = \left(\frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)^{\frac{1}{2}} S_{y/x} t_{n-2, \alpha/2}, \quad (14.12)$$

где

$$S_{y/x} = \left[\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2} \right]^{\frac{1}{2}} = \left[\left(\frac{n-1}{n-2} \right) S_y^2 (1 - r_{xy}^2) \right]^{\frac{1}{2}}, \quad (14.13)$$

а $t_{n-2, \alpha/2}$ - квантили распределения Стьюдента.

14. 4. Редукция некоторых задач к линейному анализу

Как сводить некоторые задачи регрессионного анализа к линейному регрессионному анализу.

Например, модель типа

$$\frac{1}{y} = a_1 + b_1 x$$

путем замены переменных $\frac{1}{y} = y_1$ имеем $y_1 = a_1 + b_1 x$.

Аналогично

$$y = a_1 + \frac{b_1}{x} \rightarrow x,$$

$$\frac{y}{x} = a_1 + b_1 x,$$

$$y = ax^b \lg y = \lg a + b \lg x,$$

$$\lg y \rightarrow y \lg x \rightarrow x,$$

$$y = ab^x \lg y = \lg a + x \lg b.$$

Такие преобразования, строго говоря, допустимы, если исходные величины измерены точно. Например: $y = Ax^B$ измерен с Δy : $y = Ax^B + \Delta y$.

Тогда линеаризованная форма будет отличаться от исходного за счет неясности преобразования Δy . Однако в первичном анализе Δy можно пренебречь, если это необоснованно, то вводим переменную y' путем

$$y' = y - \Delta y = Ax^B .$$

Метод МНК достаточно чувствителен к неоднородностям статистики. Наличие неоднородной статистики приводит иногда к абсурдным результатам, поэтому окончательное решение об МНК должно приниматься по однородной, очищенной статистике.

Пусть одному значению x_i из выборки объема n соответствует m значений $\{\tilde{y}_{i1}, \tilde{y}_{i2}, \dots, \tilde{y}_{ik}\}$ величины Y . Тогда, после нахождения среднего значения и среднеквадратичного отклонения, получим набор пар $\{x_i, y_i\}$ с дополнительным набором ошибок $\{\sigma_{y_1}, \dots, \sigma_{y_n}\}$. Такая ситуация возникает, например, при построении гистограмм.

Цель регрессионного анализа найти уравнение взаимосвязи случайных величин Y и X (в данном случае, парная регрессионная модель) вида: $y = \phi(x; \theta)$. В силу воздействия неучтенных случайных факторов отдельные значения y_i будут в большей или меньшей мере отклоняться от функции регрессии $\phi(x; \theta)$. В этом случае уравнение взаимосвязи двух переменных (парная регрессионная модель) может быть представлено в виде:

$$y = \phi(x; \theta) + \epsilon, \quad (15.1)$$

где $\epsilon = \{\epsilon_1, \dots, \epsilon_n\}$ случайная величина, характеризующая отклонение от функции регрессии. Эту переменную часто называют возмущением; величина $\theta = \{\theta_1, \theta_2, \dots, \theta_k\}$ определяет набор параметров регрессионной модели, которые необходимо определить.

Основные положения регрессионного анализа:

Основные положения регрессионного анализа

1. Одним из основных требований регрессионного анализа является, условие равенства нулю математического ожидания “возмущения” ϵ :

$$E \{ \epsilon_i \} = 0 . \quad (15.2)$$

2. Также предполагают, что случайные величины ϵ_i подчиняются одномерному нормальному распределению, если все y_i и y_j не коррелируют с другом или многомерному нормальному распределению, если такая корреляция имеет место.
3. Дисперсия “возмущения” ϵ_i (или зависимой переменной y для любого i) задается соотношением:

$$D \{ \epsilon_i \} = \sigma_i^2 . \quad (15.3)$$

Наиболее простой вариант (15.3) состоит в требовании того, чтобы все дисперсии отклонений были одинаковые.

В случае постоянства дисперсии и отсутствия корреляции, оценки параметров регрессии полученные, например, с помощью метода наименьших квадратов, обладают важными свойствами, а именно:

- несмещенность;
- состоятельность;
- эффективность.

В случае нелинейной зависимости, свойства постоянства дисперсии и отсутствия корреляции могут не выполняться, тогда полученные оценки параметров регрессии не будут обладать указанными характеристиками. Или связь между переменными x и y линейна, но на исследуемый показатель воздействует фактор, не включенный в модель.

Для решения этой задачи в случае линейной регрессии использовался метод наименьших квадратов. Рассмотрим некоторые методы поиска оптимальной оценки $\hat{\theta}$ в случае нелинейной регрессии. Естественно, эти методы могут быть использованы и в линейном варианте функциональной зависимости.

15. 1. Метод максимального правдоподобия

Пусть имеется выборка из n независимых результатов наблюдений $x = \{x_1, x_2, \dots, x_n\}$ с плотностью вероятности $p(x, \theta)$ (θ - набор параметров).

Определение 15.1

Совместная функция вероятности

$$L(x_1, x_2, \dots, x_n, \theta) = p(x_1, \theta) p(x_2, \theta) \cdots p(x_n, \theta) = \prod_{i=1}^n p(x_i, \theta) \quad (15.4)$$

называют **функцией максимального правдоподобия**.

Совместная функция вероятности $L(x_1, x_2, \dots, x_n, \theta)$ может быть представлена в виде (15.4) в силу независимости результатов каждого измерения. Две функции правдоподобия являются равными, если одна есть произведение второй на некоторую скалярную величину.

$L(x_1, x_2, \dots, x_n, \theta)$ рассматривают как функцию θ при фиксированных, полученных в измерениях x . Чем больше значение $L(x_1, x_2, \dots, x_n, \theta)$, тем более вероятна или правдоподобна, выборка значений $\{x_1, x_2, \dots, x_n\}$ при заданном значении θ . Поэтому $L(x_1, x_2, \dots, x_n, \theta)$ и называют функцией правдоподобия.

Метод максимального правдоподобия или метод наибольшего правдоподобия (ММП, ML, MLE — англ. maximum likelihood estimation) в математической статистике — это метод оценивания неизвестного параметра путём максимизации функции правдоподобия. Основан на предположении о том, что вся информация о статистической выборке содержится в функции правдоподобия.

Метод максимального правдоподобия состоит в поиске максимума функции $L(x_1, x_2, \dots, x_n, \theta)$, при этом $\{x_1, x_2, \dots, x_n\}$ считаются постоянными, а параметр θ - переменным. Далее находят такое значение θ , при котором функция $L(x_1, x_2, \dots, x_n, \theta)$ становится наиболее правдоподобной, т.е. **принимает максимальное значение.**

Для эффективного использования ММП требуются достаточно большие объемы выборок, точное знание анализируемого закона распределения, достаточно устойчивые распределения и не очень большое число неизвестных параметров.

Доказано, что вторая производная от функции правдоподобия $L(x_1, x_2, \dots, x_n, \theta)$ меньше нуля и, таким образом, равенство нулю первой производной дает действительно максимальное значение $L(x_1, x_2, \dots, x_n, \theta)$. Максимум определяют по стандартной методике, однако вместо самой функции для удобства вычислений берут логарифм функции и ищут его максимум. Поскольку максимумы функции правдоподобия и логарифмической функции совпадают, то для удобства вычислений берут логарифм функции и ищут его максимум.

Введем обозначение

$$\ell(x_1, x_2, \dots, x_n, \theta) = \ln L(x_1, x_2, \dots, x_n, \theta) = \sum_{i=1}^n \ln p(x_i, \theta). \quad (15.5)$$

Система уравнений для получения оценок параметров $\theta = \{\theta_1, \theta_2, \dots, \theta_k\}$ можно записать в виде

$$\left\{ \begin{array}{l} \frac{\partial \ell(x_1, x_2, \dots, x_n, \theta)}{\partial \theta_1} = 0, \\ \frac{\partial \ell(x_1, x_2, \dots, x_n, \theta)}{\partial \theta_2} = 0, \\ \dots \\ \frac{\partial \ell(x_1, x_2, \dots, x_n, \theta)}{\partial \theta_k} = 0, \end{array} \right. \quad (15.6)$$

Решение системы уравнений (15.6) дает набор оптимальных значений параметров $\tilde{\theta} = \{\tilde{\theta}_1, \tilde{\theta}_2, \dots, \tilde{\theta}_k\}$.

Наибольшие трудности возникают не при оценке параметров с помощью решения системы (15.6), а при выявлении погрешностей, с которыми эти оценки сделаны. Возможные (неявные и естественные) связи между параметрами $\theta = \{\theta_1, \theta_2, \dots, \theta_k\}$ приводят к необходимости учитывать корреляции (ковариации) между оценками различных параметров $\tilde{\theta} = \{\tilde{\theta}_1, \tilde{\theta}_2, \dots, \tilde{\theta}_k\}$.

Введем обозначения для сокращения записи соотношений:

$$\begin{aligned} \ell(x_1, x_2, \dots, x_n, \theta) &= \ell(\theta) , \\ L(x_1, x_2, \dots, x_n, \tilde{\theta}) &= L_{max} . \end{aligned} \quad (15.7)$$

Для оценки погрешностей параметров $\tilde{\theta} = \{\tilde{\theta}_1, \tilde{\theta}_2, \dots, \tilde{\theta}_k\}$ разложим логарифмическую функцию правдоподобия в ряд Тейлора в окрестностях точек $\{\tilde{\theta}_1, \tilde{\theta}_2, \dots, \tilde{\theta}_k\}$:

$$\begin{aligned} \ell(\theta) = & \ell(\tilde{\theta}) + \sum_{i=1}^k \left(\frac{\partial \ell}{\partial \theta_i} \right)_{\tilde{\theta}} (\theta_i - \tilde{\theta}_i) + \\ & + \frac{1}{2} \sum_{p=1}^k \sum_{m=1}^k \left(\frac{\partial^2 \ell}{\partial \theta_p \partial \theta_m} \right)_{\tilde{\theta}} (\theta_i - \tilde{\theta}_i) (\theta_m - \tilde{\theta}_m) + \dots \quad (15.8) \end{aligned}$$

Величины $(\theta_i - \tilde{\theta}_i)$ трактуются как дисперсии оценок максимального правдоподобия $\tilde{\theta}_i$.

Из определения оценок $\tilde{\theta}$ следует, что второй член в разложении равен нулю для всех k . Пренебрегая членами разложения выше квадратичных, логарифмическую функцию правдоподобия можно записать в матричном виде:

$$\ell(\theta) = \ln L_{max} + \frac{1}{2} (\theta - \tilde{\theta}^T) \mathbf{A} (\theta - \tilde{\theta}), \quad (15.9)$$

где

$$\mathbf{A} = \begin{pmatrix} \frac{\partial^2 \ell}{\partial \theta_1^2} & \frac{\partial^2 \ell}{\partial \theta_1 \theta_2} & \cdots & \frac{\partial^2 \ell}{\partial \theta_1 \theta_k} \\ \frac{\partial^2 \ell}{\partial \theta_1 \theta_2} & \frac{\partial^2 \ell}{\partial \theta_2^2} & \cdots & \frac{\partial^2 \ell}{\partial \theta_2 \theta_k} \\ \vdots & \cdots & \cdots & \cdots \\ \frac{\partial^2 \ell}{\partial \theta_1 \theta_k} & \frac{\partial^2 \ell}{\partial \theta_2 \theta_k} & \cdots & \frac{\partial^2 \ell}{\partial \theta_k^2} \end{pmatrix}_{\theta = \tilde{\theta}}. \quad (15.10)$$

В асимптотическом пределе $n \rightarrow \infty$, элементы матрицы практически не зависят от конкретной выборки $\{x_1, x_2, \dots, x_n\}$, и их можно заменить математическими ожиданиями:

$$\mathbf{B} = \mathbf{E} \{ \mathbf{A} \} = \begin{pmatrix} \mathbf{E} \left\{ \frac{\partial^2 \ell}{\partial \theta_1^2} \right\} & \mathbf{E} \left\{ \frac{\partial^2 \ell}{\partial \theta_1 \theta_2} \right\} & \dots & \mathbf{E} \left\{ \frac{\partial^2 \ell}{\partial \theta_1 \theta_k} \right\} \\ \mathbf{E} \left\{ \frac{\partial^2 \ell}{\partial \theta_1 \theta_2} \right\} & \mathbf{E} \left\{ \frac{\partial^2 \ell}{\partial \theta_2^2} \right\} & \dots & \mathbf{E} \left\{ \frac{\partial^2 \ell}{\partial \theta_2 \theta_k} \right\} \\ \dots & \dots & \dots & \dots \\ \mathbf{E} \left\{ \frac{\partial^2 \ell}{\partial \theta_1 \theta_k} \right\} & \mathbf{E} \left\{ \frac{\partial^2 \ell}{\partial \theta_2 \theta_k} \right\} & \dots & \mathbf{E} \left\{ \frac{\partial^2 \ell}{\partial \theta_k^2} \right\} \end{pmatrix}_{\theta = \tilde{\theta}} \quad (15.11)$$

После потенцирования (15.9) следует, что функция правдоподобия имеет вид

$$L(\boldsymbol{\theta}) = L_{max} \exp \left\{ \frac{1}{2} (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}^T) \mathbf{B} (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}) \right\}. \quad (15.12)$$

Сопоставление с известными теоретическими распределениями приводит к выводу, что (15.12) представляет собой k -мерное нормальное распределение со средними оценками $\{\tilde{\theta}_1, \tilde{\theta}_2, \dots, \tilde{\theta}_k\}$ и ковариационной матрицей

$$\mathbf{V} = -\mathbf{B}^{-1}. \quad (15.13)$$

Диагональные элементы ковариационной матрицы (15.13) - являются дисперсиями $c_{11} = \sigma_{\tilde{\theta}_1}, \dots, c_{kk} = \sigma_{\tilde{\theta}_k}$ оценок максимального правдоподобия $\{\tilde{\theta}_1, \tilde{\theta}_2, \dots, \tilde{\theta}_k\}$, а недиагональные элементы представляют собой ковариации между всевозможными парами оценок $c_{ij} = \text{cov}(\tilde{\theta}_i, \tilde{\theta}_j)$.

Коэффициент корреляции между оценками определяется формулой

$$\rho(\tilde{\theta}_i, \tilde{\theta}_j) = \frac{\text{cov}(\tilde{\theta}_i, \tilde{\theta}_j)}{\sigma_{\tilde{\theta}_i} \sigma_{\tilde{\theta}_j}}, \quad i \neq j. \quad (15.14)$$

Поскольку $L(\boldsymbol{\theta})$ в окрестности $\tilde{\boldsymbol{\theta}}$ представляет собой k -мерное нормальное распределение, то квадратичная форма

$$Q(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}) = (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})^T \mathbf{V}^{-1} (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}) \quad (15.15)$$

распределена по закону χ^2 с r степенями свободы.

Поэтому можно определить вероятностное утверждение:

$$\text{Prob}[Q(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}) \leq Q_\alpha] = \alpha, \quad (15.16)$$

где квантиль распределения χ^2 с r степенями свободы Q_α уровня α задается соотношением

$$\int_0^{Q_\alpha} \frac{2^{-r/2}}{\Gamma(r/2)} \exp[-x/2] x^{r/2-1} dx = 1 - \alpha. \quad (15.17)$$

Гиперэллипсоид, определяемый уравнением

$$Q(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}) = Q_\alpha \quad (15.18)$$

будет задавать доверительную область в пространстве переменных $\tilde{\boldsymbol{\theta}}$ с вероятностным содержанием (доверительной вероятностью) $P = 1 - \alpha$.

Для приближенного построения гиперэллипсоида (15.18), который определяет возможные значения оптимальных параметров с доверительной вероятностью P , можно использовать уравнение вида (смотри (15.9)):

$$\ln L(\boldsymbol{\theta}) \approx \ln L_{max} - \Delta(\ln L(\boldsymbol{\theta})) \quad (15.19)$$

Таким образом, как следует из (15.17) и (15.18), значение $\Delta(\ln L(\boldsymbol{\theta}))$ определяется задаваемым уровнем достоверности (С.Л.) и числом параметров, входящих в набор $\boldsymbol{\theta}$.

С помощью (15.19) для доверительной вероятности $P = 95\%$ можно найти для числа параметров 1, 2 и 3, числовое значение $\Delta(\ln L(\boldsymbol{\theta})) \approx 3,84, 5,99$ и $7,82$ соответственно. Аналогичные значения для стандартного отклонения с $P = 68,27\%$ приводят к значениям $\Delta(\ln L(\boldsymbol{\theta})) \approx 1,00, 2,30$ и $3,53$.

