

МИНИСТЕРСТВО НАУКИ И ВЫСШЕГО ОБРАЗОВАНИЯ РОССИЙСКОЙ ФЕДЕРАЦИИ

ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ  
ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ  
«САМАРСКИЙ НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ  
УНИВЕРСИТЕТ ИМЕНИ АКАДЕМИКА С.П. КОРОЛЕВА»  
(САМАРСКИЙ УНИВЕРСИТЕТ)

*А.Г. ХРАМОВ*

# МЕТОДЫ И АЛГОРИТМЫ ИНТЕЛЛЕКТУАЛЬНОГО АНАЛИЗА ДАННЫХ

Рекомендовано редакционно-издательским советом федерального государственного автономного образовательного учреждения высшего образования «Самарский национальный исследовательский университет имени академика С.П. Королева» в качестве учебного пособия для обучающихся по основной образовательной программе высшего образования по направлению подготовки 01.04.02 Прикладная математика и информатика

САМАРА  
Издательство Самарского университета  
2019

УДК 004.8(075)  
ББК 32.973.2я7  
X895

Рецензенты: зам. директора ИПУСС РАН по научной работе  
С.В. Смирнов;  
д-р техн. наук, проф. В.А. Фурсов

*Храмов, Александр Григорьевич*

**X895**      **Методы и алгоритмы интеллектуального анализа данных:** учеб. пособие / *А.Г. Храмов.* – Самара: Изд-во Самарского университета, 2019. – 176 с.: ил.

**ISBN 978-5-7883-1414-3**

В данном учебном пособии содержится описание основных методов и алгоритмов интеллектуального анализа данных (Data Mining, Machine Learning). Разбираются примеры учебных практических задач, в том числе с использованием системы R статистического анализа данных и графики.

Пособие может быть использовано для подготовки магистров по направлению 01.04.02 Прикладная математика и информатика в рамках соответствующих магистерских специальностей.

Подготовлено на кафедре технической кибернетики факультета информатики института информатики, математики и электроники Самарского университета.

УДК 004.8(075)  
ББК 32.973.2я7

ISBN 978-5-7883-1414-3

© Самарский университет, 2019

# Оглавление

<b>Введение</b> .....	6
<b>1 Корреляционный анализ зависимости</b> .....	9
1.1 Коэффициент корреляции (Пирсона) .....	9
1.2 Ранговые коэффициенты корреляции .....	13
1.2.1 Коэффициент ранговой корреляции Спирмена .....	14
1.2.2 Коэффициент ранговой корреляции Кендалла .....	16
1.3 Примеры расчёта коэффициентов ранговой корреляции и проверки гипотезы о статистической независимости факторов ..	18
<b>2 Линейный регрессионный анализ</b> .....	23
2.1 Схема регрессионного эксперимента .....	24
2.2 Классическая модель линейной регрессии .....	27
2.3 Метод наименьших квадратов .....	28
2.4 Примеры использования метода наименьших квадратов .....	29
2.5 Свойства МНК-оценки .....	35
2.6 Теорема Гаусса–Маркова .....	36
2.7 Оценивание дисперсии шума наблюдения .....	38
2.8 Интервальные оценки параметров линейной регрессионной модели и проверка гипотез о коэффициентах регрессии .....	39
<b>3 Дисперсионный анализ</b> .....	43
3.1 Вероятностные распределения хи-квадрат и Фишера .....	43
3.2 Статистика Фишера .....	44
3.3 МНК-оценка линейной векторной функции коэффициентов регрессии .....	46
3.4 Однофакторный двухуровневый дисперсионный анализ .....	48
3.5 Однофакторный многоуровневый дисперсионный анализ .....	53
3.6 Двухфакторный дисперсионный анализ при произвольном числе наблюдений .....	60
3.7 Общий подход к дисперсионному анализу .....	63
3.8 Пример двухфакторного дисперсионного анализа .....	66
<b>4 Методы и алгоритмы кластеризации данных</b> .....	70
4.1 Задача кластеризации .....	70
4.2 Метод k-средних (k-means) .....	76

4.3	Алгоритм кластеризации Isodata (ИСОМАД).....	82
4.4	Иерархическая кластеризация.....	83
<b>5</b>	<b>Распознавание образов</b> .....	<b>90</b>
5.1	Статистическая проверка гипотез.....	90
5.2	Минимизация функции риска .....	101
5.3	Критерий Неймана–Пирсона.....	102
5.4	ROC-кривая.....	108
5.5	Наивный байесовский классификатор.....	109
<b>6</b>	<b>Линейные классификаторы</b> .....	<b>113</b>
6.1	Байесовский линейный классификатор.....	113
6.2	Линейная разделяющая функция, минимизирующая вероятность ошибки .....	116
<b>7</b>	<b>Дискриминантный анализ</b> .....	<b>127</b>
7.1	Матрицы рассеяния и критерий разделимости.....	127
7.2	Разделимость выборки .....	129
7.3	Выбор признаков, максимизирующих критерий $J_1$ .....	130
7.3.1	Линейное преобразование случайных векторов.....	130
7.3.2	Диагонализирующее и декоррелирующее преобразования .....	131
7.3.3	Одновременная диагонализация двух ковариационных матриц.....	132
7.3.4	Оптимизация сокращения размерности пространства признаков .....	135
<b>8</b>	<b>Метод опорных векторов (SVM)</b> .....	<b>137</b>
8.1	Линейно разделяемая выборка .....	137
8.2	Линейно неразделяемая выборка .....	142
8.3	Ядра и спрямляющие пространства.....	146
<b>9</b>	<b>Задачи для самостоятельного решения</b> .....	<b>149</b>
	<b>Список литература и интернет-ресурсы</b> .....	<b>157</b>
	<b>Приложение А Статистические функции</b> .....	<b>158</b>
A.1	Квантиль распределения .....	158
A.2	Нормальное распределение .....	158
A.3	Распределение хи-квадрат .....	160

A.4 Распределение Стьюдента .....	163
A.5 Распределение Фишера .....	166
<b>Приложение Б Сведения из теории матриц .....</b>	<b>170</b>
Б.1 Спектральное разложение симметричной матрицы .....	170
Б.2 Свойства операции <i>trace</i> (след) .....	170
Б.3 Линейное преобразование случайного вектора .....	171
<b>Приложение В Десять основных алгоритмов анализа данных .....</b>	<b>172</b>

## ВВЕДЕНИЕ

Data Mining (дословно *добыча данных*) – на русский язык чаще переводится как «анализ данных». Для перевода «*Data Mining*» на русский язык используют словосочетания: *просев информации, добыча данных, извлечение данных, интеллектуальный анализ данных, обнаружение знаний в базах данных*. Не существует «неинтеллектуального» анализа данных [А.Г. Дьяконов; 2012].