В этом случае, можно построить контурные графики функции

$$\ln L(\boldsymbol{\theta}) - \ln L_{max} = \Delta(\ln L(\boldsymbol{\theta})) \quad (15.20)$$

для различных пар параметров и найти среднеквадратичные отклонения $\sigma_{\tilde{\theta}_1}, \dots, \sigma_{\tilde{\theta}_k}$ оценок максимального правдоподобия $\{\tilde{\theta}_1, \tilde{\theta}_2, \dots, \tilde{\theta}_k\}$.

15. 2. Метод наименьших квадратов

Метод наименьших квадратов также используется для получения “наилучших” (оптимальных) значений модельных параметров при описании экспериментальных значений.

Метод наименьших квадратов совпадает с методом максимального правдоподобия в следующем частном случае. Рассмотрим набор из n независимых измерений y_i в известных точках x_i .

Из y_i из $y = \{y_1 \pm \sigma_1, \dots, y_n \pm \sigma_k\}$ подчиняются подчиняются нормальному распределению с математическим ожиданием $\phi(x_i, \theta)$ и известной дисперсией σ_i^2 . Таким образом предполагается, что описание значений y_i с помощью модельной регрессионной кривой $\phi(x_i, \theta)$ должно быть оптимальным. Цель состоит в том, чтобы построить оценки для неизвестных параметров θ .

Для получения оптимальных значений $\tilde{\theta}$ набора модельных параметров θ используют функцию

$$\chi^2(\theta) = -2 \ln L(\theta) + \text{const} = \sum_{i=1}^n \left[\frac{y_i - \phi(x_i, \theta)}{\sigma_i} \right]^2. \quad (15.21)$$

Соотношение (15.21) можно получить из определения функции максимального правдоподобия (15.1), используя плотность нормального распределения

$$p(y_i) = \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp \left[-\frac{(y_i - \phi(x_i, \theta))^2}{2\sigma_i^2} \right]. \quad (15.22)$$

Исходя из требования минимального значения функции $\chi^2(\boldsymbol{\theta})$ при $\boldsymbol{\theta} = \tilde{\boldsymbol{\theta}}$, т. е.

$$\chi^2(\tilde{\boldsymbol{\theta}}) = \chi_{min}^2 \quad (15.23)$$

получаем систему уравнений

$$\left\{ \begin{array}{l} \frac{\partial \chi^2(\boldsymbol{\theta})}{\partial \theta_1} = 0, \\ \frac{\partial \chi^2(\boldsymbol{\theta})}{\partial \theta_2} = 0, \\ \dots \\ \frac{\partial \chi^2(\boldsymbol{\theta})}{\partial \theta_k} = 0, \end{array} \right. \quad (15.24)$$

для определения “оптимальных” значений $\tilde{\boldsymbol{\theta}}$. Процедура аналогична линейному регрессионному анализу, с той лишь разницей, что поиск ‘оптимальных’ значений $\tilde{\boldsymbol{\theta}}$, как правило, находится численно, а не аналитически.

Значение χ_{min}^2 является мерой уровня согласия между экспериментальными данными и подогаданной кривой (fit):

$$\chi_{min}^2 = \sum_{i=1}^n \left[\frac{y_i - \phi(x_i, \tilde{\theta})}{\sigma_i} \right]^2. \quad (15.25)$$

Поэтому χ_{min}^2 можно использовать как статистику соответствия предполагаемую функциональную форму $\phi(x_i, \tilde{\theta})$. Известно, что (15.25) имеет распределение χ^2 с $n_d = n - k$ степенями свободы.

В этом случае, в соответствии с разделом проверки гипотез, если $\chi_{min}^2/n_d \leq 1$, то совпадение будет “хорошим”. Или можно найти вероятность соответствия модели $\phi(x_i, \tilde{\theta})$ экспериментальным данным P (p -value) из соотношения:

$$P = \frac{1}{2^{n_d} \int_{\chi_{min}^2}^{\infty} \Gamma\{n_d/2\}} t^{n_d/2-1} \exp(-t/2) . \quad (15.26)$$

Следующим шагом решения данной задачи состоит в нахождении доверительных интервалов (ошибок $\sigma_{\tilde{\theta}_i}$) для оптимальных значений $\tilde{\theta} = \{\tilde{\theta}_1, \dots, \tilde{\theta}_k\}$. Эта задача решается так же как и для метода максимального правдоподобия, если предположить, что оценки $\tilde{\theta}_i, i = 1, \dots, k$ подчиняются k -мерному нормальному распределению (15.12).

Отличием, состоит в том, что приближенного построения гиперэллипсоида (15.18), который определяет возможные значения оптимальных параметров с вероятностным содержанием С.Л., используется уравнение вида:

$$\chi^2(\boldsymbol{\theta}) = \chi_{min}^2 + \Delta\chi_{crit}^2. \quad (15.27)$$

Значение $\Delta\chi_{crit}^2$ определяется задаваемым уровнем достоверности p и числом параметров, входящих в набор $\boldsymbol{\theta}$ и совпадает с $\Delta(\ln L(\boldsymbol{\theta}))$.

Замечание.

Если y_i не являются независимыми, и имеют ковариационную матрицу $\mathbf{V}_{ij} = \text{cov}(y_i, y_j)$, то наилучшие оценки параметров $\tilde{\boldsymbol{\theta}} = \{\tilde{\theta}_1, \dots, \tilde{\theta}_k\}$ путем поиска минимума функционала

$$\chi^2(\boldsymbol{\theta}) = \sum_{i=1}^n \sum_{j=1}^n (y_i - \phi(x_i, \boldsymbol{\theta})) (\mathbf{V})_{ij}^{-1} (y_j - \phi(x_j, \boldsymbol{\theta})). \quad (15.28)$$

16. Сравнение моделей

При поиске зависимости между измеряемыми переменными по выборочным данным наиболее важной задачей является поиск модели регрессии. Построение эмпирического уравнения регрессии - начальный этап анализа. При этом, на первом этапе не всегда удается получить “оптимальную” модель регрессии. Тогда для описания зависимости между случайными величинами может использоваться несколько моделей (уравнений регрессии). Поэтому возникает задача об оценке “качества” полученных моделей.

Качество модели оценивается по следующим направлениям:

- 1 Содержательная оценка качества модели.
- 2 Статистическая оценка качества модели.

Содержательный анализ подразумевает анализ физического (экономического и т.д.) смысла модели. Результат данной работы должен быть ответы на вопросы вида:

- ✓ *Являются ли те факторы, которые используются в модели, значимыми для описания данного физического явления или процесса?*
- ✓ *Какой физический смысл параметров регрессионного уравнения (модели)?*

Статистическая оценка качества модели включает следующие этапы:

- ✓ Проверка статистической значимости параметров уравнения регрессии при помощи различных критериев.
- ✓ Проверку общего качества уравнения регрессии.
- ✓ Проверку свойств данных, выполнение которых предполагалось при оценивании уравнения.

16. 1. Проверка статистической значимости параметров уравнения регрессии

После того, как найдено уравнение регрессии (например, линейной регрессии (14.1)), проводится оценка значимости как уравнения в целом, так и отдельных его параметров.

F-критерий Фишера

Оценка значимости уравнения регрессии в целом дается с помощью *F*-критерия Фишера. При этом выдвигается нулевая гипотеза о том, что коэффициент регрессии *B* равен нулю, т. е. $H_0 : B = 0$. Следовательно, фактор случайной величины *x* не оказывает влияния на результат *y*.

Перед расчетом критерия проводятся анализ дисперсии. Общая сумма квадратов отклонений (СКО) y от среднего значения \bar{y} раскладывается на две части - “объясненную” и “необъясненную”:

$$\underbrace{\sum_{i=1}^n (y_i - \bar{y}_i)^2}_{\text{Общая СКО}} = \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y}_i)^2}_{\text{Факторная СКО}} + \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i)^2}_{\text{Остаточная СКО}} . \quad (16.1)$$

(Объясненная) (Необъясненная)

Проиллюстрируем соотношение (16.1) на примере двух крайних вариантов. Если общая СКО в точности равна остаточной, то случайный фактор x не оказывает влияния на результат, вся дисперсия y обусловлена воздействием прочих факторов. При этом прямая для линейной регрессии параллельна оси Ox .

Когда общая СК0 равна факторной, то никакие другие (прочие) факторы не влияют на результат. Величина y связана с x с вероятностью 100%, и остаточная СК0 равна нулю.

Однако на практике в правой части (16.1) присутствуют оба слагаемых. Пригодность линии регрессии для прогноза зависит от того, какая часть общей вариации y приходится на объясненную вариацию. Если объясненная СК0 будет больше остаточной СК0, то уравнение регрессии статистически значимо и фактор x оказывает существенное воздействие на результат y . Это равносильно тому, что коэффициент корреляции будет приближаться к единице.

Для использования критерия Фишера для оценки значимости регрессии необходимо ввести понятие число степеней свободы для различных величин.

Определение 16.1

Число степеней свободы (df-degrees of freedom) - это минимально необходимое число значений зависимой переменной, которых достаточно для получения искомой характеристики выборки и которые могут свободно изменяться с учетом того, что для этой выборки известны все другие величины, используемые для расчета искомой характеристики. Или другими словами - это **число независимо варьируемых значений признака**.

Уравнение для определения F -статистики в случае многомерной регрессии имеет вид:

$$F_{\phi} = \left(\frac{\sum_{i=1}^n (\hat{y}_i - \bar{y}_i)^2}{m} \right) \left\{ 1 / \frac{\left(\sum_{i=1}^n (y_i - \hat{y}_i)^2 \right)}{(n - m - 1)} \right\}, \quad (16.2)$$

где n - объем выборки; m - число независимых переменных в факторной части СКО.

Поясним на примере линейной регрессии. Факторную СКО в этом случае можно выразить так:

$$\begin{aligned} \sum_{i=1}^n (\hat{y}_i - \bar{y}_i)^2 &= \sum_{i=1}^n ([A + Bx_i] - [A + B\bar{x}])^2 = \\ &= \sum_{i=1}^n (Bx_i - B\bar{x})^2 = B^2 \sum_{i=1}^n (x_i - \bar{x})^2 . \end{aligned} \quad (16.3)$$

СКО (16.3) зависит только от одного параметра B . Следовательно, факторная СКО в случае линейной регрессии имеет одну степень свободы, т. е. $m = 1$.

Можно и рассчитать число t и другим способом. Для того вспомним выражение (14.7) для регрессионной прямой

$$\hat{y} = A + Bx = \bar{y} + B(x - \bar{x}) . \quad (16.4)$$

Как видно из (16.4) прямая регрессии определяется только одним параметром.

Разделив каждую СКО на свое число степеней свободы, получим среднее квадратичные отклонения (или дисперсии на одну степень свободы):

$$\begin{aligned}
 D_{\text{общ.}} &= \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y}_i)^2, \\
 D_{\text{факт.}} &= \frac{1}{m} \sum_{i=1}^n (\hat{y}_i - \bar{y}_i)^2, \\
 D_{\text{ост.}} &= \frac{1}{n-m-1} \sum_{i=1}^n (y_i - \hat{y}_i)^2.
 \end{aligned} \tag{16.5}$$

С учетом (16.5) F -критерий (16.2) можно записать в виде:

$$F_{\phi} = \frac{D_{\text{факт.}}}{D_{\text{ост.}}} \tag{16.6}$$

В теории вероятности доказано, что для выборки из генеральной совокупности у которой отсутствует связь между зависимой y и независимой x (или x_1, x_2, \dots, x_k) переменной имеет распределение Фишера:

$$p_{FR}(x, f_1, f_2) = \frac{f_1^{f_1/2} f_2^{f_2/2} x^{f_1/2-1} (f_1 x + f_2)^{-\frac{1}{2}(f_1+f_2)}}{B\left(\frac{f_1}{2}, \frac{f_2}{2}\right)}, \quad x \geq 0,$$

$$F_R(x, f_1, f_2) = I_{\frac{f_1 x}{f_1 x + f_2}}\left(\frac{f_1}{2}, \frac{f_2}{2}\right). \quad (16.7)$$

где $f_{1,2}$ - параметры распределения; функция $I_n(a, b)$ - регуляризованная неполная бэ́та функция.

Благодаря этому для осуществления статистической проверки значимости уравнения регрессии формулируется нулевая гипотеза об отсутствии связи между переменными (все коэффициенты при переменных равны нулю) и выбирается уровень значимости α .

Уровень значимости α - это вероятность совершить ошибку первого рода, т. е. отвергнуть в результате проверки верную нулевую гипотезу. В рассматриваемом случае совершить ошибку первого рода означает признать по выборке наличие связи между переменными в генеральной совокупности, **когда на самом деле ее там нет.**

Вычисленное значение F_{Φ} признается достоверным (отличным от единицы), если оно больше табличного квантиля распределения Фишера F_{α}

$$\int_{F_{\alpha}}^{\infty} p_{F_R}(x, m, n - m - 1) dx = \alpha, \quad (16.8)$$

т.е.

$$F_{\Phi} \geq F_{\alpha}. \quad (16.9)$$

В этом случае нулевая гипотеза об отсутствии связи между переменными отклоняется и делается вывод о существенности превышения $D_{\text{факт}}$ над $D_{\text{ост.}}$. Или другими словами делается вывод о существенной статистической зависимости между зависимой y и независимыми величинами $x = x_1, \dots, x_k$.

Если $F_{\text{ф}} \leq F_{\alpha}$, то вероятность нулевой гипотезы выше заданного уровня α (например, $0,05 \div 0,10$), и эта гипотеза не может быть отклонена без серьезного риска сделать неправильный вывод о наличии связи между y и x . Уравнение регрессии считается статистически незначимым, гипотеза H_0 не отклоняется.

С помощью компьютерных вычислений, практически всегда можно найти α путем решения уравнения:

$$\int_{F_{\phi}}^{\infty} p_{FR}(x, m, n - m - 1) dx = 1 - F_R(F_{\phi}, m, n - m - 1) = \alpha. \quad (16.10)$$

В если α относительно велико ($\alpha > 0,1$), то о наличии корреляции говорит сложно (или с большими оговорками).

Как и в случае парной регрессии, статистическая значимость коэффициентов множественной линейной регрессии с m объясняющими переменными проверяется на основе t -критерия Стьюдента. Величина стандартной ошибки совместно с t -распределением Стьюдента при $n - m$ степенях свободы применяется для проверки существенности коэффициента регрессии и для расчета его доверительных интервалов. Для этого для каждого параметра регрессии a_i , полученного в результате вычислений рассчитывается величина

$$t_{\text{exp}} = \frac{a_i}{\sigma_{a_i}}, \quad (16.11)$$

т.е. величина параметров регрессии сравнивается с его стандартной ошибкой.

Значение (16.11) сравнивается с табличным значением при определенном уровне значимости α и числе степеней свободы. По сути проверяется нулевая гипотеза в виде $H_0: a_i = 0$. Если $t_{\text{exp}} > t_{n-m, \alpha/2}$, то гипотеза $H_0: a_i = 0$ должна быть отклонена, а статистическая связь y с x считается установленной. В случае $t_{\text{exp}} > t_{n-m, \alpha/2}$ нулевая гипотеза не может быть отклонена, и влияние x на y признается несущественным.

16. 2. Проверка общего качества уравнения регрессии

Оценить общее качество уравнения регрессии означает, установить соответствует ли математическая модель, выражающая зависимость между переменными экспериментальным данным и достаточно ли включенных в модель переменных, объясняющих поведение предсказанной величины (y). Оценить общее качества модели = оценить надежность модели = оценить достоверность уравнения регрессии.

16. 3. Коэффициент детерминации

Для оценки качества различных моделей используют коэффициент детерминации.

Определение 16.2

Коэффициент детерминации (R^2 — R -квадрат) — это доля дисперсии зависимой переменной, объясняемая рассматриваемой моделью зависимости, то есть объясняющими переменными.

Более точно — это единица минус доля необъяснённой дисперсии (дисперсии случайной ошибки модели, или условной по факторам дисперсии зависимой переменной) в дисперсии зависимой переменной. Его рассматривают как универсальную меру зависимости одной случайной величины от множества других. В частном случае линейной зависимости R^2 является квадратом так называемого множественного коэффициента корреляции между зависимой переменной и объясняющими переменными.

В частности, для модели парной линейной регрессии коэффициент детерминации равен квадрату обычного коэффициента корреляции между y и x .

Истинный коэффициент детерминации модели зависимости случайной величины y от факторов x определяется следующим образом:

$$R^2 = 1 - \frac{D\{y|x\}}{D\{y\}} = 1 - \frac{\sigma^2}{\sigma_y^2}, \quad (16.12)$$

где $D\{y|x\} = \sigma^2$ — условная (по факторам x) дисперсия зависимой переменной (дисперсия случайной ошибки модели).

В данном определении используются истинные параметры, характеризующие распределение случайных величин. Если использовать выборочную оценку значений соответствующих дисперсий, то получим формулу для выборочного коэффициента детерминации (который обычно и подразумевается под коэффициентом детерминации):

Коэффициент детерминации

$$R^2 = 1 - \frac{\hat{\sigma}^2}{\hat{\sigma}_y^2} = 1 - \frac{SS_{res}/n}{SS_{tot}/n} = 1 - \frac{SS_{res}}{SS_{tot}}, \quad (16.13)$$

где

$$SS_{res} = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (16.14)$$

сумма квадратов остатков регрессии, y_i , \hat{y}_i — фактические (“экспериментальные”) и расчётные значения объясняемой переменной, а величина

$$SS_{tot} = \sum_{i=1}^n (y_i - \bar{y})^2 = n\hat{\sigma}_y^2 \quad (16.15)$$

так называемая **общая сумма квадратов**, а

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i .$$

В случае линейной регрессии с константой $SS_{tot} = SS_{reg} + SS_{res}$, где

$$SS_{reg} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \quad (16.16)$$

объяснённая сумма квадратов, поэтому получаем более простое определение в этом случае — коэффициент детерминации — это доля объяснённой суммы квадратов в общей:

$$R^2 = \frac{SS_{reg}}{SS_{tot}}. \quad (16.17)$$

Необходимо подчеркнуть, что формула (16.17) справедлива только для модели с константой, в общем случае необходимо использовать предыдущую формулу.

Интерпретация

1. Коэффициент детерминации для модели с константой принимает значения от 0 до 1. Чем ближе значение коэффициента к 1, тем сильнее зависимость. При оценке регрессионных моделей это интерпретируется как соответствие модели данным. Для приемлемых моделей предполагается, что коэффициент детерминации должен быть хотя бы не меньше 50% (в этом случае коэффициент множественной корреляции превышает по модулю 70%). Модели с коэффициентом детерминации выше 80% можно признать достаточно хорошими (коэффициент корреляции превышает 90%). Значение коэффициента детерминации 1 означает функциональную зависимость между переменными.

Интерпретация

2. При отсутствии статистической связи между объясняемой переменной и факторами, статистика nR^2 для линейной регрессии имеет асимптотическое распределение $\chi^2(k-1)$, где $k-1$ — количество факторов модели. В случае линейной регрессии с нормально распределёнными случайными ошибками статистика $F = \frac{R^2/(k-1)}{(1-R^2)/(n-k)}$ имеет точное (для выборок любого объёма) распределение Фишера $F(k-1, n-k)$ (так называемый F-тест). Информация о распределении этих величин позволяет проверить статистическую значимость регрессионной модели исходя из значения коэффициента детерминации. Фактически в этих тестах проверяется гипотеза о равенстве истинного коэффициента детерминации нулю.
3. В общем случае коэффициент детерминации может быть и отрицательным, это говорит о крайней неадекватности модели: простое среднее приближает лучше.

16. 4. Недостаток R^2 и альтернативные показатели

Основная проблема применения (выборочного) R^2 заключается в том, что его значение увеличивается (не уменьшается) от добавления в модель новых переменных, даже если эти переменные никакого отношения к объясняемой переменной не имеют! Поэтому сравнение моделей с разным количеством факторов с помощью коэффициента детерминации, вообще говоря, некорректно. Для этих целей можно использовать альтернативные показатели.

Скорректированный (adjusted) R^2

Для того, чтобы была возможность сравнивать модели с разным числом факторов так, чтобы число факторов не влияло на статистику R^2 обычно используется скорректированный коэффициент детерминации, в котором используются несмещённые оценки дисперсий:

$$R_{adj}^2 = 1 - \frac{s^2}{s_y^2} = 1 - \frac{SS_{res}/(n-k)}{SS_{tot}/(n-1)} = 1 - (1 - R^2) \frac{(n-1)}{(n-k)} \leq R^2 \quad (16.18)$$

который даёт “штраф” за дополнительно включённые факторы, где n — количество наблюдений, а k — количество параметров.

Данный показатель всегда меньше единицы, но теоретически может быть и меньше нуля (только при очень маленьком значении обычного коэффициента детерминации и большом количестве факторов). Поэтому теряется интерпретация показателя как “доли”. Тем не менее, применение показателя в сравнении вполне обоснованно.

Для моделей с одинаковой зависимой переменной и одинаковым объёмом выборки сравнение моделей с помощью скорректированного коэффициента детерминации эквивалентно их сравнению с помощью остаточной дисперсии $s^2 = SS_{res}/(n - k)$ или стандартной ошибки модели s . Разница только в том, что последние критерии чем меньше, тем лучше.

Для оценки качества различных моделей при описании данных также используют информационные критерии.

Определение 16.3

Информационный критерий — применяемая в статистике мера относительного качества статистических моделей, учитывающая степень “подгонки” модели под данные с корректировкой (“штрафом”) на используемое количество оцениваемых параметров. То есть критерии основаны на некоем компромиссе между точностью и сложностью модели. Критерии различаются тем, как они обеспечивают этот баланс.

Информационный характер критериев связан с концепцией информационной энтропии и расстоянием Кульбака-Лейблера, на основе которой был разработан исторически первый критерий — критерий Акаике (AIC), предложенный в 1974 году Хиротсугу Акаике (H. Akaike, “A new look at the statistical model identification” IEEE Transactions on Automatic Control, vol. 19, № . 6, pp. 716-723, 1974.)

Информационные критерии используются исключительно для сравнения моделей между собой, без содержательной интерпретации значений этих критериев. Они не позволяют тестировать модели в смысле проверки статистических гипотез. Обычно чем меньше значения критериев, тем выше относительное качество модели.

Информационный критерий Акаике (AIC)

Предложен Хиротугу Акаике в 1971 году, описан и исследован им же в 1973, 1974, 1983 годах. Первоначально аббревиатура AIC, предложенная автором, расшифровывалась как an information criterion (“некий информационный критерий”), однако последующие авторы называли его Akaike information criterion. Исходная расчетная формула критерия имеет вид:

$$AIC = 2k - 2L \quad (16.19)$$

где L - значение логарифмической функции правдоподобия построенной модели, k -количество использованных (оцененных) параметров.

Многие современные авторы, а также во многих программных продуктах применяется несколько иная формула, предполагающая деление на объем выборки n , по которой строилась модель:

$$AIC = \frac{2k - 2L}{n}. \quad (16.20)$$

Данный подход позволяет сравнивать модели, оцененные по выборках разного объема.

Чем меньше значение критерия, тем лучше модель.

Многие другие критерии являются модификациями AIC.

Байесовский информационный критерий (BIC) или критерий Шварца (SC)

Байесовский информационный критерий (Bayesian information criterion — BIC) предложен Шварцем в 1978 году, поэтому часто он называется также критерием Шварца (Schwarz criterion — SC). Он разработан исходя из байесовского подхода и является наиболее часто используемой модификацией AIC:

$$\text{BIC} = \text{SC} = k \ln n - 2L . \quad (16.21)$$

Как видно из формулы, данный критерий налагает больший штраф на увеличение количества параметров по сравнению с AIC, так как $\ln n$ больше 2 уже при количестве 8 наблюдений.

Прочие информационные критерии

Состоятельный критерий Акаике (Consistent AIC — CAIC) предложенный в 1987 году Боздоганом:

$$\text{CAIC} = (1 + \ln n)k - 2L. \quad (16.22)$$

Данный критерий асимптотически эквивалентен BIC. Тот же автор в 1994 году предложил модификации, увеличивающие коэффициент при количестве параметров (вместо 2-3 или 4 для AIC_3 и AIC_4).

Скорректированный критерий Акаике (Corrected AIC – AIC_c), который рекомендуется применять на малых выборках (предложен в 1978 году Sugiura):

$$AIC_c = AIC + \frac{2k(k+1)}{n-k-1}. \quad (16.23)$$

Данный критерий, наряду с AIC и BIC выдается в результатах оценки моделей в Wolfram Mathematica при использовании оператора `NonlinearModelFit`.

Критерий Ханнана-Куинна (Hannan-Quinn, HQ) предложен авторами в 1979 году

$$HQ = 2k \ln \ln n/n - 2L/n. \quad (16.24)$$

Имеются также модификации AIC, использующие более сложные штрафные функции, зависящие от различных характеристик.

Замечание

Высокие значения коэффициента детерминации, вообще говоря, не свидетельствуют о наличии причинно-следственной зависимости между переменными (также как и в случае обычного коэффициента корреляции). Например, если объясняемая переменная и факторы, на самом деле не связанные с объясняемой переменной, имеют возрастающую динамику, то коэффициент детерминации будет достаточно высок. **Поэтому логическая и смысловая адекватность модели имеют первостепенную важность.** Кроме того, необходимо использовать критерии для всестороннего анализа качества модели.

17. Рекомендуемая литература

Рекомендуемая литература

1. Лавренчик, В.Н. Постановка физического эксперимента и статистическая обработка его результатов/ В.Н. Лавренчик.. – М.: Энергоатомиздат, 1986. - 272 с.
2. Бендат Дж., Пирсол А. Прикладной анализ случайных данных/ Дж.Бендат, А.Пирсол. – М.: Мир,1989. -504 с.
3. Новицкий, П.В. Оценка погрешностей результатов измерений/ П.В.Новицкий, И.А..Зограф - Л.: Энергоатомиздат, 1991. - 304 с
4. Тейлор Дж. Введение в теорию ошибок/ Дж.Тейлор. - М.: Мир, 1985. - 45 с.
5. Гмурман, В. Е. Теория вероятностей и математическая статистика: Учеб. пособие для вузов/В. Е. Гмурман. - 9-е изд., стер. - М.: Высш. шк., 2003. - 479 с.

Рекомендуемая литература

6. Дьяконов, В.П. Mathematica 4: учебный курс/ В.П. Дьяконов. - СПб.: Питер, 2001 . - 654 с.
7. Воробьев Е. М. Введение в систему Mathematica./ Е.М.Воробьев. - М.: Финансы и статистика, 1998. - 345 с.
8. Львовский, Б.Н. Статистические методы построения эмпирических формул: Учеб. пособие для втузов/ Б.Н.Львовский. – М.: Высш. шк., 1988 – 239 с.
9. Кобзарь А. И. Прикладная математическая статистика. Для инженеров и научных работников/ А. И. Кобзарь. – М.: Физматлит, 2006. – 816 с.
10. Боровков Л. Л. Математическая статистика/ Л. Л.Боровков-Учебник.- М.: Наука. Главная редакция физико-математической литературы, 1984. – 472 с.

11. Кассандрова, О. Н. Обработка результатов наблюдений/ О. Н. Кассандрова, В. В. Лебедев- М.:Наука, Главная редакция физ.-мат. литературы, 1970 г. – 104 с.
12. Ивченко, Г. И. Математическая статистика/ Г.И.Ивченко, Ю. И. Медведев Учеб. пособие для втузов. – М.: Высш. шк., 1984. – 248 с.

Учреждение образования
“ГОМЕЛЬСКИЙ ГОСУДАРСТВЕННЫЙ
УНИВЕРСИТЕТ
ИМЕНИ ФРАНЦИСКА СКОРИНЫ”
Кафедра теоретической физики

Статистические методы обработки данных

Лабораторный практикум
Специальность
1-31 04 08 Компьютерная физика

Материал подготовил
Андреев
Виктор Васильевич
доктор физ.-мат. наук, доцент

Гомель, 2018

ОГЛАВЛЕНИЕ

Введение	4
Лабораторная работа № 1 Введение в Wolfram Mathematica	6
1.1 Общая информация	7
1.2 Ячейки	9
1.3 Вычисления	10
1.4 Палитры	12
1.5 Работа со списками	12
1.6 Порядок выполнения работы. Задания	13
Лабораторная работа № 2 Построение графиков статических распределений распределений	17
2.1 Построение графиков в Wolfram Mathematica	17
2.2 Порядок выполнения работы. Задания	18
Лабораторная работа № 3 Точечные и интервальные оценки случайной величины	22
3.1 Доверительные интервалы для математического ожидания и дисперсии	22
3.2 Порядок выполнения работы. Задания	23
Лабораторная работа № 4 Графическое представление экспериментальных данных в системе Mathematica	25
4.1 Графическое представление эмпирических данных	25
4.2 Оптимальное число интервалов для получения гистограммы	26
4.3 Порядок выполнения работы. Задания	29
Лабораторная работа № 5 Промахи и методы их исключения	31
5.1 Промахи и методы их исключения	31
5.2 Порядок выполнения работы. Задания	33
Лабораторная работа № 6 Идентификация форм распределения экспериментальных данных	34
6.1 Критерий χ^2 (критерий Пирсона)	34
6.2 Порядок выполнения работы. Задания	36

Лабораторная работа № 7 Фитирование экспериментальных данных	38
7.1 Средневзвешенное значение	38
7.2 Поиск параметров регрессионной модели	39
7.3 Порядок выполнения работы. Задания	40
Список использованных источников	43

Репозиторий ГГУ им. Ф. Скорины

Введение

Современный инженер должен быть подготовлен к организационно-управленческой, научной, производственной и преподавательской деятельности в области экспериментального исследования физических процессов на различных уровнях структурной организации материи при различных физических условиях; теоретического анализа эффектов и явлений и предсказания новых физических закономерностей на основе современных теоретических представлений, математических и компьютерных методов; разработке приборов на основе новых материалов и физических принципов, созданию новых технологий, использующих математические методы и компьютерную технику; работе по математическому моделированию разнообразных процессов и объектов.

Одним из обязательных требований такого рода специалистов является их способность планировать, организовывать и проводить физические и радиофизические эксперименты, применяя вероятностные методы анализа и используя технические средства автоматизации эксперимента; обрабатывать и анализировать полученные результаты, создавать математические модели и программные средства, составлять отчеты и вести научно-техническую документацию.

Как известно статистическая обработка информации в физических исследованиях играет одну из важнейших ролей. Новые неизвестные явления открываются в настоящее время путем тщательной обработки результатов измерений. Современный физик обязан знать основы статистической обработке данных и с помощью программных продуктов уметь проводить анализ и обработку физической информации на персональном компьютере. Исходя из этих требований, следует что, такой курс, как “Статистические методы обработки данных” является необходимой составной частью образования современного инженера-физика.

Целью курса “Статистические методы обработки данных” является овладение студентами основами математической статистики и теории вероятностей и их использование при обработке и анализе экспериментальных данных. Задачами курса “Статистические методы обработки данных” являются приобретение навыков и изучение способов, приёмов работы по анализу и статистической обработке информации на персональном компьютере, наиболее часто используемых в физических исследованиях (обратив внимание на смысл понятий критериев, возможности алгоритмизации), приобретение

навыков по решению теоретических и экспериментальных задач различных физических дисциплин, таких как квантовая механика, физика ядра и атома, радиационная безопасность и др., овладение новыми методами теоретических и экспериментальных исследований (решение практических задач, методы получения и обработки результатов опытов).

Материал курса “Статистические методы обработки данных” базируется на ранее полученных знаниях по курсу “Теория вероятностей”. Полученные навыки могут быть использованы при изучении дисциплин учебного плана, таких как “Физика ядра”, “Квантовая механика”. Данный курс формирует необходимую базу для научных исследований студентов и выполнения курсовых и дипломных работ. В результате изучения специального курса “Статистические методы обработки данных” студент должен знать и уметь:

- способы, приёмы работы для процесса обработки физической информации;
- применять системы аналитических вычислений для решения задач математической статистики, использовать технические средства автоматизации эксперимента;
- обрабатывать и анализировать полученные результаты, создавать математические модели и программные средства;
- планировать и организовать научные исследования, применять соответствующие экспериментальные и теоретические методы интерпретации результатов экспериментальных исследований.

Для углубленного изучения системы **Wolfram Mathematica** рекомендуем книги [1–12] и конечно книгу Стивена Вольфрама [13]. Конечно представленные книги не исчерпывают всего, что имеется и пытливый студент может с легкостью дополнить его, исследовав информационные просторы Интернета.

Также полезно использовать литературу по теории вероятностей и математической статистике [14–26].

Примечание.

При выполнении лабораторных работ в блоке вида

Операторы WM[®]

Operator1, Operator2, ...

приводятся операторы **Wolfram Mathematica[®]**, которые могут быть использованы для выполнения заданий.

Лабораторная работа № 1

Введение в Wolfram Mathematica

Цель работы: Овладение некоторыми навыками работы в системе Wolfram Mathematica

Краткие теоретические сведения

Система Wolfram Mathematica является очень популярной во всем мире. Так, например, в США официально зарегистрировано свыше миллиона пользователей. Считается, и не без основания, что Wolfram Mathematica - лидер среди систем символьной математики. Высокие интеллектуальные возможности системы Wolfram Mathematica позволяют решать задачи в аналитическом виде. Преобразования математических выражений осуществляются на таком высоком уровне, что система позволяет получить решения большинства математических задач в аналитическом виде, выводить формулы.