**Data Mining** – собирательное название для обозначения совокупности методов обнаружения в данных, доступных для интерпретации, ранее неизвестных, нетривиальных и практически полезных знаний, необходимых для принятия решений. Методы Data Mining: *классификация, моделирование, прогнозирование, деревья решений, искусственные нейронные сети, генетические алгоритмы, эволюционное программирование, ассоциативная память, нечёткая логика, статистические методы* (**дескриптивный анализ, корреляционный и регрессионный анализ, факторный анализ, дисперсионный анализ, компонентный анализ, дискриминантный анализ, анализ временных рядов, анализ выживаемости, анализ связей**). Выделенные разделы будут изучаться в рамках настоящего курса «*Интеллектуальный анализ данных*» или изучались ранее в курсе «*Тория случайных процессов*».

Интеллектуальный анализ данных – это раздел информатики, изучающий процессы обработки данных с целью получения полезной информации и принятия решений. Основная польза разрабатываемых методов анализа данных заключается в некоторой предсказательной способности: проанализировав неко-

торый набор данных, информационная система анализа данных должна обучиться для дальнейшего распознавания или прогнозирования некоторых участков данных в ситуациях, когда часть данных утеряна или неизвестна. Кроме этого, системы анализа данных могут решать задачи редуцирования объёма данных с целью устранения избыточности, визуализации данных для их удобного восприятия человеком, моделирования новых данных по имеющимся данным и др.

Впервые понятие Data Mining появилось в 1989 году. Изначально оно было связано с автоматизацией и оптимизацией запросов к крупным базам данных. Между тем понятие анализ данных (*англ.* Data Analysis) существовало намного раньше и означало обработку и интерпретацию данных, полученных в ходе экспериментов, в основном научных. С развитием науки и техники эти понятия расширились и обобщались, стали очень близки друг к другу и в настоящий момент тесно связаны как с анализом больших объёмов данных (*англ.* Big Data), так и с понятием машинного обучения (*англ.* Machine Learning).

Интеллектуальный анализ данных – это во многом прикладная теория, число приложений которой к реальным промышленным задачам растёт с каждым годом. В настоящее время методы и средства интеллектуального анализа данных используются при веб-разработке, в биоинформатике, в системах компьютерного зрения, в разработке компьютерных игр, в маркетинге, в медицинской диагностике, в методах оптимизации, при разработке поисковых систем, при распознавании образов, изображений, речи и сигналов и т.д. Востребованность специалистов по интеллектуальному анализу данных постоянно возрастает, как и доля финансирования разработок в этой области. Появляется всё больше программных решений для анализа данных, в том числе с открытым исходным кодом. Всё это свидетельствует о необходимости включения курса интеллектуального анализа данных в учебные программы по боль-

шинству технических специальностей, связанных с информатикой.

Смежная область – **Машинное обучение** (*Machine Learning*) – раздел искусственного интеллекта, математическая дисциплина, использующая математическую статистику, численные методы оптимизации, теорию вероятностей, выделяющая знания из данных. Различают два типа обучения: 1. *Обучение по прецедентам*, или *индуктивное обучение*, – основано на выявлении закономерностей в эмпирических данных. 2. *Дедуктивное обучение* – предполагает формализацию знаний экспертов и их перенос в компьютер в виде базы знаний, относится к области экспертных систем, поэтому термины *машинное обучение* и *обучение по прецедентам* можно считать синонимами.

В табл. В.1 (приложение В) приведены десять основных алгоритмов анализа данных. Мы их изучим в течение настоящего курса.



# 1 КОРРЕЛЯЦИОННЫЙ АНАЛИЗ ЗАВИСИМОСТИ

## 1.1 Коэффициент корреляции (Пирсона)

Простейшим способом проверки наличия статистической зависимости между двумя факторами является анализ их коэффициента корреляции (**Пирсона**). Допустим, что известна выборка  $\{x^i, y^i\}_{i=1}^N$  из генеральной совокупности  $(X, Y)$ , имеющей двумерное нормальное распределение с плотностью вероятности:

$$f_{XY}(x, y) = \frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-r^2}} \exp\left[-\frac{(x-m_x)^2}{2\sigma_x^2} + \frac{r(x-m_x)(y-m_y)}{\sigma_x\sigma_y} - \frac{(y-m_y)^2}{2\sigma_y^2}\right],$$

где  $m_x = \mathbf{MX}$ ,  $m_y = \mathbf{MY}$ ,  $\sigma_x^2 = \mathbf{DX}$ ,  $\sigma_y^2 = \mathbf{DY}$ ,

$$r = \frac{\mathbf{MXY} - \mathbf{MX} \cdot \mathbf{MY}}{\sqrt{\mathbf{DX} \cdot \mathbf{DY}}}.$$