Система позволяет:

- определять вещественные и комплексные корни алгебраических уравнений;
- решать алгебраические и дифференциальные уравнения;
- вычислять неопределенные и определенные интегралы;
- осуществлять интегральные преобразования и решать задачи оптимизации;
- осуществлять разложение функции в степенной ряд;
- находить пределы, вычислять суммы и произведения математических функций;
- осуществлять упрощения сложных математических выражений до таких уровней, когда выражение становится формулой осуществлять самопроверку результатов решения задач.

Wolfram Mathematica это мощная вычислительная система. Она позволяет без программирования получать численные решения большинства задач прикладной математики. Система поражает объемом вычислений. Например,

функции $\pi, e, n!$ вычисляет практически с любым числом знаков. Она способна выполнять математические действия с абсолютной точностью. При этом количество цифр не ограничено. **Wolfram Mathematica** позволяет также выполнять вычисления с произвольной точностью.

Wolfram Mathematica - справочная математическая система. В считанные секунды пользователь получит таблицы логарифмов, элементарных и специальных функций, таблицы производных, интегралов, сумм и произведений и т.д. Решение задач осуществляется в режиме диалога и не требует программирования. Язык общения системы - язык функционального программирования высокого уровня. Его можно отнести к классу интерпретаторов, когда система анализирует (интерпретирует) введенное выражение и сразу его исполняет. Система компьютерной алгебры является математическим справочником высоко уровня. Существуют и поддерживаются версии **Wolfram Mathematica** на немецком, французском и японском языках (с целым интерфейсом, включающим в себя меню, палитры, окна диалога, сообщения и предупреждения об ошибках, и более тысячи страниц помощи).

1.1 Общая информация

Mathematica включает в себя редактор (**FrontEnd**) и вычислительное ядро (**Kernel**).

Запуск системы **Wolfram Mathematica** в среде операционной системы **WINDOWS** осуществляется традиционными для этой системы способами: 1) Поиск соответствующего ярлыка на рабочем столе и запуск его двумя ударами (желательно пальцами) по правой кнопки мыши, если ваша мышь и вы правша и соответственно по левой кнопке, если вы левша и ваша мышь для левши. 2) Активизация позиции **Wolfram Mathematica** в главном меню программ.

Wolfram Mathematica программно состоит из нескольких отдельных частей: интеллектуального ядра (**MathKernel**), интерфейсного процессора (**front end**).

Выход из системы осуществляется командой **EXIT** из меню **FILE** или нажатием кнопки с крестом (**x**) в правом верхнем углу.

После того как вы научились входить и что самое главное, выходить из системы, приступим к знакомству с интерфейсом системы. Запустите систему, если вы еще не сделали этого и приступим к знакомству:

Главное меню системы (**front end**) (смотри на экран) содержит следующие позиции:

- ◇ **File** – работа с файлами: создание нового файла, выбор файла из каталога, закрытие файла, запись текущего файла, запись файла с изменением имени, печать документа и завершение работы;
- ◇ **Edit** – основные операции редактирования (отмена операции, копирование выделенных участков документа в буфер с их удалением и без удаления, перенос выделенных участков, их стирание);
- ◇ **Format** – управление форматом документов;
- ◇ **Cell** – работа с ячейками (объединение и разъединение ячеек, установка статуса ячейки, открытие и закрытие);
- ◇ **Graphics** – работа с графическими объектами;
- ◇ **Insert** – задание элементов ввода (графиков, матриц, гиперссылок и т. д.);
- ◇ **Palettes** – работа с палитрами (инструментарий, облегчающий формирование документа);
- ◇ **Evaluation** - управление ядром системы и процессом вычислений;
- ◇ **Window** – операции с окнами и их расположением;
- ◇ **Help** – управление справочной системой. Часть команд может быть в данный момент невыполнима, например, нельзя вычислить значение выражения, если его самого нет в окне редактирования или если ячейка с ним не выделена. Названия таких команд выделяются характерным серым расплывчатым шрифтом. Четкий шрифт, напротив, характерен для тех команд, которые в данный момент могут исполняться. Управление главным меню самое обычное.

Важнейшим источником информации является опция **HELP**. Выбор пункта **Documentation Center** введет вас в мир полезной, как для сдачи зачета, так и для желающих научиться использовать систему для решения инженерно-технических и научных задач, информации. Единственным неудобным моментом является, что данный помощник написан на английском языке. Браузер справки вызывается командой **Help Browser (Shift+F1)**. Справочная система позволяет:

- получить сведения обо всех командах главного меню;
- изучить правила записи и набора математических выражений;

- уточнить назначение любой функции или оператора;
- ознакомиться с примерами и приспособить их к интересам пользователя (примеры имеются по каждой функции);
- получить доступ к пакетам расширения;
- получить справку о системе **Wolfram Mathematica** и фирме ее разработавшей;
- воспользоваться электронной книгой разработчика системы Стивена Вольфрама (раздел справки **The Mathematica Book**).

Справку можно получить достаточно просто по имени функции или алфавитному указателю. Таким образом, **Wolfram Mathematica** обладает возможностями компьютерной алгебры. Можно назвать её также системой компьютерной математики. **Wolfram Mathematica** позволяет создавать отчеты, презентации и другие документы, что актуально в настоящее время.

1.2 Ячейки

Визуально каждый файл с расширением ***.nb** (notebook или записная книжка) состоит из набора ячеек (синие скобки в правой части).

Форматами ячеек могут быть следующие:

- **InputForm** (**Shift+Ctrl+I**) – формат ввода;
- **OutputForm** – формат вывода;
- **StandardForm** (**Shift+Ctrl+N**) – стандартный формат;
- **TradinationalForm** (**Shift+Ctrl+T**) – традиционный формат (близкий к обычному математическому стилю);
- **Bitmap** – растровый формат изображений;
- **Metafile** – векторный графический формат **Windows Metafile**.

При работе с текстами с большим числом математических знаков целесообразно использовать стандартный формат. Подменю **Cell Properties** устанавливает свойства ячеек. Оно содержит следующие команды:

- **Cell Open** – устанавливает ячейку открытой или закрытой;

– **Cell Editable** – устанавливает ячейку редактируемой или не редактируемой..

Ячейка ввода и соответствующая ей ячейка вывода обрамляются справа скобками: одиночными и общей. Активизируя скобку двойным щелчком, можно скрывать и снова выводить на экран выходную ячейку. Это полезно в том случае, если результат в ней слишком громоздкий. Редактировать можно содержимое как входной, так и выходной ячеек. Для этого выходную ячейку необходимо сделать редактируемой, установив свойство **Cell Editable**.

Ячейка создается автоматически когда вы начинаете вводить текст. Каждая ячейка может содержать либо выражение, либо его вывод, либо текст.

Стиль ячейки можно выбрать в меню **Format > Style**. Основные стили ячеек:

- **Input** — выражение для вычисления;
- **Output** — вывод результата вычисления;
- **Title** — заголовок файла;
- **Section** — раздел (секция) файла;
- и множество других.

Ячейки можно удалять, переносить, окрашивать в разные цвета и т. д. Все основные рации вы найдете в меню **Format** и **Cell**.

1.3 Вычисления

Для вычисления выражения нужно стать в его ячейку и нажать **Shift+Enter** или **Enter** на калькуляторе клавиатуры (можно через меню **Evaluation**):

```
In[1] := (2+3)*4
```

```
Out[1]= 20
```

Часто бывает удобно вводить несколько выражений в одну ячейку. Например, создадим матрицу и вычислим ее определитель.

```
In[2] := Range[9]
```

```
y=Partition[%,3]
```

```
y//MatrixForm
```

```
Det[y]
```

```
Out[2]= {1, 2, 3, 4, 5, 6, 7, 8, 9}
```

```
Out[3]= {{1,2,3},{4,5,6},{7,8,9}}
```

```
Out[4] =  $\begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix}$ 
```

```
Out[5]= 0
```

Знак процента (%) означает последний результат. В данном случае мы сначала создали упорядоченный список от 1 до 9, потом разбили его на три части (то есть получили матрицу) и присвоили результат переменной y , потом получили ответ для y в матричном виде (для наглядности), а затем вычислили определитель матрицы. Применение процента позволяет видеть промежуточные результаты. Иначе все можно записать в одну строку (без матричного вида):

```
In[6]:= Det[Partition[Range[9],3]]
```

```
Out[6]= 0
```

Чтобы не видеть промежуточные результаты, нужно поставить в конце выражения точку с запятой.

```
In[7]:= Range[9];
```

```
Partition[%,3];
```

```
Det[%]
```

```
Out[8]= 0
```

Все переменные, как обычно, сохраняются в память пока вы работаете с программой. Но в *Mathematica* есть понятие вычислительной сессии. Обычно

сессия начинается при запуске программы и заканчивается при ее закрытии. Однако, вы можете прекратить сессию когда угодно — для этого нужно выбрать меню `Evaluation > Quit Kernel > Local`. При этом Mathematica “забудет” значения всех сохраненных переменных. Это бывает полезно, если посыпались ошибки и вам нужно пересчитать все заново.

Для очистки полезно использовать операторы, которые “очищают” данные о переменных и функциях, которые вы используете при вычислениях. Новые вычисления полезно начинать с блока вида

```
ClearAll;
```

```
ClearAll["Global'*"];
```

1.4 Палитры

В Mathematica для ввода выражений очень удобно пользоваться палитрами. Все палитры собраны в меню `Palettes`. Базовой является `Palettes > Basic Math Assistant`, которая позволяет:

- ◇ Вводить арифметические выражения, корни, степени, математические константы;
- ◇ Вводить тригонометрические функции, а также функции для упрощения выражений, суммы, производные, интегралы;
- ◇ Позволяет создавать списки и матрицы;
- ◇ Дает доступ к графическим функциям
- ◇ и многое другое.

Со временем вы научитесь многие операторы и объекты системы вводит без палитр.

1.5 Работа со списками

Операторы WM[®]

{}, Range, Table, Part, Span, Position, First, Last, Rest, Take, Insert, Append, Prepend, Drop, Delete, Rest, Length, Count, Total, Join, Union, Intersection, Complement, Sum, Product, Select, MemberQ, FreeQ, Flatten

Создание списка Создать небольшой список проще всего вручную.

```
In[30]:= {1, 2, 3, 4, 5}
```

```
Out[30]= {1, 2, 3, 4, 5}
```

Внутри списка могут быть и подсписки (и подподсписки, и подподподсписки, и т.д.).

```
In[31]:= {1, {2, {3, 4, {5, 6, 7}}}}
```

```
Out[31]= {1, {2, {3, 4, {5, 6, 7}}}}
```

Самая общая функция для задания списков и матриц — `Table`. Она позволяет задавать сколько угодно измерений + формулу для общего члена. В качестве формулы можно использовать любое выражение (встроенную функцию, свою функцию, комбинацию функций и т.п.). Например, вот так можно получить таблицу умножения чисел от 1 до 5.

Важно: Функция `MatrixForm` используется для красивого отображения матриц. Никогда не сохраняйте ее в переменную, так как вы потом не сможете работать с такой переменной как с матрицей/списком.

1.6 Порядок выполнения работы. Задания

1. Внимательно прочитайте теоретические сведения по работе с САВ `Mathematica`.
2. Отработайте вход и выход из САВ.
3. Ознакомьтесь с палитрами математических операций и функций.
4. Изучите возможности `front end` для работы с документами.
5. Выполните задания ниже.

Задание 1.1

Используя справочную систему найдите информацию по следующим операторам САВ `Mathematica`:

1. `Pi`;

2. N;
3. List;
4. Table;
5. ArcSin;
6. String;
7. Print;
8. Import;
9. Export.

В отчете приведите информацию о каждом перечисленном операторе: для чего предназначен, какие варианты имеются (синтаксис оператора, сокращенная запись) для его использования (с краткими комментариями на русском языке), примеры использования (2-3 примера на каждый оператор).

Операторы WM®

Precision (')

Только истинный ценитель и специалист программирования может понять возможности системы, выполнив

Задание 1.2

Вычислите:

1. 1234463^4
2. 2^{200}
3. 2.0^{200} (Сравните результат с пунктом 2. и сделайте выводы об отличии результатов (в тексте лабораторной работы))
4. $3^{200}/123446^3$
5. $200!$ (трижды различными способами)

6. Вычислите e^{81} с точностью до 27 цифр после запятой (смотри операторы `N`, `Precision`, `SetPrecision`, `SetAccuracy`).
7. Выяснить, что больше 4100109999999999990^5 или $64!$ (сделайте 2 способа проверки).

Операторы WM®

`TableForm`, `MatrixForm`, `Import`, `Export`, `SetDirectory`, `NotebookDirectory`, `StringJoin (<>)`

Задание 1.3

1. $y(x) = \sin^2(x) + x^3$ для значений $x = 0,1, 2, 3, 4,5, 7,32, \pi, e$.
2. $y(x) = \ln(|x|)$ для значений $x = \{-5, -1, 0, 5,5, ; 7,3222, 10\}$;
3. $z(y,x) = \sin(x^2 + y^2)$ для значений $(x,y) = (-5, -1), (0,5,5), (7,32, 22,1), (10,\pi)$

Результаты этого задания запишите на ПК в формате, совместимым с Excel (*.xls или *.xlsx).

Операторы WM®

`Take`, `Length`, `Text`, `Grid`, `TableForm`

Задание 1.4

Считать данные для своего варианта из таблицы 1.1 (номер вашей фамилии в журнале группы соответствует номеру варианта) и построить таблицу из двух колонок: Первая колонка (условно назовем колонкой x) – это данные из первого файла, а вторая колонка – данные из второго файла (условно назовем колонкой y). Если файлы разной длины, то сделайте их одинаковой длины

Таблица 1.1–
Данные для считывания по вариантам

Вариант	Номера файлов
1	d1.dat, d11.dat, d16.dat
2	d2.dat, d12.dat, d17.dat
3	d3.dat, d13.dat, d18.dat
4	d4.dat, d14.dat, d19.dat
5	d5.dat, d15.dat, d20.dat
6	d6.dat, d23.dat, d21.dat
7	d7.dat, d24.dat, d22.dat
8	d8.dat, d25.dat, d16.dat
9	d9.dat, d11.dat, d17.dat
10	d10.dat, d12.dat, d18.dat
11	d10.dat, d13.dat, d19.dat
12	d9.dat, d14.dat, d20.dat
13	d8.dat, d15.dat, d21.dat
14	d7.dat, d23.dat, d22.dat
15	d6.dat, d24.dat, d16.dat
16	d5.dat, d25.dat, d17.dat
17	d4.dat, d25.dat, d18.dat
18	d3.dat, d24.dat, d19.dat
19	d2.dat, d23.dat, d20.dat
20	d1.dat, d15.dat, d21.dat
21	d1.dat, d14.dat, d22.dat
22	d2.dat, d13.dat, d16.dat
23	d3.dat, d12.dat, d17.dat
24	d4.dat, d11.dat, d18.dat
25	d5.dat, d23.dat, d19.dat

Рекомендация: Чтобы нечаянно “не переписать” свои dat-файлы, сохраните копию отдельно.

Лабораторная работа № 2

Построение графиков статических распределений распределений

Цель работы: Изучить графические возможности системы Wolfram Mathematica отображению данных экспериментов и функций распределения вероятностей. Расчет моментов теоретических распределений с помощью встроенных функций Wolfram Mathematica.

Краткие теоретические сведения

2.1 Построение графиков в Wolfram Mathematica

Операторы WM®

Plot, ListPlot, ListLinePlot, DiscretePlot, ParametricPlot, LogPlot, LogLinearPlot, LogLogPlot, PolarPlot, Plot3D, ContourPlot, Graphics, Show

Концептуально графики в системе Wolfram Mathematica являются графическими объектами, которые создаются (возвращаются) соответствующими графическими функциями. Их немного, около десятка, и они охватывают построение практически всех типов математических графиков. Как уже отмечалось, достигается это за счет применения опций и директив. Поскольку графики являются объектами, то они могут быть значениями переменных.

Поэтому Wolfram Mathematica допускает следующие конструкции:

1. `Plot[Cos[x], {x, 0, 20}]` – построение графика косинуса для аргумента x пределах от 0 до 20;
2. `g:=Plot[Cos[x], {x, 0, 20}];` – задание графического объекта - графика косинуса с отложенным выводом (**будет вычисляться при следующем появлении переменной g**);
3. `g=Plot[Cos[x], {x, 0, 20}]` – задание объекта - графика косинуса с немедленным выводом (вычисляется сразу в данном месте).

Начнем рассмотрение графических возможностей системы с построения простейших графиков функций одной переменной вида $y = f(x)$ (функция задана явно) или просто $f(x)$. График таких функций строится на плоскости, то есть в двумерном пространстве. При этом используется прямоугольная (декартова) система координат. График представляет собой геометрическое положение точек (x, y) при изменении независимой переменной (абсциссы) в заданных пределах, например от минимального значения x_{\min} до максимального x_{\max} с шагом dx . По умолчанию строятся и линии координатной системы.