Наличие линейной статистической зависимости между факторами  $X$  и  $Y$  характеризуется отличием от нуля коэффициента корреляции  $r$ . Точечной оценкой коэффициента корреляции  $r$  является выборочный коэффициент корреляции

$$\hat{r} = \frac{\sum_{i=1}^N (x^i - \bar{x})(y^i - \bar{y})}{\sqrt{\sum_{i=1}^N (x^i - \bar{x})^2 \sum_{i=1}^N (y^i - \bar{y})^2}} = \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{s_x \cdot s_y}, \quad (1.1)$$

$$\text{где } \bar{x} = \frac{1}{N} \sum_{i=1}^N x^i, \quad \bar{y} = \frac{1}{N} \sum_{i=1}^N y^i, \quad \overline{xy} = \frac{1}{N} \sum_{i=1}^N x^i y^i, \quad \overline{x^2} = \frac{1}{N} \sum_{i=1}^N (x^i)^2, \\ \overline{y^2} = \frac{1}{N} \sum_{i=1}^N (y^i)^2, \quad s_x^2 = \overline{x^2} - \bar{x}^2, \quad s_y^2 = \overline{y^2} - \bar{y}^2.$$

Точное вероятностное распределение выборочного коэффициента корреляции  $\hat{r}$  получено Р. Фишером в 1915 году [1]. Плотность вероятности имеет вид:

$$f_{\hat{r}}(\rho) = \frac{(1-r^2)^{(N-1)/2}}{\pi \Gamma(N-2)} (1-\rho^2)^{(N-4)/2} \frac{d^{N-2}}{d(\rho r)^{N-2}} \left\{ \frac{\arccos(-\rho r)}{\sqrt{1-r^2 \rho^2}} \right\}, |\rho| < 1.$$

В частности, при  $r=0$  (некоррелированность факторов  $X$  и  $Y$ ) плотность распределения  $\hat{r}$  имеет вид:

$$f_{\hat{r}}(\rho) = \frac{1}{B(N/2-1, 1/2)} (1-\rho^2)^{(N-4)/2}, \quad |\rho| < 1.$$

Здесь  $\Gamma(s) = \int_0^{\infty} t^{s-1} e^{-t} dt$  – гамма-функция,  $B(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}$  –

бета-функция.

Сложный вид вероятностного распределения выборочного коэффициента корреляции не позволяет использовать его в практических целях для получения интервальных оценок коэффициента корреляции и для проверки гипотезы о некоррелированности. Поэтому для этих целей используются различные асимптотические приближения [1]. В частности, для проверки гипотезы о независимости факторов  $X$  и  $Y$  ( $H_0: r=0$ ) используется статистика:

$$t = \frac{\hat{r}\sqrt{N-2}}{\sqrt{1-\hat{r}^2}} \sim t(N-2), \quad (1.2)$$

которая имеет  $t$ -распределение Стьюдента с числом степеней свободы  $(N - 2)$ . Факторы  $X$  и  $Y$  считаются линейно зависимыми, когда величина выборочного коэффициента корреляции  $\hat{r}$  отлична от нуля на уровне значимости  $\alpha$ , то есть если выполняется неравенство

$$\hat{r}^2 > \left[ 1 + \frac{N-2}{(t_{N-2}^{1-\alpha/2})^2} \right]^{-1}, \quad (1.3)$$

где  $t_{N-2}^{1-\alpha/2}$  – квантиль на уровне  $(1-\alpha/2)$   $t$ -распределения Стьюдента с  $(N - 2)$  степенями свободы. Чем с меньшим уровнем значимости отвергается гипотеза о некоррелированности, тем меньше вероятность того, что это сделано ошибочно. Таким образом, *уровень значимости* статистической гипотезы – это вероятность отвергнуть эту гипотезу при условии, что на самом деле она верна.

На рис. 1.1 изображены критические области для проверки гипотез (см. курс математической статистики).

Если решить неравенство (1.3) относительно  $\alpha$ , то можно найти граничный уровень значимости  $\alpha_0$ , для которого гипотеза о некоррелированности всё ещё отвергается при том, что при чуть меньшем уровне значимости она уже была бы принята при заданных выборочном коэффициенте корреляции  $\hat{r}$  и объёме выборки  $N$ :

$$\alpha < \alpha_0 = 2 - 2F_{t(N-2)} \left( \frac{|\hat{r}| \sqrt{N-2}}{\sqrt{1-\hat{r}^2}} \right),$$

где  $F_{t(N-2)}(\cdot)$  – интегральная функция  $t$ -распределения Стьюдента с  $(N - 2)$  степенями свободы. Значение граничного уровня значимости  $\alpha_0$  в западной литературе носит название **p-value**.

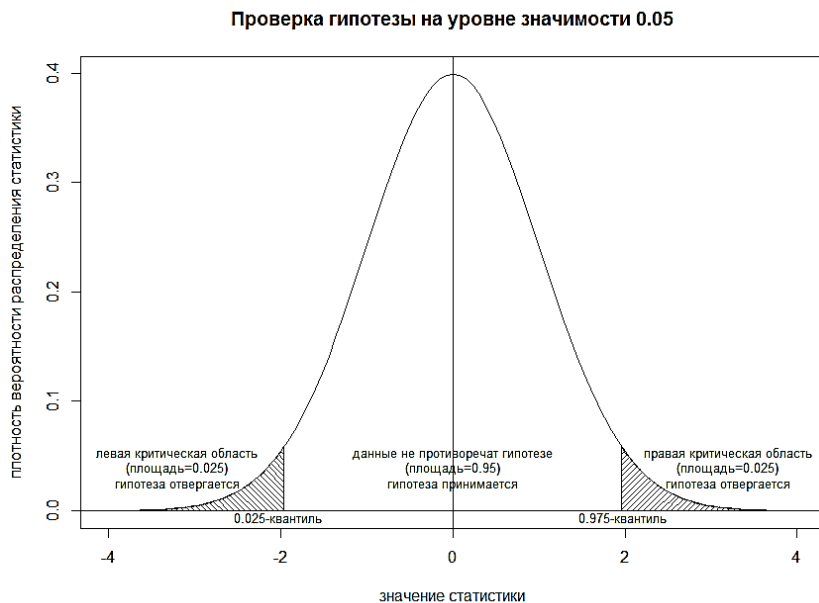


Рис. 1.1. Проверка гипотезы при двусторонней критической области

Эту же статистику (1.2) можно использовать для построения доверительных интервалов коэффициента корреляции, однако при значениях  $r$ , существенно отличных от нуля, более точные значения границ интервалов получаются при использовании  $z$ -преобразования Фишера коэффициента корреляции:

$$z = \operatorname{arcth} r = \frac{1}{2} \ln \frac{1+r}{1-r}; \quad r = \operatorname{th} z = \frac{e^{2z} - 1}{e^{2z} + 1},$$

где  $\operatorname{th}$  и  $\operatorname{arcth}$  – функции гиперболического тангенса и гиперболического арктангенса.

Статистика Фишера

$$z = (N - 3)^{1/2} (\operatorname{arcth} \hat{r} - \operatorname{arcth} r) \sim N(0;1) \quad (1.4)$$

распределена по стандартному нормальному закону  $N(0;1)$  асимптотически при  $N \rightarrow \infty$ .

Из (1.4) получаем доверительный интервал для коэффициента корреляции:

$$\Pr\left\{|\operatorname{arcth} \hat{r} - \operatorname{arcth} r| < (N-3)^{-1/2} u^{(1+p)/2}\right\} = p,$$

$$\Pr\left\{\operatorname{th}\left(\operatorname{arcth} \hat{r} - (N-3)^{-1/2} u^{(1+p)/2}\right) < r < \operatorname{th}\left(\operatorname{arcth} \hat{r} + (N-3)^{-1/2} u^{(1+p)/2}\right)\right\} = p, \quad (1.5)$$

где  $p$  — доверительная вероятность;  $u^{(1+p)/2}$  — квантиль на уровне  $(1+p)/2$  стандартного нормального распределения.

## 1.2 Ранговые коэффициенты корреляции

Иногда возникает потребность в статистическом анализе нечисловых факторов. Например, требуется проверить наличие статистической зависимости между результатами двух кругов чемпионата по футболу. Так может быть и в случае, когда признаками являются некоторые натуральные порядковые номера или другие целые числа, абсолютные значения которых несут только информацию о порядке объектов. В этом случае целесообразно вычислять корреляцию не между самими признаками, а между порядковыми номерами объектов, упорядоченных по этим признакам.

В этом случае вместо коэффициента корреляции Пирсона используются коэффициенты ранговой корреляции **Кендалла** и **Спирмена**.

### 1.2.1 Коэффициент ранговой корреляции Спирмена

Коэффициент ранговой корреляции Спирмена  $\rho$  (англ. Spearman's rank correlation coefficient) является частным случаем коэффициента корреляции Пирсона, если в качестве значений факторов  $X$ ,  $Y$  использовать значения их **рангов**  $x^i$ ,  $y^i$ , то есть индексов в упорядоченных последовательностях их значений:

$$X'_1 \leq X'_2 \leq X'_3 \leq \dots \leq X'_{N-1} \leq X'_N, \quad x^i \in [1; N],$$

$$Y'_1 \leq Y'_2 \leq Y'_3 \leq \dots \leq Y'_{N-1} \leq Y'_N, \quad y^i \in [1; N].$$

Докажем, что

$$\rho = 1 - \frac{6}{N(N^2 - 1)} \sum_{i=1}^N (x^i - y^i)^2,$$

где  $x^i$  – ранг фактора  $X_i$ ,  $y^i$  – ранг фактора  $Y_i$ ,  $i$  – номер объекта (наблюдения).

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x^i = \frac{1}{N} (1 + 2 + 3 + \dots + N) = \frac{N+1}{2},$$

$$\bar{y} = \frac{1}{N} \sum_{i=1}^N y^i = \frac{1}{N} (1 + 2 + 3 + \dots + N) = \frac{N+1}{2},$$

$$\overline{x^2} = \frac{1}{N} \sum_{i=1}^N (x^i)^2 = \frac{1}{N} (1^2 + 2^2 + 3^2 + \dots + N^2) =$$

$$= \frac{1}{N} \cdot \frac{N(N+1)(2N+1)}{6}, \quad \overline{y^2} = \overline{x^2} = \frac{(N+1)(2N+1)}{6},$$

$$s_x^2 = s_y^2 = \frac{(N+1)(2N+1)}{6} - \left( \frac{N+1}{2} \right)^2 = \frac{N^2 - 1}{12},$$

$$\begin{aligned} \sum_{i=1}^N (x^i - y^i)^2 &= \sum_{i=1}^N (x^i)^2 + \sum_{i=1}^N (y^i)^2 - 2 \sum_{i=1}^N x^i y^i = \\ &= \frac{N(N+1)(2N+1)}{6} + \frac{N(N+1)(2N+1)}{6} - 2 \sum_{i=1}^N x^i y^i. \end{aligned}$$

Из последнего уравнения следует, что

$$\begin{aligned} \sum_{i=1}^N x^i y^i &= \frac{N(N+1)(2N+1)}{6} - \frac{1}{2} \sum_{i=1}^N (x^i - y^i)^2, \\ \overline{xy} &= \frac{1}{N} \sum_{i=1}^N x^i y^i = \frac{(N+1)(2N+1)}{6} - \frac{1}{2N} \sum_{i=1}^N (x^i - y^i)^2. \end{aligned}$$

Тогда по определению выборочного коэффициента корреляции Пирсона из (1.1) получаем

$$\begin{aligned} \rho &= \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{s_x \cdot s_y} = \frac{\left[ \frac{(N+1)(2N+1)}{6} - \frac{1}{2N} \sum_{i=1}^N (x^i - y^i)^2 \right] - \left( \frac{N+1}{2} \right)^2}{\frac{N^2 - 1}{12}} = \\ &= 1 - \frac{6}{N(N^2 - 1)} \sum_{i=1}^N (x^i - y^i)^2. \end{aligned}$$

Для проверки гипотезы о независимости факторов  $X$  и  $Y$ , как и для коэффициента корреляции Пирсона, используется статистика (1.2)  $t = \rho \sqrt{N-2} / \sqrt{1-\rho^2}$ , которая имеет  $t$ -распределение Стьюдента с числом степеней свободы  $(N-2)$ .