Для построения двумерных графиков функций вида $f(x)$ используется встроенная в ядро функция `Plot`:

1. `Plot[f, {x, xmin, xmax}]` – возвращает объект, представляющий собой график функции f аргумента x в интервале от x_{\min} до x_{\max} ;
2. `Plot[{f1, f2, ...}, {x, xmin, xmax}]` – возвращает объект в виде графиков ряда функций $f1, f2, \dots$ аргумента x .

2.2 Порядок выполнения работы. Задания

Операторы WM[®]

PDF, CDF, Mean, Variance, Expectation, PoissonDistribution, BinomialDistribution, Total, Sum

Задание 2.1

Постройте графики функций (см. конспект лекций или Document Center):

1. $y(x) = \operatorname{tg}^3(x)$;
2. Построить эллипс с полуосями $a = 2$ и $b = 5$, используя полярную систему координат;
3. $z(y, x) = \sin(x^2 + y^2)$.

Области изменения x и y выберите таким образом, чтобы графики имели эстетический приемлемый вид. Подпишите оси графиков.

Задание 2.2

Построить графики экспоненциальных распределений, где плотность распределения задается формулой

$$p(x) = A(\alpha) \exp(-|x|^\alpha),$$
$$A(\alpha) = \frac{\alpha}{2\Gamma(1/\alpha)} \quad (2.1)$$

для $\alpha = 1, 2, 3, 5, 10$ на одном графике, **отметив каждый их них своей меткой**.

Сделать вывод о зависимости формы экспоненциальных распределений от параметра α (описать словесно чем отличаются распределения друг от друга).

Задание 2.3

Построить графики дискретных и непрерывных с помощью встроенных функций Mathematica:

1. плотности биномиального распределения с параметрами $n = 100, p = 0.4$ (относится к классу дискретных распределений);
2. плотности распределения Пуассона с параметрами $n = 10, 20, 50$ (относится к классу дискретных распределений);
3. плотности нормального распределения с параметрами $a = 0, \sigma = 1, 2, 3$ (один график) и $a = 10, \sigma = 5, 6$ (один график);
4. интегрального закона нормального распределения с параметрами $a = 0, \sigma = 1, 2, 3$;
5. интегральный и дифференциальные законы распределения χ^2 -квадрат для 2, 3 и 7 степеней свободы.
6. интегральный и дифференциальные законы распределения Стьюдента (t -распределение) для 2, 3 и 7 степеней свободы.

Задание 2.4

Распределение Максвелла для молекул идеального газа по скоростям можно разделить на

- распределение по проекции скорости;
- распределение по модулю скоростей.

Плотность распределения Максвелла для вектора скорости молекулы $\{v_x, v_y, v_z\}$ является произведением распределений для каждого из трех направлений:

$$p_v(v_x, v_y, v_z) = p_v(v_x)p_v(v_y)p_v(v_z), \quad (2.2)$$

где распределение по одному направлению (x) имеет вид:

$$p_v(v_i) = \sqrt{\frac{m}{2\pi kT}} \exp\left[\frac{-mv_i^2}{2kT}\right]. \quad (2.3)$$

В (2.3): m — масса одной молекулы газа (не путайте с молекулярной массой, выраженной в атомных единицах); T — термодинамическая температура и k — постоянная Больцмана.

Обычно, в учебниках по молекулярной физике приводят распределение по абсолютному значению $v = |\mathbf{v}|$. Модуль скорости, v определяется как:

$$v = \sqrt{v_x^2 + v_y^2 + v_z^2}, \quad (2.4)$$

функция плотности вероятности для модуля скорости равна

$$p(v) = 4\pi v^2 \left(\frac{m}{2\pi kT}\right)^{3/2} \exp\left(\frac{-mv^2}{2kT}\right). \quad (2.5)$$

1. Построить графики распределения плотности вероятности (2.3) и интегрального закона вероятности для проекции скорости v_x для водорода;
2. Построить графики распределения плотности вероятности (2.5) и интегрального закона вероятности для модуля скорости для водорода.
3. Найти для распределения по модулю скоростей:
 - (a) наиболее вероятную скорость v_p ;
 - (b) среднюю арифметическую скорость $\langle v \rangle$;
 - (c) среднеквадратичную скорость $\langle v^2 \rangle = E\{v^2\}$;
 - (d) найти численные значения этих скоростей для водорода.

4. **Ответить на вопрос (в отчете):** К каким распределениям в математической статистике можно отнести плотности вероятности (2.3) и (2.5)?

Задание 2.5

С помощью встроенных функций Mathematica:

1. Найти математическое ожидание, дисперсию, коэффициент асимметрии и эксцесс нормального распределения с параметрами a и σ .
2. Найти математическое ожидание, медиану, дисперсию, коэффициент асимметрии и эксцесс распределения χ^2 с ν степенями свободы.

Примечание. Поскольку уместить всю необходимую информацию в лабораторной работе невозможно, то часть нужной Вам информации **возьмите из конспекта лекций, которые у Вас, как всех прилежных студентов конечно имеется.** Выполнение этой и последующих работ обязательно должно сопровождаться текстовыми комментариями.

Лабораторная работа № 3

Точечные и интервальные оценки случайной величины

Цель работы: Вычисление точечных и интервальных оценок для одномерной случайной величины с помощью встроенных функций Wolfram Mathematica и их сравнение с формулами из литературы.

Краткие теоретические сведения

3.1 Доверительные интервалы для математического ожидания и дисперсии

Доверительный интервал для математического ожидания $E\{x\} = \xi$ (двусторонний), если неизвестна дисперсия распределения:

$$\left[\bar{x} - \frac{t_{n-1, 1-\alpha/2} S}{\sqrt{n}} < \xi < \bar{x} + \frac{t_{n-1, 1-\alpha/2} S}{\sqrt{n}} \right], \quad (3.1)$$

где $t_{n,P}$ определяется соотношением:

$$\int_{-\infty}^{t_{n,P}} p_{ST}(t) dt = P.$$

Коэффициенты $t_{n,P}$ носят название **коэффициентов Стьюдента**.

Для построения доверительного интервала для дисперсии $D\{x\} = \sigma^2$ используется, то, что функция

$$\frac{(n-1) S^2}{\sigma^2} \quad (3.2)$$

имеет распределение χ^2 с $n-1$ степенями свободы.

Тогда доверительный интервал для дисперсии при неизвестном математическом ожидании имеет вид:

$$\left[\frac{(n-1) S^2}{\chi_{n-1, \alpha/2}^2} < D\{x\} < \frac{(n-1) S^2}{\chi_{n-1, 1-\alpha/2}^2} \right], \quad (3.3)$$

где $\chi_{n,\alpha}^2$ – квантиль распределения χ^2 :

$$\int_{\chi_{n,\alpha}^2}^{\infty} p(\chi^2) d\chi^2 = \alpha$$

с плотностью

$$p(\chi^2) = \frac{(\chi^2)^{(\frac{n}{2}-1)} e^{-\frac{\chi^2}{2}}}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})}, \quad (\chi^2 > 0).$$

Значения $E\{x\}$ и $D\{x\}$ лежат в интервале с доверительной вероятностью $P = 1 - \alpha$.

3.2 Порядок выполнения работы. Задания

Задание 3.1

Загрузить из файла выборку для своего варианта (например, если вариант 7, то извлекаются данные из файла d7.dat).

С помощью формул в конспекте найти точечные оценки физической случайной величины (несмещенные и состоятельные) для заданной выборки:

1. математического ожидания
2. дисперсии
3. среднеквадратичного отклонения
4. коэффициента асимметрии
5. эксцесса

Задание 3.2

Прodelать **Задание 3.1** с помощью встроенных функций Mathematica.

Задание 3.3

Провести анализ и сравнение полученных данных из двух заданий. **Сделать вывод о соответствии или несоответствии формул из конспекта и встроенных в WM функций.**

Задание 3.4

С помощью формул из конспекта и встроенных функций Mathematica найти интервальные оценки физической случайной величины для заданной выборки для доверительной вероятности $p = 1 - \alpha$, $p = 0,99$, $p = 0,95$, $p = 0,683$, т.е. найти:

1. доверительный интервал для математического ожидания;
2. доверительный интервал для дисперсии;
3. исследовать как изменяется доверительный интервал с изменением доверительной вероятности $p = 1 - \alpha$.

Сделать вывод о соответствии или несоответствии формул из конспекта и встроенных в WM функций.

Репозиторий ГГУ им. Ф. Скорины

Лабораторная работа № 4

Графическое представление экспериментальных данных в системе Mathematica

Цель работы: Изучить основные этапы и овладеть навыками построения гистограмм с помощью различных критериев выбора числа столбцов. Изучить графическое представление экспериментальных данных с помощью встроенных функций Wolfram Mathematica.

Краткие теоретические сведения

4.1 Графическое представление эмпирических данных

Цель обработки данных заключается в выявлении вида распределений случайных величин и оценки параметров установленного распределения.

Полученные экспериментальные данные представляют, как правило, в виде таблиц. Полученные таблицы удобно представить графически. Используя набор независимых наблюдений x_1, x_2, \dots, x_n случайной величины X , полезным первым шагом в исследовании поведения случайной величины является организация и представление их таким образом, чтобы их можно было легко интерпретировать и оценивать. Для достаточно большого количества наблюдаемых данных, полигон частот (распределения), гистограмма и кумулятивная линия является отличным графическим представлением данных, что облегчает оценку адекватности предполагаемой модели и оценку параметров распределения.

Гистограмма и полигон распределений являются графическим отображением частот, которые, в свою очередь, представляют собой оценки плотностей вероятностей $p(x)$. Кумулятивная линия - это график накопленных частот, в свою очередь оценивающих интегральную функцию распределения $F(x)$.

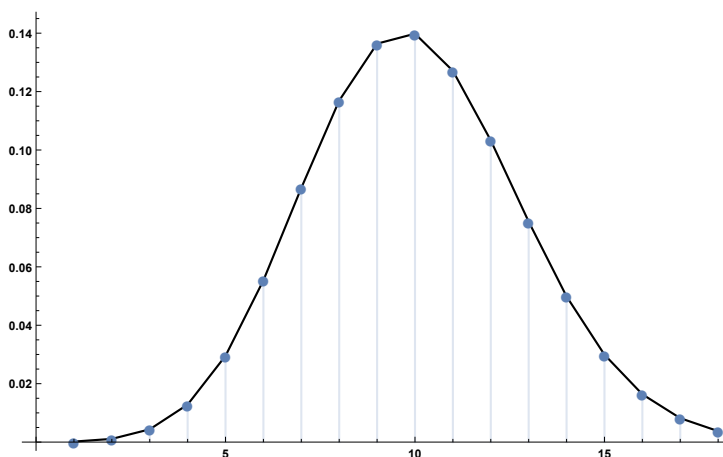


Рисунок 4.1– Пример полигона частот

Как строить полигон частот

1. Построить вариационный ряд для выборки, т. е. упорядочить значения случайной величины так, чтобы выполнялось условие: $x_1 \leq x_2 \leq \dots \leq x_n$.
2. Разбить область на m интервалов (бинов).
3. Построить точки на плоскости $x-\nu$ с координатами $\{\tilde{x}_i, \nu_i\}$, ($i = 1, \dots, n$), где

$$\tilde{x}_i = \frac{x_i + x_{i+1}}{2} \quad (4.1)$$

является серединой i -того интервала, ν_i - частота попадания случайной величины X в i -тый интервал.

Это же распределение можно представить в виде гистограммы. Для построения гистограммы необходимо над каждым отрезком оси абсцисс, соответствующим интервалу значений измеряемой величины, построить прямоугольник, площадь которого пропорциональна частоте попадания в этот интервал. Обычно выбирают интервалы одинаковой ширины, поэтому высота прямоугольников различна.

4.2 Оптимальное число интервалов для получения гистограммы

Как выбрать m и d .



Рисунок 4.2– Пример гистограммы

Оптимальное число

Оптимальное число существует!!

Оптимальное число интервалов группирования это такое число, когда ступенчатая огибающая гистограммы наиболее близка к плавной кривой распределения случайной величины.

Рекомендации по выбору m .

I группа: эвристические критерии (без доказательства).

Формула Старджеса

$$m = \log_2 n + 1 . \quad (4.2)$$

Формула Брукса и Каррузера

$$m = 5 \lg n . \quad (4.3)$$

Формула если $n > 100$

$$m = \sqrt{n} . \quad (4.4)$$

Эти три формулы являются наиболее часто встречающимися в литературе по математической статистике.

II группа: с использованием критерия χ^2 .

В ней используется рассмотрение интервалов не с равной длиной, а с **равной вероятностью** в соответствии с принимаемой моделью, т. е. предположением о законе распределения. В данном подходе неявно учитывается форма распределения.

Число интервалов с равной вероятностью, которые мы обозначили как K , отличаются от числа m с равной длиной d .

Г. Манн и А. Ваальд установили, что при $n \rightarrow \infty$ оптимальное число K равновероятных интервалов задается соотношением:

Критерий Мана-Ваальда

$$K \sim b\sqrt{2} \left(\frac{n}{Z_\alpha} \right)^{2/5}, \quad (4.5)$$

где $b = 2 \div 4$.

Здесь Z_α – квантиль нормального распределения, соответствующий вероятности $P = 1 - \alpha$, α – принятый уровень значимости.

$$Z_\alpha = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{Z_\alpha} e^{-\frac{x^2}{2}} dx = 1 - \alpha.$$

Критерий Мана-Ваальда

На практике часто берут $\alpha = 0.1$, тогда

$$K \simeq 1,9n^{2/5}. \quad (4.6)$$

В итоге приходим к таким рекомендациям при использования равновероятностных бинов K :

1. найти число ячеек, используя (4.6);
2. если окажется, что n/K мало, уменьшить K чтобы выполнялось неравенство $n/K \geq 5$;
3. Сформировать K равновероятностных ячеек гистограммы на основе данных. Отметим, что если измеряемая случайная величина многомерная, то существуют различные способы формирования ячеек с одинаковым вероятностным содержанием в K ячейках.

III группа.

Поскольку для K интервалы получаются не равной длины, то это приводит к ряду неудобств при построении гистограмм, но зато при этом мы неявно закладываем при использовании χ^2 выбор K в зависимости от формы распределения.

III группа рекомендаций “устраняет” недостаток II группы возвращаясь к интервалам m с равной длиной d , но при этом и учитывает, в отличие от I группы, форму распределения (форма характеризуется эксцессом ε или контрэксцессом κ).

Примером такого подхода является соотношение:

III группа (формула И.У.Алексеевой)

$$m = \frac{4}{\kappa} \lg \frac{n}{10}, \quad \kappa = \frac{1}{\sqrt{\varepsilon}}. \quad (4.7)$$

На практике эту формулу в зависимости от n при $\kappa = const$ удобнее аппроксимировать выражением

$$m = \frac{\varepsilon + 1,5}{6} n^{2/5} \quad (4.8)$$

4.3 Порядок выполнения работы. Задания

Для своего варианта:

Задание 4.1

Построить полигоны частот для данных выборок

Задание 4.2

Построить кумулятивные линии для данных выборок

Задание 4.3

Построить соответствующие заданию 4.1 гистограммы.

Примечание. Построение гистограмм провести с помощью встроенных функций Wolfram Mathematica 2 различными способами (не меньше).

Задание 4.4

Построить гистограммы, используя эвристические критерии и критерии с учетом формы распределения. Сравнить полученные гистограммы с гистограммами, полученными в задании 4.3 и **сделать выводы.**

Задание 4.5

Построить гистограммы, используя интервалы равной вероятности K для столбиков гистограммы.

Репозиторий ГГУ им. Ф. Скорины

Лабораторная работа № 5

Прوماхи и методы их исключения

Цель работы: Овладеть методами исключения промахов из значений случайной величины.

Краткие теоретические сведения

5.1 Прوماхи и методы их исключения

Одним из условий правомерности статистической выборки является требование ее однородности, т.е. принадлежности всех ее членов к одной и той же генеральной совокупности.

Однако на практике это требование очень часто нарушается. И, если скажем, при обработке вручную еще можно вспомнить как (при каких условиях) были получены “подозрительные” данные, то при автоматической обработке данных необходимы методы исключения “чужих” для данной выборки результатов.

Определение 5.1

*Отсчёты, резко отклоняющиеся по своим значениям от большинства других отсчетов принято называть **промахами** и исключать их из выборки.*

Важно. Если серия из небольшого числа измерений содержит грубую погрешность — промах, то наличие этого промаха может сильно исказить как среднее значение измеряемой величины, так и границы доверительного интервала. Поэтому из окончательного результата необходимо исключить этот промах.

Обычно промах имеет резко отличающееся от других измерений значение. Однако это отклонение от значений других измерений не дает еще права исключить это измерение как промах, пока не проверено, не является ли это отклонение следствием статистического разброса.