Граничный (предельный) уровень значимости:

$$\alpha_0 = 2 - 2F_{t(N-2)} \left( \frac{|\rho| \sqrt{N-2}}{\sqrt{1-\rho^2}} \right).$$

### 1.2.2 Коэффициент ранговой корреляции Кендалла

Коэффициент ранговой корреляции Кендалла (англ. Kendall's rank correlation coefficient) вычисляют на основе ранжирования признаков.

Если рассмотреть все возможные пары объектов  $(i, j)$ ,  $i, j = 1, 2, 3, \dots, N$ , то про каждую пару можно сказать, одинаково ли или нет упорядочены признаки в этой паре. Обозначим через  $P$  – количество пар объектов, у которых признаки  $x$  и  $y$  упорядочены одинаково, а через  $Q$  – количество пар объектов, у которых признаки  $x$  и  $y$  упорядочены по-разному. Тогда

$$P - Q = \sum_{i=1}^N \sum_{j=1}^N \operatorname{sgn}(x_k - x_i) \operatorname{sgn}(y_k - y_i),$$

где  $x_i$  – ранг признака  $x$ ,  $y_i$  – ранг признака  $y$ ,  $i$  – номер объекта,

$$\operatorname{sgn}(u) = \begin{cases} 1, & u > 0; \\ 0, & u = 0; \\ -1, & u < 0. \end{cases}$$

Коэффициент ранговой корреляции вычисляют с использованием показателя  $(P - Q)$ :

$$\tau = \frac{2(P - Q)}{N(N - 1)}.$$

При «ручном» расчёте параметров  $P$  и  $Q$  можно использовать следующий алгоритм:

1. Значения фактора  $X$  выставляют в порядке возрастания и присваивают ранги.



2. Ранжируют значения показателя  $Y$ .
3. Рассчитывают:  $P$  – суммарное число наблюдений, следующих за текущими наблюдениями, с большим значением рангов  $Y$ ;  $Q$  – суммарное число наблюдений, следующих за текущими наблюдениями с меньшим значением рангов  $Y$  (равные ранги не учитываются).
4. Рассчитывают коэффициент корреляции:

$$\tau = \frac{2(P-Q)}{N(N-1)}.$$

$$\text{sgn}(u) = \begin{cases} 1, & u > 0; \\ 0, & u = 0; \\ -1, & u < 0. \end{cases}$$

При отсутствии зависимости факторов  $X$  и  $Y$  коэффициент ранговой корреляции  $\tau$  Кендалла имеет при  $N \geq 10$  имеет распределение, близкое к нормальному с параметрами  $\mathbf{M}\tau = 0$ ,

$\mathbf{D}\tau = \frac{2(2N+5)}{9N(N-1)}$ , то есть статистика  $u = \tau \sqrt{\frac{9N(N-1)}{2(2N+5)}}$  имеет приблизительно стандартное нормальное распределение.

Коэффициенты корреляции  $\tau$  и  $\rho$  изменяются от  $-1$  до  $+1$ . Если коэффициент корреляции равен  $+1$ , то это означает, что ранжировки одинаковы; если он равен  $-1$ , то противоположны (ранжировки обратны друг другу). Равенство коэффициента корреляции нулю означает, что ранжировки линейно независимы (некоррелированы).

### 1.3 Примеры расчёта коэффициентов ранговой корреляции и проверки гипотезы о статистической независимости факторов

**Пример 1.** Требуется найти коэффициенты ранговой корреляции Спирмена и Кендалла между местом команды после первого круга и местом команды по результатам сезона. Проверить гипотезу о независимости рассматриваемых факторов на уровне значимости **0.1**.

Таблица 1.1 – Итоги чемпионата России по футболу  
2002 года

	Локомотив	ЦСКА	Спартак	Торпедо	Крылья Советов	Сатурн
Итог сезона	1	2	3	4	5	6
Итог I круга	1	2	3	10	6	4
	Шинник	Динамо	Ротор	Зенит		
	7	8	9	10		
	8	7	5	9		

Количество объектов (команд):  $N = 10$ .

*Коэффициент ранговой корреляции Спирмена*

$$\sum_{i=1}^N (x_i - y_i)^2 = 0^2 + 0^2 + 0^2 + 6^2 + 1^2 + 2^2 + 1^2 + 1^2 + 4^2 + 1^2 = 60,$$

$$\rho = 1 - \frac{6}{N(N^2 - 1)} \sum_{i=1}^N (x_i - y_i)^2 = 1 - \frac{6}{10 \cdot 99} \cdot 60 = \frac{7}{11} \approx 0.636.$$

Значение  $t$ -статистики Стьюдента:

$$t = \frac{\rho \sqrt{N-2}}{\sqrt{1-\rho^2}} = \frac{(7/11) \sqrt{10-2}}{\sqrt{1-(7/11)^2}} = \frac{7\sqrt{8}}{\sqrt{72}} = \frac{7}{3} \approx 2.333333.$$

Критическое значение – квантиль на уровне 0.95  $t$ -распределения Стьюдента с 8 степенями свободы  $t_8^{0.95} \approx 1.859548$ . Так как  $2.333333 > 1.859548$ , то гипотезу о независимости рассматриваемых факторов можно отвергнуть на уровне значимости **0.1**, то есть факторы считаем статистически зависимыми на этом уровне значимости.

Граничный уровень значимости ( $p$ -value):

$$\alpha_0 = 2 - 2F_{N-2}^t \left( \frac{|\rho| \sqrt{N-2}}{\sqrt{1-\rho^2}} \right) = 2 - 2F_8^t \left( \frac{7}{3} \right) \approx 0.0479,$$

где  $F_8^t(\cdot)$  – интегральная функция  $t$ -распределения Стьюдента с 8 степенями свободы.

### ***Коэффициент ранговой корреляции Кендалла***

$$S = P - Q = (9 + 8 + 7 + 0 + 3 + 4 + 1 + 1 + 1 + 0) - (0 + 0 + 0 + 6 + 2 + 0 + 2 + 1 + 0 + 0) = 34 - 11 = 23,$$

$$\tau = \frac{2S}{N(N-1)} = \frac{2 \cdot 23}{10 \cdot 9} = \frac{23}{45} \approx \mathbf{0.511}.$$

Находим значение статистики

$$u = \tau \sqrt{\frac{9N(N-1)}{2(2N+5)}} = \frac{23}{45} \sqrt{\frac{90 \cdot 9}{2 \cdot 25}} = \frac{23}{5\sqrt{5}} \approx 2.057183.$$

Критическое значение – квантиль на уровне 0.95 стандартного нормального распределения  $u^{0.95} \approx 1.644854$ . Так как  $2.057183 > 1.644854$ , то гипотезу о независимости рассматриваемых факторов можно отвергнуть на уровне значимости 0.1.

Граничный уровень значимости (*p-value*):

$$\alpha_0 = 2 - 2F_u \left( \tau \sqrt{\frac{9N(N-1)}{2(2N+5)}} \right) = 2 - 2F_u \left( \frac{23}{5\sqrt{5}} \right) \approx 0.0397,$$

где  $F_u(\cdot)$  – интегральная функция стандартного нормального распределения.

Видим, что *p-value* при проверке с помощью ранговых коэффициентов корреляции гипотезы о некоррелированности между факторами (местом команды после первого круга и местом команды по результатам сезона) лежит в пределах 4–5%.

**Пример 2.** Вычислить выборочные коэффициенты корреляции Пирсона  $r$ , Спирмена  $\rho$  и Кендалла  $\tau$  двух **бинарных** факторов:  $X$  – пациенту проведена или не проведена вакцинация от гриппа,  $Y$  – клиент заболел или не заболел гриппом. Протокол прививок и наблюдений задан следующей таблицей (число объектов  $N = 9$ ):

№ пациента		1	2	3	4	5	6	7	8	9
$X$	Прививка	1	1	1	0	1	0	0	1	0
$Y$	Заболевание	0	1	1	1	1	0	0	1	0

**Коэффициент корреляции Пирсона  $r$**

$$r = \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{s_x \cdot s_y} = \frac{\frac{4}{9} - \frac{5}{9} \cdot \frac{5}{9}}{\frac{5}{9} - \left(\frac{5}{9}\right)^2} = \frac{11}{20} = 0.55.$$

**Коэффициент корреляции Спирмена  $\rho$**

Для вычисления ранговой корреляции Спирмена проведём ранжирование факторов  $X$  и  $Y$ :

№ пациента		1	2	3	4	5	6	7	8	9
x	Прививка (ранг)	5	6	7	1	8	2	3	9	4
y	Заболевание (ранг)	1	5	6	7	8	2	3	9	4

$$\rho = 1 - \frac{6}{N(N^2 - 1)} \sum_{i=1}^N (x_i - y_i)^2 =$$

$$= 1 - \frac{6}{9 \cdot 80} (4^2 + 1^2 + 1^2 + 6^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2) = 1 - \frac{6 \cdot 54}{9 \cdot 80} = \frac{11}{20} = 0.55.$$

### **Коэффициент корреляции Кендалла $\tau$**

Для вычисления ранговой корреляции Кендалла расположим один из факторов (например, фактор X) в порядке возрастания рангов и запишем соответствующие ранги фактора Y:

№ пациента		4	6	7	9	1	2	3	5	8
x	Прививка (ранг)	1	2	3	4	5	6	7	8	9
y	Заболевание (ранг)	7	2	3	4	1	5	6	8	9

$$\tau = \frac{2(P - Q)}{N(N - 1)} = \frac{2 \cdot 18}{9 \cdot 8} = 0.5;$$

$$P = 2 + 6 + 5 + 4 + 4 + 3 + 2 + 1 + 0 = 27,$$

$$Q = 6 + 1 + 1 + 1 + 0 + 0 + 0 + 0 + 0 = 9.$$

**Пример 3.** На уровне значимости  $\alpha = 0.1$  проверить гипотезу о статистической независимости рабочего стажа и показателя травматизма с использованием рангового коэффициента корреляции Спирмена.

X – Рабочий стаж в годах	До 1 года	1–2 года	3–4 года	5–6 лет	7 лет и более
Y – Травматизм на 100 работающих	24	16	12	12	6

(Данные взяты на сайте <http://forex365.ru/indicators/korreljacionnyj-analiz-spirmena.html> – Корреляционный анализ Спирмена).

Запишем ранги факторов (по возрастанию значений факторов, с корректировкой одинаковых рангов):

$x$ – Ранги рабочего стажа	1	2	3	4	5
$y$ – Рани травматизма	5	4	2,5	2,5	1

$$\rho = 1 - \frac{6}{N(N^2 - 1)} \sum_{i=1}^N (x_i - y_i)^2 = 1 - \frac{6}{5 \cdot 24} (4^2 + 2^2 + 0.5^2 + 1.5^2 + 4^2) =$$

$$= -\frac{37}{40} = -0.925.$$

Значение  $t$ -статистики Стьюдента:

$$t = \frac{|\rho| \sqrt{N-2}}{\sqrt{1-\rho^2}} = \frac{(37/40) \sqrt{5-2}}{\sqrt{1-(37/40)^2}} = \frac{37}{\sqrt{77}} \approx 4.217.$$

Критическое значение – квантиль на уровне 0.95  $t$ -распределения Стьюдента с 3 степенями свободы  $t_3^{0.95} \approx 2.353363$ . Так как  $4.217 > 2.353363$ , то гипотезу о независимости рассматриваемых факторов можно отвергнуть на уровне значимости **0.1**, то есть факторы считаем статистически зависимыми на этом уровне значимости, то есть травматизм зависит от величины рабочего стажа.

Граничный уровень значимости ( $p$ -value):

$$\alpha_0 = 2 - 2F_{N-2}^t \left( \frac{|\rho| \sqrt{N-2}}{\sqrt{1-\rho^2}} \right) = 2 - 2F_3^t \left( \frac{37}{\sqrt{77}} \right) \approx 0.0244,$$

где  $F_3^t(\cdot)$  – интегральная функция  $t$ -распределения Стьюдента с 3 степенями свободы.