Особую неприятность доставляют отсчеты, которые и не входят в компактную группу отсчетов, но и не удалены от нее на значительное расстояние. Такой отсчет называют предполагаемым промахом. В экспериментальной

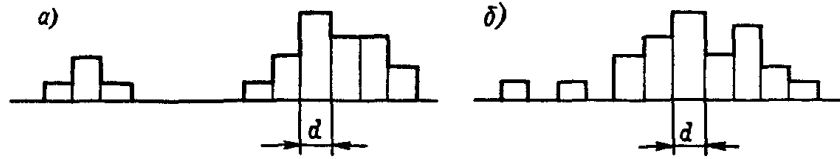


Рисунок 5.1– Возможные промахи

практике исследователи просто отбрасывали крайние, “слишком удаленные от центра наблюдения”. Эта процедура получила название **цензурирование выборки**.

Однако для принятия решения необходимы какие-либо формальные критерии.

Простейший метод заключается в использовании “правила 3σ ”, когда по выборке с удаленными отсчётами (предполагаемыми промахами) вычисляется оценка среднеквадратичного отклонения S и граница $|X_{\text{group}}| = 3\sigma$, а все $|x_i| \pm 3S$ отбрасываются.

Для расчета количества σ для цензурирования можно использовать аппроксимационные формулы:

Формулы для расчета количества σ ()

$$\begin{aligned}
 t_{\text{гр.}} &= 1,2 + 3,6 (1 - 1/\sqrt{\varepsilon}) \lg \left(\frac{n}{10} \right) , \\
 t_{\text{гр.}} &= 1,55 + 0,8\sqrt{\varepsilon - 1} \lg \left(\frac{n}{10} \right)
 \end{aligned}
 \tag{5.1}$$

Тогда интервал цензурирования выглядит следующим образом:

$$\bar{x} \pm t_{\text{гр.}} S ,
 \tag{5.2}$$

где S - точечная оценки среднеквадратичного отклонения (??).

5.2 Порядок выполнения работы. Задания

Задание 5.1

Сгенерировать выборки объема $n = 100$ и $n = 500$ для равномерного ($a = 1$, $b = \text{№}$ своего варианта) и нормальное ($\mu = \text{№}$ своего варианта, $\sigma = 6$) распределения при помощи встроенных функций Mathematica.

С помощью правила “ 3σ ” провести операцию цензурирования полученных выборок.

Задание 5.2

Для данных из задания 5.1 установите при каком количестве “ σ ” - $t_{\text{гр}}$. в результате процедуры цензурирования число элементов выборки уменьшится не более, чем на 1%. Сделайте качественный вывод о зависимости $t_{\text{гр}}$. от объема выборки n и формы распределения (ε).

Задание 5.3

Провести операцию цензурирования полученных выборок с помощью критериев, учитывающих форму распределения и объем выборки. Построить гистограммы до цензурирования и после. Сравнить результаты. Сделайте выводы.

Лабораторная работа № 6

Идентификация форм распределения экспериментальных данных

Цель работы: Овладеть методами идентификация форм распределения вероятностей.

Краткие теоретические сведения

6.1 Критерий χ^2 (критерий Пирсона)

Рассмотрим этапы необходимые для проверки гипотез на примере критерия χ^2 .

Процедура проверки гипотез с использованием критериев типа χ^2 предусматривает группирование наблюдений.

1. Выбираем “модель”: обычно это теоретическое дифференциальное распределение вероятностей $p(x, \theta)$, где θ – один (или несколько) параметров распределения.
2. Область определения случайной величины разбивают на m непересекающихся интервалов граничными точками $x_0, x_1, \dots, x_{m-1}, x_m$, где $x_0 < x_1 < \dots < x_{m-1} < x_m$.
3. В соответствии с заданным разбиением подсчитывают число n_i выборочных значений, попавших в i -й интервал и вероятности попадания в i -й интервал

$$\begin{aligned} p_i &= \int_{x_{i-1}}^{x_i} p(x, \theta) dx = \\ &= F(x_i) - F(x_{i-1}) \end{aligned} \quad (6.1)$$

соответствующие теоретическому закону с интегральной функцией распределения $F(x, \theta)$ для всех m интервалов.

При этом $\sum_{i=1}^m n_i = n$ и $\sum_{i=1}^m p_i = 1$.

4. Проводим расчет статистики критерия согласия χ^2 Пирсона с помощью соотношения

$$\chi_{\text{exp}}^2 = n \sum_{i=1}^m \frac{(n_i/n - p_i)^2}{p_i} \quad (6.2)$$

При проверке гипотезы при $n \rightarrow \infty$ для которой известны, как вид закона $p(x, \theta)$, так и все его параметры θ (простая гипотеза) функция χ_{exp}^2 подчиняется распределению χ_r^2 с $r = m - 1$ степенями свободы (доказано в Pearson, Karl (1900). On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling. Philosophical Magazine Series 5 50 (302): 157—175.).

5. Далее находим из уравнения величину

$$\int_{\chi_{\text{exp}}^2}^{\infty} p(s) ds = \int_{\chi_{\text{exp}}^2}^{\infty} \frac{(s)^{\frac{r}{2}-1} e^{-\frac{s}{2}}}{2^{\frac{r}{2}} \Gamma(\frac{r}{2})} ds = p, \text{ или} \\ 1 - F_r(\chi_{\text{exp}}^2) = p, \quad (6.3)$$

где $F_r(x)$ – интегральная функция вероятности распределения χ_r^2 с r степенями свободы.

6. После вычисления p получаем ответ: **Данная выборка с вероятностью p соответствует теоретическому распределению $p(x, \theta)$ с параметрами θ .**

Или: **Данная выборка с вероятностью $\alpha = 1 - p$ не соответствует теоретическому распределению $p(x, \theta)$ с параметрами θ**

Примечания.

- ★ На практике (например с помощью точечных оценок) удается оценить (рассчитать) все или часть параметров распределения. Тогда статистика χ_{exp}^2 при справедливости проверяемой гипотезы подчиняется χ_r^2 -распределению с $r = m - k - 1$ степенями свободы, где k количество оцененных по выборке параметров.

- ★ Некорректное использование критериев согласия (не построен вариационный ряд для выборки, неверно выбрано число интервалов) может приводить к необоснованному принятию (чаще всего) или необоснованному отклонению проверяемой гипотезы.
- ★ Существуют и другие критерия согласия : критерий Колмогорова-Смирнова, критерий Мизеса и т.д.

Существует несколько более удобный способ понимания критерия χ^2 . Введем понятие приведенного значения $\tilde{\chi}^2$ (или $\tilde{\chi}^2$ на одну степень свободы), которое определим как

$$\tilde{\chi}^2 = \frac{\chi^2}{r} \quad (6.4)$$

Тогда, каким бы ни было число степеней свободы, наш критерий можно сформулировать следующим образом: если мы получаем значение $\tilde{\chi}^2$ порядка 1 или меньше, то у нас нет оснований сомневаться в нашем ожидаемом распределении; если мы получаем значение $\tilde{\chi}^2$ много большее, чем единица, то невероятно, чтобы наше ожидаемое распределение было верным. Т.е. если

$$\tilde{\chi}^2 \leq 1, \quad (6.5)$$

то наша выборка соответствует теоретическому распределению.

6.2 Порядок выполнения работы. Задания

Для своего варианта:

Задание 6.1

С помощью критерия Пирсона χ^2 определить является ли случайная физическая величина нормальным распределением, т. е. определить с какой вероятностью данная выборка будет нормальным распределением с соответствующими параметрами.

Задание 6.2

Проделать задание 6.1 с помощью критерия Колмогорова-Смирнова.

Задание 6.3

Построить графики и гистограммы. Сделать визуальное сравнение выбранной модели и экспериментальной кривой. Сделать выводы.

Репозиторий ГГУ им.Ф.Скорины

Лабораторная работа № 7

Фитирование экспериментальных данных

Цель работы: Овладеть методами линейного и нелинейного регрессионного анализа в Wolfram Mathematica

Краткие теоретические сведения

Обобщенная электрическая поляризуемость π^\pm -мезонов $\bar{\alpha}$ была определена из экспериментов по комптоновскому рассеянию и в процессах фотон-фотонных столкновений. Экспериментальные данные представлены в таблице 7.1.

Таблица 7.1— Экспериментальные данные для комптоновской электрической поляризуемости π^\pm -мезонов

Эксперименты	$\bar{\alpha}_{\pi^\pm}/10^{-4}\Phi_M^3$
$\pi^- Z \rightarrow \gamma \pi^- Z$, Серпухов (1983)	$6,8 \pm 1,4 \pm 1,2$
$\gamma p \rightarrow \gamma \pi^+ n$, Физ.Ин-тут им.Лебедева (1984)	20 ± 12
$\gamma\gamma \rightarrow \pi^+\pi^-$: PLUTO (1984)	$19,1 \pm 4,8 \pm 5,7$
DM 1 (1986)	$17,2 \pm 4,6$
DM 2 (1986)	$26,3 \pm 7,4$
Mark II (1990)	$2,2 \pm 1,6$

7.1 Средневзвешенное значение

Средневзвешенное значение \hat{x}_{sw} и ее среднеквадратичное отклонение $\Delta\hat{x}_{sw}$ для набора независимых значений

$$\{x_i, \Delta x_i\}, i = 1, \dots, n, \quad (7.1)$$

где Δx_i – среднеквадратичные отклонения x_i вычисляется с помощью формул:

$$\hat{x}_{sw} = \frac{1}{w} \sum_{i=1}^n w_i x_i, \quad (7.2)$$

$$\Delta \hat{x}_{sw} = \frac{1}{\sqrt{w}}, \quad (7.3)$$

$$w = \sum_{i=1}^n w_i, \quad w_i = \frac{1}{(\Delta x_i)^2}. \quad (7.4)$$

7.2 Поиск параметров регрессионной модели

Для нахождения параметров модели, которые наиболее лучше описывают экспериментальные данные используется оператор `NonlinearModelFit` вида:

`NonlinearModelFit[data, model, pars, vars]`,

где

1. блок `data` определяется “экспериментальными” данными. Например, набор `data={{0,1},{1,0},{3,2},{5,4},{6,4},{7,5}}` формирует список данных вида $\{x_i, y_i\}$ для двумерной случайной величины.
2. блок `model` содержит функциональную зависимость выбранной для фитирования модели вида $y = f(vars, pars)$. Например,

$$\text{Log}[a + bx^2]. \quad (7.5)$$

3. блок `pars` содержит набор параметров для фитирования вида $\{a, b, \dots\}$. Для модели (7.5) это набор из двух параметров $\{a, b\}$.
4. В блоке `vars` задаются переменные от которых зависит модель `model`. Для модели (7.5) это переменная x .

Оператор `NonlinearModelFit` возвращает символьный объект `FittedModel`, содержащий формулы для нелинейной модели, которую он создает. Свойства и диагностику модели можно получить из `model["property"]`.

Оператор `NonlinearModelFit` содержит набор опций, позволяющих влиять на процесс получения оптимальных параметров модели. Информацию об их количестве можно найти с помощью команды `Options[NonlinearModelFit]`. Результатом работы будет список опций с установленными по умолчанию значениями:

```
{AccuracyGoal->Automatic, ConfidenceLevel->19/20,  
EvaluationMonitor->None, Gradient->Automatic,  
MaxIterations->Automatic,  
Method->Automatic, PrecisionGoal->Automatic,  
StepMonitor->None, Tolerance->Automatic,  
VarianceEstimatorFunction->Automatic, Weights->Automatic,  
WorkingPrecision->Automatic}.
```

При наличии ошибок $\Delta y_i, i = 1, \dots, n$ для значений $y_i, i = 1, \dots, n$ величины y в нашем примере для эффективного поиска фитируемых параметров a, b используйте опцию `Weights` вида: `Weights ->{1/(\Delta y_1)^2, 1/(\Delta y_2)^2, \dots, 1/(\Delta y_n)^2}`. Данный расчет рекомендуется дополнить опцией `VarianceEstimatorFunction->(1&)`.

Большое множество примеров применения оператора `NonlinearModelFit` и других операторов для построения эффективных моделей, описывающих различные многомерные случайные величины можно найти в `Documentation Center` системы `Wolfram Mathematica`.

7.3 Порядок выполнения работы. Задания

Задание 7.1

Из данных таблицы 7.1 получить средневзвешенное значение и среднее арифметическое значение электрической поляризуемости пиона и их ошибки. Полученные результаты сравнить и сделать выводы.

Задание 7.2

Среднеквадратичные электромагнитные радиусы заряженных пионов $\langle r^2 \rangle$, можно определить, используя экспериментальные значения форм факторов $F(t)$ в зависимости от квадрата переданного импульса t .

В файле FormFactorPion2.xls представлены экспериментальные значения квадрата форм фактора пиона (безразмерная величина) $|F(t)|^2$ (вторая колонка) в зависимости от квадрата переданного импульса t , выраженного в ГэВ^2 (первая колонка). В третьей колонке представлены экспериментальные абсолютные ошибки $|\Delta F(t)|^2$ измерения квадрата форм фактора пиона.

На основе экспериментальных данных определить значение среднеквадратичного радиуса заряженного пиона $\langle r^2 \rangle$ и его ошибку используя три модели для поведения форм фактора в зависимости от квадрата переданного импульса t :

1. Линейная модель

$$F(t) = 1 - \langle r^2 \rangle \frac{t}{6}; \quad (7.6)$$

2. Полосная модель

$$F(t) = \frac{1}{1 + \langle r^2 \rangle \frac{t}{6}}; \quad (7.7)$$

3. Дипольная модель

$$F(t) = \frac{1}{\left(1 + \langle r^2 \rangle \frac{t}{12}\right)^2}. \quad (7.8)$$

Найти значения $\langle r^2 \rangle$ в Ферми² $\equiv \Phi_{\text{м}}^2$, используя соотношение

$$0,1973269631 \Phi_{\text{м}} \text{ ГэВ} = 1. \quad (7.9)$$

Провести графический сравнительный анализ экспериментальных данных $F(t)$ и трех моделей с использованием полученного значения $\langle r^2 \rangle$. Выбрать наиболее оптимальную модель.

Примечание. Для построения в Wolfram Mathematica экспериментальных значений с ошибками удобно использовать пакет ErrorBarPlots. Это дополнительный пакет, который не загружается автоматически. Для его загрузки необходима команда Needs, т. е.

Needs [ErrorBarPlots[""];"

Следующий шаг состоит в импорте экспериментальных данных из фай-

ла FormFactorPion2.xls. Программный блок, позволяющий импортировать отдельные колонки из этого файла отображен на рисунке 7.1 (не забудьте изменить путь файлу. Он находится вместе с данными.).

```
Needs["ErrorBarPlots`"];  
|необходимо  
SetDirectory["d:\\Users\\Andreev\\PREDMET\\C_K\\СтатОбработка\\Компьютерная физика\\Lab6\\"];  
|задать рабочую директорию  
datagpn = Import["FormFactorPion2.xls"][[1]];  
|импорт  
(* Извлечь 1 столбик и удалить первый элемент *)  
t2 = Delete[datagpn[[All, 1]], 1];  
|удалить элемент |всё  
(* Извлечь 2 столбик и удалить первый элемент *)  
f2 = Delete[datagpn[[All, 2]], 1];  
|удалить элемент |всё  
(* Извлечь 3 столбик и удалить первый элемент *)  
ErFp = Delete[datagpn[[All, 3]], 1];  
|удалить элемент |всё
```

Рисунок 7.1– Блок программы с импортом данных из таблицы

Репозиторий ГГУ им. Ф. Скорины

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

1. Baumann, G. *Mathematica for Theoretical Physics: Classical Mechanics and Nonlinear Dynamics*/ G. Baumann. – Springer Science+Business Media, Inc., 2006. – 544 P.
2. Baumann, G. *Mathematica for Theoretical Physics: Electrodynamics, Quantum Mechanics, General Relativity and Fractals*/ G. Baumann. – Springer Science+Business Media, Inc., 2005. – 544 P.
3. Boccara, N. *Essentials of Mathematica: With Applications to Mathematics and Physics*/ N. Boccara. – Springer Science+Business Media, LLC, 2007. – 540 P.
4. Kythe, P. K. *Partial differential equations and Mathematica*/ P. K. Kythe, P. Puri, M. R. Schäferkötter. – CRC Press, Inc., 1997. – 398 P.
5. Thaller, B. *Visual Quantum Mechanics: Selected Topics with Computer-Generated Animations of Quantum-Mechanical Phenomena Includes electronic media*/ B. Thaller. – TELOS/Springer-Verlag, 2000.
6. Trott, M. *The Mathematica GuideBook for Numerics*/ M. Trott. – Springer, 2005. – P. 1208.
7. Trott, M. *The Mathematica guidebook for symbolics*/ M. Trott. – Springer, 2005. – P. 1208.
8. Trott, M. *The Mathematica Guidebook: Programming*/ M. Trott. – Springer, 2002. – P. 905.
9. Wellin, P. R. *An Introduction to Programming with Mathematica*/ P. R. Wellin, R. Gaylord, S. Kamin. – Cambridge University Press, 2005. – 570 P.
10. Zimmerman, R. L. *Mathematica for Physics*/ R. L. Zimmerman, F. I. Olness. – Second edition edition. – Addison-Wesley, 2002. – 646 P.
11. Воробьёв, Е. *Введение в систему Mathematica*/ Е. Воробьёв. – Москва: Финансы и статистика, 1998. – 345 с.
12. Половко, А. М. *Matematica для студентов*/ А. М. Половко. – Санкт-Петербург: БХВ-Петербург, 2007. – 368 с.
13. Wolfram, S. *The Mathematica book*/ S. Wolfram. – 4th edition. – Addison-Wesley, 1999.
14. Тейлор, Д. *Введение в теорию ошибок*/ Д. Тейлор. – Москва: Мир, 1985. – С. 272.
15. Новицкий, П. В. *Оценка погрешностей результатов измерений*/ П. В. Новицкий, И. А. Зограф. – Ленинград: Энергоатомиздат, 1985. – С. 248.