## 2 ЛИНЕЙНЫЙ РЕГРЕССИОННЫЙ АНАЛИЗ

**Регрессионный анализ** (*англ.* regression analysis) – это статистический метод исследования влияния одной или нескольких независимых переменных  $\mathbf{x} = (x_1 \ x_2 \ x_3 \ \dots \ x_n)^T$  на зависимую переменную  $y$ . Независимые переменные иначе называют **регрессорами** (*англ.* regressors) или **предикторами** (*англ.* predictors), а зависимую переменную – **критериальной** (*англ.* criterion variable). В частности, регрессионный анализ помогает нам понять, как значение критериальной переменной изменяется, если меняется какая-либо из независимых переменных, а другие независимые переменные фиксированы. В регрессионном анализе оценивается условное математическое ожидание (среднее значение) зависимой переменной при заданных значениях предикторов:

$$f(x_1, x_2, x_3, \dots, x_n) = \mathbf{M}(y | x_1, x_2, x_3, \dots, x_n).$$

Функция  $y = f(x_1, x_2, x_3, \dots, x_n)$  называется функцией регрессии  $y$  на  $x_1, x_2, x_3, \dots, x_n$ .

Для построения функции регрессии  $f$  необходимо задать её общий вид и параметры. В зависимости от вида функции и её неизвестных параметров различают **линейную** регрессию и **нелинейную** регрессию.

## 2.1 Схема регрессионного эксперимента

Для построения функции регрессии  $f$  статистическим методом зададим *схему регрессионного эксперимента*, которую в общем случае можно описать уравнением наблюдения:

$$y^i = f(\mathbf{x}^i, \boldsymbol{\theta}, \xi^i), \quad i = \overline{1, N}, \quad (2.1)$$

где  $i$  – номер опыта;  $N$  – число опытов;  $\mathbf{x} = (x_1 \ x_2 \ x_3 \ \dots \ x_n)^T$  – вектор независимых переменных,  $\mathbf{x}^i \in \mathbf{R}^n$ ;  $n$  – размерность пространства независимых переменных;  $\mathbf{x}^1, \mathbf{x}^2, \mathbf{x}^3, \dots, \mathbf{x}^N$  – точки наблюдения:  $\mathbf{x}^i = (x_1^i \ x_2^i \ x_3^i \ \dots \ x_n^i)^T$ ,  $i = \overline{1, N}$ ;  $y^i$  – значение выходного (наблюдаемого, измеряемого) фактора в  $i$ -м опыте;  $\xi^i$  – случайное значение шумового фактора (помехи) в  $i$ -м опыте;  $\boldsymbol{\theta} = (\theta_1 \ \theta_2 \ \theta_3 \ \dots \ \theta_m)^T$  – вектор параметров объекта исследования, подлежащий оцениванию;  $m$  – число параметров. Параметры  $\theta_1, \theta_2, \theta_3, \dots, \theta_m$  считаются детерминированными (постоянными, неслучайными) величинами.

На рис. 2.1 приведена иллюстрация нелинейного регрессионного эксперимента. Здесь  $\mathbf{x}^i = x_1^i = x^i$  ( $n=1$ ),  $\boldsymbol{\theta} = (\omega \ a \ b)^T$ ,  $y^i = a \sin(\omega x^i) + b + \xi^i$ . Заметим, что если  $\omega = \text{const}$  – известная детерминированная величина, то рассматриваемая регрессия является линейной по параметрам  $\boldsymbol{\theta} = (a \ b)^T$ .

На рис. 2.2 приведена иллюстрация простейшего линейного регрессионного эксперимента. Здесь  $\mathbf{x}^i = x_1^i = x^i$  ( $n=1$ ),  $\boldsymbol{\theta} = (a \ b)^T$ ,  $y^i = ax^i + b + \xi^i$ .