16. Львовский, Е. Н. Статистические методы построения эмпирических формул: Учеб. пособие для вузов./ Е. Н. Львовский. – 2-е изд., перераб. и доп. изд. . – Москва: Высшая школа, 1988. – С. 239.
17. Лавренчик, В. Н. Постановка физического эксперимента и статистическая обработка его результатов/ В. Н. Лавренчик. – Москва: Энергоатомиздат, 1986. – С. 272.
18. Кремер, Н. Ш. Теория вероятностей и математическая статистика: Учебник для вузов./ Н. Ш. Кремер. – 2-е изд., перераб. и доп. изд. . – Москва: ЮНИТИ- ДАНА, 2004. – С. 573.
19. Кобзарь, А. И. Прикладная математическая статистика. Для инженеров и научных работников./ А. И. Кобзарь. – Москва: - М.: ФИЗМАТЛИТ, 2006. - 816 с. -, 2006. – С. 816. – ISBN 5-9221-0707-0.
20. Статистические методы в экспериментальной физике/ В. Идье, Д. Драйард, Ф. Джеймс , [и др.] под ред. А. А. Тяпкина. – Москва: Атомиздат, 1976. – С. 335.
21. Гмурман, В. Е. Теория вероятностей и математическая статистика: Учеб. пособие для вузов/ В. Е. Гмурман. – 9-е издание изд. . – Москва: Высшая школа, 2003. – С. 479.
22. Бендат, Д. Прикладной анализ случайных данных/ Д. Бендат, А. Пирсол. – Москва: Мир, 1989. – С. 540.
23. Taylor, J. An introduction to error analysis/ J. Taylor. – Second edition. – University Science Books, 1996. – P. 327.
24. Soong, T. Fundamentals of Probability and Statistics for Engineers/ T. Soong. – John Wiley & Sons, Ltd, 2004. – P. 406. – ISBN: 9780470868157.
25. Lyons, L. A practical guide to data analysis for physical science students/ L. Lyons. – Cambridge University Press, 1991. – P. 95.
26. Statistical methods in experimental physics/ W. T. Eadie, D. Dryard, F. E. James [et al.]. – Amsterdam , London: North–Holland; Publishing Company, 1971.

СТАТИСТИЧЕСКИЕ МЕТОДЫ ОБРАБОТКИ ДАННЫХ

учебная дисциплина для специальности:
1-31 04 08 Компьютерная физика

1. Частота попадания случайной величины.
2. Интегральный закон распределения вероятности.
3. Дифференциальный закон распределения вероятности.
4. Примеры наиболее распространенных в физике функций распределения.
5. Моменты распределения вероятностей.
6. Математическое ожидание. Геометрическая интерпретация математического ожидания.
7. Медиана. Мода распределения.
8. Дисперсия. Разброс распределения.
9. Полуширота распределения. Геометрическая интерпретация дисперсии.
10. Третий центральный момент. Геометрическая интерпретация третьего центрального момента.
11. Четвертый центральный момент. Геометрическая интерпретация четвертого центрального момента. Центральные моменты более высоких порядков.
12. Типы оценок. Точечные оценки.
13. Среднее арифметическое и математическое ожидание. Среднее квадратичное отклонение.
14. Точечная оценка третьего центрального момента. Коэффициент асимметрии.
15. Точечная оценка четвертого центрального момента. Эксцесс. Контрэксцесс. С
16. Лучший характер точечных оценок. Разброс точечных оценок.
17. Интервальные оценки. Доверительные интервалы для математического ожидания и дисперсии.
18. Квантили распределений случайных чисел. Выбор доверительной вероятности для доверительного интервала.
19. Графическое отображение законов распределений. Полигон частот. Гистограмма и ее построение.
20. Необходимость гистограмм в обработке физического эксперимента. Оптимальное число интервалов в гистограмме.
21. Критерии выбора длины интервала и числа столбцов. Форма распределения вероятности. Зависимость оптимального числа столбцов гистограммы от формы распределения физической величины.
22. Практические рекомендации для построения гистограммы. Предварительная обработка данных.

23. Цензурирование выборки. Критерии для нахождения возможных промахов. Сравнение гистограмм.
24. Критерии согласия. Критерий Колмогорова-Смирнова. Условия использования критерия Колмогорова-Смирнова.
25. Критерий согласия Пирсона (хи-квадрат критерий). Условия использования критерия Пирсона.
26. Выбор объема выборки случайной величины для критерия Пирсона.
27. Многомерные случайные величины и статистические задачи, возникающие при их обработке.
28. Двумерные случайные величины. Некоторые сведения из теории вероятностей для двумерных случайных функций.
29. Коэффициент корреляции. Статистически независимые двумерные случайные величины.
30. Анализ двумерных случайных величин. Задача корреляционного анализа. Линейный корреляционный анализ.
31. Задачи регрессионного анализа. Линейный регрессионный анализ.
32. Метод наименьших квадратов.
33. Нелинейный корреляционный анализ. Нелинейный регрессионный анализ. Задачи, сводимые к линейным задачам корреляционного и регрессионного анализа.
34. Работа с файлами в системе MATHEMATICA. Создание файлов физических данных системе MATHEMATICA.
35. Средства графического отображения функций в системе MATHEMATICA. Построение двумерных и трехмерных графиков.
36. Встроенные операторы для расчета моментов распределений. Расчет медианы различных распределений.
37. Расчет среднего арифметического и дисперсии. Вычисление точечной оценки третьего центрального момента. Встроенная функция для коэффициента асимметрии.
38. Вычисление точечной оценки четвертого центрального момента. Вычисление моментов высоких порядков в системе MATHEMATICA.
39. Встроенные операторы для расчета доверительных интервалов.
40. Встроенные функции системы MATHEMATICA для построения графиков функций распределений случайных величин.
41. Генерация псевдослучайных чисел в системе MATHEMATICA. Расчет квантилей и интегральных функций распределений вероятностей. Способы построения гистограмм. Типы гистограмм.
42. Встроенные функции системы MATHEMATICA для построения гистограмм. Фурье-анализ в системе MATHEMATICA.
43. Встроенные функции системы MATHEMATICA для построения фурье-образов.
44. Линейный корреляционный анализ в системе MATHEMATICA.
45. Элементы регрессионного анализа в пакете MATHEMATICA.
46. Обработка данных с помощью критериев согласия с использованием системы MATHEMATICA.

Учреждение образования
«Гомельский государственный университет имени Франциска Скорины»

УТВЕРЖДАЮ

Проректор по учебной работе
УО «ГГУ им. Ф. Скорины»

_____ И.В. Семченко
(подпись)

(дата утверждения)
Регистрационный № УД-_____

**СТАТИСТИЧЕСКИЕ МЕТОДЫ
ОБРАБОТКИ ДАННЫХ**

Учебная программа учреждения высшего образования
по учебной дисциплине для специальности:

1-31 04 08 Компьютерная физика

Репозиторий ГГУ им. Ф. Скорины

2017 г

Учебная программа составлена на основе образовательного стандарта высшего образования ОСВО 1-31 04 08-2013 Компьютерная физика и учебного плана учреждения высшего образования «ГГУ им. Ф. Скорины», регистрационный № G 31-01-16, дата утверждения 12.01.2016

СОСТАВИТЕЛЬ:

В.В. АНДРЕЕВ, заведующий кафедрой теоретической физики, доктор физико-математических наук, доцент

РЕЦЕНЗЕНТЫ:

А.А. БАБИЧ, заведующий кафедрой «Высшая математика» УО «ГГТУ им. П.О. Сухого», кандидат физико-математических наук, доцент;
Е.Б. ШЕРШНЕВ, заведующий кафедрой общей физики УО «ГГУ им. Ф.Скорины», кандидат технических наук, доцент

РЕКОМЕНДОВАНА К УТВЕРЖДЕНИЮ:

Кафедрой теоретической физики учреждения образования «Гомельский государственный университет им. Ф. Скорины» (протокол № 10 от 23.05.2017);

Научно-методическим советом учреждения образования «Гомельский государственный университет им. Ф. Скорины» (протокол № 8 от 07.06.2017)

ПОЯСНИТЕЛЬНАЯ ЗАПИСКА

Современный инженер должен быть подготовлен к организационно-управленческой, научной, производственной и преподавательской деятельности в области экспериментального исследования физических процессов на различных уровнях структурной организации материи при различных физических условиях; теоретического анализа эффектов и явлений и предсказания новых физических закономерностей на основе современных теоретических представлений, математических и компьютерных методов; разработке приборов на основе новых материалов и физических принципов, созданию новых технологий, использующих математические методы и компьютерную технику; работе по математическому моделированию разнообразных процессов и объектов.

Одним из обязательных требований такого рода специалистов является их способность планировать, организовывать и проводить физические и радиофизические эксперименты, применяя вероятностные методы анализа и используя технические средства автоматизации эксперимента; обрабатывать и анализировать полученные результаты, создавать математические модели и программные средства, составлять отчеты и вести научно-техническую документацию.

Как известно статистическая обработка информации в физических исследованиях играет одну из важнейших ролей. Новые неизвестные явления открываются в настоящее время путем тщательной обработки результатов измерений. Современный физик обязан знать основы статистической обработки данных и с помощью программных продуктов уметь проводить анализ и обработку физической информации на персональном компьютере.

Исходя из этих требований, следует что, такая учебная дисциплина, как «Статистические методы обработки данных» является необходимой составной частью образования современного инженера-физика.

Целью учебной дисциплины «Статистические методы обработки данных» является овладение студентами основами математической статистики и теории вероятностей и их использование при обработке и анализе экспериментальных данных.

Задачами учебной дисциплины «Статистические методы обработки данных» являются

- приобретение навыков и изучение способов, приёмов работы по анализу и статистической обработке информации на персональном компьютере, наиболее часто используемых в физических исследованиях (обратив внимание на смысл понятий критериев, возможности алгоритмизации),

- приобретение навыков по решению теоретических и экспериментальных задач различных физических дисциплин, таких как квантовая механика, физика ядра и атома, радиационная безопасность и другие.

- овладение новыми методами теоретических и экспериментальных исследований (решение практических задач, методы получения и обработки результатов опытов).

Материал учебной дисциплины «Статистические методы обработки данных» базируется на ранее полученных знаниях по курсу «Теория вероятностей». Полученные навыки могут быть использованы при изучении дисциплин учебного плана, таких как «Физика ядра», «Квантовая механика» и специального курса «Метод Монте-Карло в физике элементарных частиц».

Учебная дисциплина «Статистические методы обработки данных» формирует необходимую базу для научных исследований студентов и выполнения курсовых и дипломных работ.

В результате изучения учебной дисциплины «Статистические методы обработки данных» студент должен

знать:

- способы, приёмы работы для процесса обработки физической информации.

уметь:

- применять системы аналитических вычислений для решения задач квантовой механики и квантовой теории поля
- использовать технические средства автоматизации эксперимента
- обрабатывать и анализировать полученные результаты, создавать математические модели и программные средства

иметь опыт:

- планирования и организации научного исследования, применения соответствующих экспериментальных и теоретических методов
- интерпретации результатов экспериментальных исследований.

Освоение данной образовательной программы должно обеспечить формирование следующих компетенций:

АК-1. Уметь применять базовые научно-теоретические знания для решения теоретических и практических задач.

АК-2. Владеть системным и сравнительным анализом.

АК-3. Владеть исследовательскими навыками.

АК-4. Уметь работать самостоятельно.

АК-7. Иметь навыки, связанные с использованием технических устройств, управлением информацией и работой с компьютером.

АК-8. Иметь лингвистические навыки (устная и письменная коммуникация).

АК-9. Уметь учиться, повышать свою квалификацию в течение всей жизни

СЛК-1. Обладать качествами гражданственности.

СЛК-2. Быть способным к социальному взаимодействию.

СЛК-3. Обладать способностью к межличностным коммуникациям.

СЛК-5. Быть способным к критике и самокритике (критическое мышление).

СЛК-6. Уметь работать в команде.

ПК-1. Применять знания теоретических и экспериментальных основ физики и математики, методов измерения физических величин, методов планирования, организации и ведения научно-производственной, научно-педагогической, производственно-технической, опытно-конструкторской работы, средств автоматизации, правового обеспечения хозяйственной деятельности и налоговой системы, государственного регулирования экономики и экономической политики.

ПК-2. Владеть современными методами программирования, компьютерными методами сбора, хранения и обработки информации, системами автоматизированного проектирования.

ПК-3. Оценивать конкурентоспособность и экономическую эффективность разрабатываемого программного обеспечения.

ПК-4. Пользоваться глобальными информационными ресурсами, новой научной, технической и патентной литературой по физике, математике, информатике, экономике и инновационным технологиям, основами психолого-педагогических знаний, навыками самообразования и самосовершенствования.

ПК-6. Применять полученные знания фундаментальных положений физики, экспериментальных, теоретических и компьютерных методов исследования, планирования, организации и ведения научно-технической и научно-педагогической работы.

ПК-7. Использовать новейшие открытия в естествознании, методы научного анализа, информационно-образовательные технологии, физические основы современных технологических процессов.

ПК-8. Пользоваться государственными языками Республики Беларусь и иными иностранными языками как средством делового общения.

ПК-9. Реализовывать методы защиты производственного персонала и населения в условиях возникновения аварий, катастроф, стихийных бедствий и обеспечения радиационной безопасности при осуществлении научной, производственной и педагогической деятельности.

ПК-11. Осуществлять поиск, систематизацию и анализ информации по перспективным направлениям развития отрасли, инновационным технологиям, проектам и решениям.

ПК-12. Определять цели инноваций и способы их достижения.

ПК-13. Применять методы анализа и внедрения инноваций в научно-производственной, научно-педагогической и научно-технической деятельности.

Учебная программа дисциплины «Статистические методы обработки данных» рекомендуется для подготовки специалистов по специальности 1-31 04 08 Компьютерная физика.

Дисциплина изучается на 2 курсе в 3 семестре.

Форма обучения – дневная.

Общее количество часов – 96; аудиторное количество часов – 64.

Распределение аудиторного времени по видам занятий: лекции – 24 (из них УСР-6), лабораторные занятия – 34.

Форма отчётности – зачет в 3 семестре.

Репозиторий ГГУ им. Ф. Скорины

СОДЕРЖАНИЕ УЧЕБНОГО МАТЕРИАЛА

Раздел 1. Основные понятия теории вероятности и математической статистики

Тема 1. Вероятность и характеристики распределения вероятностей

Частота попадания случайной величины. Интегральный закон распределения вероятности. Дифференциальный закон распределения вероятности. Примеры наиболее распространенных в физике функций распределения.

Моменты распределения вероятностей. Математическое ожидание. Геометрическая интерпретация математического ожидания. Медиана. Мода распределения. Дисперсия. Разброс распределения. Полуширота распределения. Геометрическая интерпретация дисперсии. Третий центральный момент. Геометрическая интерпретация третьего центрального момента. Четвертый центральный момент. Геометрическая интерпретация четвертого центрального момента. Центральные моменты более высоких порядков.

Тема 2. Основы теории оценок

Типы оценок. Точечные оценки. Среднее арифметическое и математическое ожидание. Среднее квадратичное отклонение. Точечная оценка третьего центрального момента. Коэффициент асимметрии. Точечная оценка четвертого центрального момента. Эксцесс. Контрэксцесс. Случайный характер точечных оценок. Разброс точечных оценок. Интервальные оценки. Доверительные интервалы для математического ожидания и дисперсии. Квантили распределений случайных чисел. Выбор доверительной вероятности для доверительного интервала.

Тема 3. Гистограммы и операции над ними

Графическое отображение законов распределений. Полигон частот. Гистограмма и ее построение. Необходимость гистограмм в обработке физического эксперимента. Оптимальное число интервалов в гистограмме. Критерии выбора длины интервала и числа столбцов. Форма распределения вероятности. Зависимость оптимального числа столбцов гистограммы от формы распределения физической величины. Практические рекомендации для построения гистограммы. Предварительная обработка данных. Цензурирование выборки. Критерии для нахождения возможных промахов. Сравнение гистограмм. Фурье-преобразование гистограмм. Быстрое преобразование Фурье. Примеры использования преобразование Фурье в технике.

Тема 4. Идентификация формы распределения экспериментальных данных

Критерии согласия. Критерий Колмогорова-Смирнова. Условия использования критерия Колмогорова-Смирнова. Критерий согласия Пирсона (хи-квадрат критерий). Условия использования критерия Пирсона. Использование гистограмм для критерия Пирсона. Критерий согласия Мизеса. Критерий согласия Мизеса для дискретных распределений. Выбор объема выборки случайной величины для критерия Пирсона. Выбор объема выборки случайной величины для критерия Колмогорова-Смирнова. Примеры использования в физике критериев согласия.

Тема 5. Регрессионный и корреляционный анализ

Многомерные случайные величины и статистические задачи, возникающие при их обработке. Двумерные случайные величины. Некоторые сведения из теории вероятностей для двумерных случайных функциях. Коэффициент корреляции. Статистически независимые двумерные случайные величины. Анализ двумерных случайных величин. Задача корреляционного анализа. Линейный корреляционный анализ. Задачи регрессионного анализа. Линейный регрессионный анализ. Метод наименьших квадратов. Нелинейный корреляционный анализ. Нелинейный регрессионный анализ. Задачи, сводимые к линейным задачам корреляционного и регрессионного анализа.

Раздел 2. Обработка результатов измерения на персональном компьютере

Тема 6. Точечные и интервальные оценки распределений в системе МАТНЕМАТІСА

Работа с файлами в системе МАТНЕМАТІСА. Создание файлов физических данных в системе МАТНЕМАТІСА. Средства графического отображения функций в системе МАТНЕМАТІСА. Построение двумерных и трехмерных графиков. Встроенные операторы для расчета моментов распределений. Расчет медианы различных распределений. Расчет среднего арифметического и дисперсии. Вычисление точечной оценки третьего центрального момента. Встроенная функция для коэффициента асимметрии. Вычисление точечной оценки четвертого центрального момента. Вычисление моментов высоких порядков в системе МАТНЕМАТІСА. Встроенные операторы для расчета доверительных интервалов.

Тема 7. Функции распределений случайных величин и средства статистического анализа в системе МАТНЕМАТІСА

Встроенные функции системы МАТНЕМАТІСА для построения графиков функций распределений случайных величин. Генерация псевдослучайных чисел в системе МАТНЕМАТІСА. Расчет квантилей и интегральных функций распределений вероятностей. Способы построения гистограмм. Типы гистограмм. Встроенные функции системы МАТНЕМАТІСА для построения гистограмм. Фурье-анализ в системе МАТНЕМАТІСА. Встроенные функции системы МАТНЕМАТІСА для построения фурье-образов. Линейный корреляционный анализ в системе МАТНЕМАТІСА. Элементы регрессионного анализа в пакете МАТНЕМАТІСА. Обработка данных с помощью критериев согласия с использованием системы МАТНЕМАТІСА.

УЧЕБНО-МЕТОДИЧЕСКАЯ КАРТА ДИСЦИПЛИНЫ

Номер раздела, темы, занятия	Название раздела, темы, занятия; перечень изучаемых вопросов	Количество аудиторных часов				Материальное обеспечение занятия (наглядные, методические пособия и др.)	Литература	Формы контроля знаний
		лекции	практические (семинарские) занятия	лабораторные занятия	управляемая самостоятельная работа студента			
1	2	3	4	5	6	7	8	9
1	Основные понятия теории вероятности и математической статистики	16		18	4			
1.1	Вероятность и характеристики распределения вероятностей. 1. Частота попадания случайной величины. 2. Интегральный и дифференциальный законы распределений вероятностей. 3. Примеры наиболее распространенных в физике функций распределения. 4. Математическое ожидание и дисперсия. 5. Центральные моменты высоких порядков и их смысл.	4		2		Мультимедийный проектор, презентация	[1] [2] [3]	
1.2	Основы теории оценок 1. Типы оценок. 2. Точечные оценки и их случайный характер. 3. Разброс точечных оценок. 4. Интервальные оценки. 5. Доверительные интервалы для математического ожидания и дисперсии. Выбор доверительной вероятности для доверительного интервала.	2		6		Мультимедийный проектор	[1] [2] [3]	
1.3	Гистограммы и операции над ними 1. Полигон частот. 2. Гистограмма и ее построение. Необходимость гистограмм в обработке физического эксперимента. 3. Критерии выбора длины интервала и числа столбцов. 4. Зависимость оптимального числа столбцов гистограммы от формы распределения физической величины.	4		4		Мультимедийный проектор	[1] [3] [4] [6]	

	5. Практические рекомендации для построения гистограммы.							
1.4	Идентификация формы распределения экспериментальных данных 1. Критерий Колмогорова-Смирнова. 2. Критерий согласия Пирсона. 3. Критерий согласия Мизеса.	2		6	2	Мультимедийный проектор	[1] [2] [4] [9]	
1.5	Регрессионный и корреляционный анализ Критерий 1. Многомерные случайные величины и статистические задачи, возникающие при их обработке. 2. Двумерные случайные величины. 3. Коэффициент корреляции. Статистически независимые двумерные случайные величины. Анализ двумерных случайных величин. 4. Задача корреляционного анализа. Линейный корреляционный анализ. 5. Задачи регрессионного анализа. Линейный регрессионный анализ. 6. Метод наименьших квадратов. Нелинейный корреляционный анализ.	4		6	2	Мультимедийный проектор	[1] [2] [4] [6]	
2	Обработка результатов измерения на персональном компьютере	8		16	2			
2.1.	Точечные оценки распределений в системе МАТНЕМАТІСА 1. Встроенные операторы для расчета моментов распределений. 2. Встроенные операторы для расчета доверительных интервалов. 3. Работа с файлами в в системе МАТНЕМАТІСА.	4		8		Мультимедийный проектор	[1] [3] [7] [8]	
2.2	Функции распределений случайных величин в системе МАТНЕМАТІСА 1. Встроенные функции системы МАТНЕМАТІСА для построения функций распределений случайных величин, 2. Встроенные функции системы МАТНЕМАТІСА для построения гистограмм. 3. Встроенные функции системы МАТНЕМАТІСА для построения фурье-образов.	4		8	2	Мультимедийный проектор	[1] [3] [7] [8]	
	Итого	24		34	6			зачет

ИНФОРМАЦИОННО-МЕТОДИЧЕСКАЯ ЧАСТЬ

Перечень лабораторных работ

1. Начальные навыки работы в системе Wolfram Mathematica.
2. Построение графиков различных распределений.
3. Точечные и интервальные оценки случайной величины.
4. Построение гистограмм в системе Mathematica.
5. Промахи и методы их исключения.
6. Идентификация форм распределения экспериментальных данных.
7. Линейный корреляционный анализ.
8. Линейный регрессионный анализ.
9. Фильтрация экспериментальных данных.
10. Встроенные операторы математической статистики в системе Wolfram Mathematica
11. Преобразование Фурье.
12. Встроенные функции системы Wolfram Mathematica для построения функций распределений случайных величин, гистограмм и фурье-образов.

Темы для управляемой самостоятельной работы

1. Критерий Колмогорова-Смирнова.
2. Критерий согласия Мизеса.
3. Метод наименьших квадратов.
4. Нелинейный корреляционный анализ.
5. Встроенные функции системы Wolfram Mathematica для построения фурье-образов.

Рекомендуемая литература

Основная

1. Новицкий, П.В. Оценка погрешностей результатов измерений/ П.В.Новицкий, И.А..Зограф — Л.: Энергоатомиздат, 1991. — 304 с.
2. Бендат Дж., Пирсол А. Прикладной анализ случайных данных/ Дж.Бендат, А.Пирсол. — М.: Мир, 1989. — 504 с.
3. Лавренчик, В.Н. Постановка физического эксперимента и статистическая обработка его результатов/ В.Н. Лавренчик.. — М.: Энергоатомиздат, 1986. — 272 с.
4. Тейлор Дж. Введение в теорию ошибок/ Дж.Тейлор. — М.: Мир, 1985. — 45 с.
5. Кембровский, Г.С. Приближенные вычисления и методы обработки результатов измерений в физике/ Г.С.Кембровский. — Мн.: Университетское, 1990. — 67 с.
6. Злоказов В.Б. Математические методы анализа экспериментальных спектров и спектроподобных экспериментов/ В.Б. Злоказов. // ЭЧАЯ, 1985, Т.16, Вып.5, С.1126-1163.
7. Дьяконов, В.П. Mathematica 4: учебный курс/ В.П. Дьяконов. — СПб.: Питер, 2001. — 654 с.
8. Воробьев Е.М. Введение в систему Mathematica./ Е.М.Воробьев. — М.: Финансы и статистика, 1998. — 345 с.
9. Гмурман, В. Е. Теория вероятностей и математическая статистика: Учеб. пособие для вузов/В. Е. Гмурман. — 9-е изд., стер. — М.: Высш. шк., 2003. — 479 с.

Дополнительная

1. Львовский, Б.Н. Статистические методы построения эмпирических формул: Учеб. пособие для втузов/ Б.Н.Львовский. — М.: Высш. шк., 1988— 239 с..
2. Кобзарь А. И. Прикладная математическая статистика. Для инженеров и научных работников/ А. И. Кобзарь. — М.: ФИЗМАТЛИТ, 2006. — 816 с.
3. Боровков Л. Л. Математическая статистика/ Л. Л.Боровков— Учебник.— М.: Наука. Главная редакция физико-математической литературы, 1984.— 472 с.
4. Кассандрова, О. Н. Обработка результатов наблюдений,/ О. Н. Кассандрова, В. В. Лебедев— М.:Наука, Главная редакция физ.-мат. литературы, 1970 г. —104 с.
5. Ивченко, Г. И. Математическая статистика/ Г.И.Ивченко, Ю. И. Медведев Учеб. пособие для втузов. — М.: Высш. шк., 1984. — 248 с.

Репозиторий ГГУ им. Ф.Скоринны

ПРОТОКОЛ СОГЛАСОВАНИЯ УЧЕБНОЙ ПРОГРАММЫ
ПО ИЗУЧАЕМОЙ УЧЕБНОЙ ДИСЦИПЛИНЕ
С ДРУГИМИ ДИСЦИПЛИНАМИ СПЕЦИАЛЬНОСТИ

Название дисциплины, с которой требуется согласование	Название кафедры	Предложения об изменениях в содержании учебной программы по изучаемой учебной дисциплине	Решение, принятое кафедрой, разработавшей учебную программу (с указанием даты и номера протокола)

Репозиторий ГГУ им. Ф. Скорины

ДОПОЛНЕНИЯ И ИЗМЕНЕНИЯ К УЧЕБНОЙ ПРОГРАММЕ
ПО ИЗУЧАЕМОЙ УЧЕБНОЙ ДИСЦИПЛИНЕ
на ____ / ____ учебный год

№№ пп	Дополнения и изменения	Основание

Учебная программа пересмотрена и одобрена на заседании кафедры
теоретической физики
(протокол № ____ от _____ 201 _ г.)

Заведующий кафедрой

теоретической физики
д.ф.-м.н., доцент

_____ В.В. Андреев

УТВЕРЖДАЮ

Декан факультета физики и ИТ УО «ГГУ им. Ф. Скорины»
к.ф.-м.н., доцент

_____ Д.Л. Коваленко

СТАТИСТИЧЕСКИЕ МЕТОДЫ ОБРАБОТКИ ДАННЫХ

учебная дисциплина для специальности:
1-31 04 08 Компьютерная физика

Рекомендуемая литература

Основная

1. Новицкий, П.В. Оценка погрешностей результатов измерений/ П.В.Новицкий, И.А.Зограф — Л.: Энергоатомиздат, 1991. — 304 с.
2. Бендат Дж., Пирсол А. Прикладной анализ случайных данных/ Дж.Бендат, А.Пирсол. —М.: Мир,1989. —504 с.
3. Лавренчик, В.Н. Постановка физического эксперимента и статистическая обработка его результатов/ В.Н. Лавренчик.. — М.: Энергоатомиздат, 1986. — 272 с.
4. Тейлор Дж. Введение в теорию ошибок/ Дж.Тейлор. — М.: Мир, 1985. — 45 с.
5. Дьяконов, В.П. Mathematica 4: учебный курс/ В.П. Дьяконов. — СПб.: Питер, 2001 . — 654 с.
6. Воробьёв Е.М. Введение в систему Mathematica./ Е.М.Воробьёв. — М.: Финансы и статистика, 1998. — 345 с.
7. Гмурман, В. Е. Теория вероятностей и математическая статистика: Учеб. пособие для вузов/В. Е. Гмурман. — 9-е изд., стер. — М.: Высш. шк., 2003. — 479 с.
8. Львовский, Б.Н. Статистические методы построения эмпирических формул: Учеб. пособие для втузов/ Б.Н.Львовский. — М.: Высш. шк., 1988— 239 с.
9. Кобзарь А. И. Прикладная математическая статистика. Для инженеров и научных работников/ А. И. Кобзарь. — М.: ФИЗМАТЛИТ, 2006. — 816 с.

Дополнительная

1. Боровков Л. Л. Математическая статистика/ Л. Л.Боровков— Учебник.— М.: Наука. Главная редакция физико-математической литературы, 1984.— 472 с.
2. Кассандрова, О. Н. Обработка результатов наблюдений,/ О. Н. Кассандрова, В. В. Лебедев— М.:Наука, Главная редакция физ.-мат. литературы, 1970 г. —104 с.
3. Ивченко, Г. И. Математическая статистика/ Г.И.Ивченко, Ю. И. Медведев Учеб. пособие для втузов. — М.: Высш. шк., 1984. — 248 с.

4. Baumann, G. *Mathematica for Theoretical Physics: Classical Mechanics and Nonlinear Dynamics*/ G. Baumann. — Springer Science+Business Media, Inc., 2006. — 544 P.
5. Baumann, G. *Mathematica for Theoretical Physics: Electrodynamics, Quantum Mechanics, General Relativity and Fractals*/ G. Baumann. — Springer Science+Business Media, Inc., 2005. — 544 P.
6. Boccara, N. *Essentials of Mathematica: With Applications to Mathematics and Physics*/ N. Boccara. — Springer Science+Business Media, LLC, 2007. — 540 P.
7. Kythe, P. K. *Partial differential equations and Mathematica*/ P. K. Kythe, P. Puri, M. R. Schaferkotter. — CRC Press, Inc., 1997. — 398 P.
8. Thaller, B. *Visual Quantum Mechanics: Selected Topics with Computer-Generated Animations of Quantum-Mechanical Phenomena Includes electronic media*/ B. Thaller. — TELOS/Springer-Verlag, 2000.
9. Trott, M. *The Mathematica GuideBook for Numerics*/ M. Trott. — Springer, 2005. — P. 1208.
10. Trott, M. *The Mathematica guidebook for symbolics*/ M. Trott. — Springer, 2005. — P. 1209.
11. Trott, M. *The Mathematica Guidebook: Programming*/ M. Trott. — Springer, 2002. — P. 905.
12. Wellin, P. R. *An Introduction to Programming with Mathematica*/ P. R. Wellin, 13. R. Gaylord, S. Kamin. — Cambridge University Press, 2005. — 570 P.
14. Zimmerman, R. L. *Mathematica for Physics*/ R. L. Zimmerman, F. I. Olness. — Second edition. — Addison-Wesley, 2002. — 646 P.
15. Wolfram, S. *The Mathematica book*/ S. Wolfram. — 4th edition. — Addison-Wesley, 1999.
16. Кремер, Н. Ш. *Теория вероятностей и математическая статистика: Учебник для вузов*/ Н. Ш. Кремер. — 2-е изд., перераб. и доп. изд. . — Москва: ЮНИТИ- ДАНА, 2004. — С. 573
17. *Статистические методы в экспериментальной физике*/ В. Идье, Д. Драйард, Ф. Джеймс, [и др.] под ред. А. А. Тяпкина. — Москва: Атомиздат, 1976. — С. 335. 2003. — С. 479.
18. Taylor, J. *An introduction to error analysis*/ J. Taylor. — Second edition. — University Science Books, 1996. — P. 327.
19. Soong, T. *Fundamentals of Probability and Statistics for Engineers*/ T. Soong.— John Wiley & Sons, Ltd, 2004. — P. 406. — ISBN: 9780470868157.
20. Lyons, L. *A practical guide to data analysis for physical science students*/ L. Lyons. — Cambridge University Press, 1991. — P. 95.
21. *Statistical methods in experimental physics*/ W. T. Eadie, D. Dryard, F. E. James [et al.]. — Amsterdam, London: North—Holland; Publishing Company, 1971.