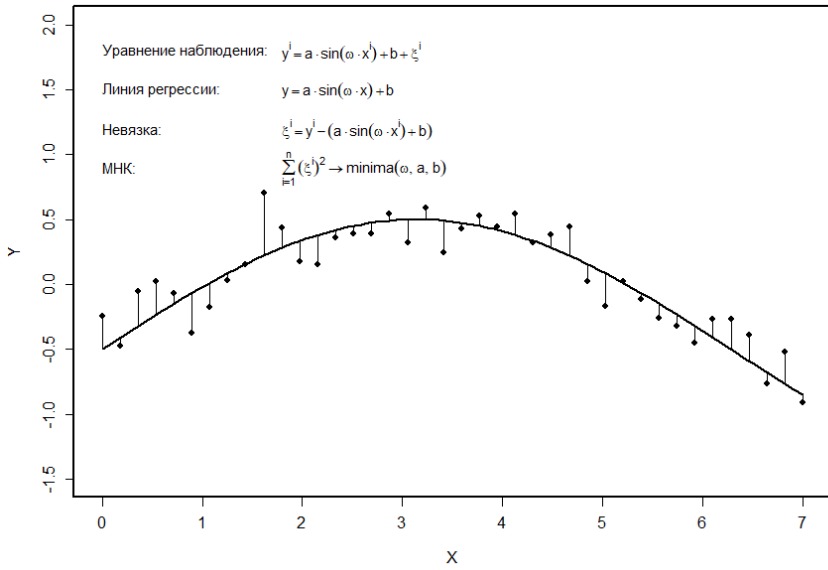


Рис. 2.1. Нелинейная регрессия

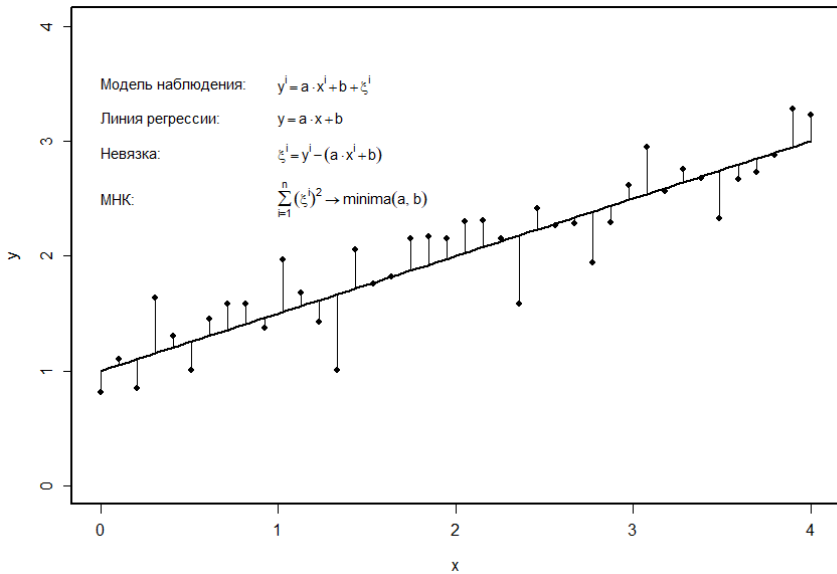


Рис. 2.2. Линейная регрессия

В общем случае модель наблюдения (2.1) с линейной относительно параметров функцией регрессии  $f$  и с аддитивным шумом наблюдения  $\xi$  можно единственным образом представить в виде:

$$y = f(\mathbf{x}, \boldsymbol{\theta}, \xi) = f_1(\mathbf{x})\theta_1 + f_2(\mathbf{x})\theta_2 + \dots + f_m(\mathbf{x})\theta_m + \xi, \quad (2.2)$$

где  $f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_m(\mathbf{x})$  – некоторые функции, не зависящие от неизвестных параметров  $\theta_1, \theta_2, \theta_3, \dots, \theta_m$  и являющиеся коэффициентами перед этими неизвестными параметрами. Такие  $m$  функций будем называть *базовыми*.

Заметим, что в таком представлении базовые функции  $f_k(\mathbf{x})$  и сама функция регрессии  $f(\mathbf{x}, \boldsymbol{\theta}, \xi)$  могут быть нелинейными от независимых переменных  $\mathbf{x}$ . Например,  $f(\mathbf{x}, \boldsymbol{\theta}, \xi) = a \sin(\omega x) + b + \xi$  – нелинейная функция  $x$ , когда параметры  $\boldsymbol{\theta} = (a \ b)^T$ , а  $\omega = \text{const}$  (см. рис. 2.1).

Объединяя базовые функции в вектор-столбец  $\mathbf{f}(\mathbf{x}) = (f(\mathbf{x}) \ f(\mathbf{x}) \ \dots \ f(\mathbf{x}))^T$ , приходим к матричной форме записи уравнения наблюдения:

$$\mathbf{Y} = \mathbf{F}\boldsymbol{\theta} + \boldsymbol{\xi}, \text{ то есть}$$

$$\begin{pmatrix} y^1 \\ y^2 \\ y^3 \\ \dots \\ y^N \end{pmatrix} = \begin{bmatrix} f_1^1 & f_2^1 & \dots & f_m^1 \\ f_1^2 & f_2^2 & \dots & f_m^2 \\ f_1^3 & f_2^3 & \dots & f_m^3 \\ \dots & \dots & \dots & \dots \\ f_1^N & f_2^N & \dots & f_m^N \end{bmatrix} \begin{pmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \dots \\ \theta_m \end{pmatrix} + \begin{pmatrix} \xi^1 \\ \xi^2 \\ \xi^3 \\ \dots \\ \xi^N \end{pmatrix}. \quad (2.3)$$

Матрица

$$\mathbf{F} = \begin{bmatrix} f_1^1 & f_2^1 & \dots & f_m^1 \\ f_1^2 & f_2^2 & \dots & f_m^2 \\ f_1^3 & f_2^3 & \dots & f_m^3 \\ \dots & \dots & \dots & \dots \\ f_1^N & f_2^N & \dots & f_m^N \end{bmatrix} \quad (2.4)$$

размера  $N \times m$  называется матрицей эксперимента.

Здесь  $f_k^i = f_k(\mathbf{x}^i) = f_k(x_1^i, x_2^i, x_3^i, \dots, x_n^i)$ ;  $i = \overline{1, N}$  – значение  $k$ -й базовой функции при  $i$ -м наблюдении;  $\mathbf{Y} = (y^1, y^2, y^3, \dots, y^N)^T$  – вектор наблюдения;  $\xi = (\xi^1, \xi^2, \xi^3, \dots, \xi^N)^T$  – случайный вектор шума наблюдения.

## 2.2 Классическая модель линейной регрессии

Будем называть *классической моделью линейной регрессии* схему эксперимента (2.3), когда шумы наблюдения  $\xi^i$  являются некоррелированными (независимыми) для различных наблюдений случайными величинами, имеющими нулевое математическое ожидание и одинаковую дисперсию  $\sigma^2$ :

$$\mathbf{M}\xi^i = 0, \quad \mathbf{M}\xi^i \xi^j = \sigma^2 \delta_{ij} \quad (\delta_{ij} \text{ – символ Кронекера}).$$

Нулевое математическое ожидание шума наблюдения означает одинаковые в среднем отклонения наблюдаемых значений от линии регрессии, одинаковость дисперсии шума наблюдения в различных опытах, одинаковый в среднем разброс наблюдаемых значений от линии регрессии. Некоррелированность (независимость) шума наблюдения в различных опытах говорит о невлинии опытов друг на друга. Таким образом, классическую модель наблюдения можно записать в матричном виде следующим образом:

$$\mathbf{Y} = \mathbf{F}\boldsymbol{\theta} + \boldsymbol{\xi}, \quad \mathbf{M}\boldsymbol{\xi} = \mathbf{0}, \quad \mathbf{D}\boldsymbol{\xi} = \sigma^2 \mathbf{I}_N, \quad (2.5)$$

где  $\mathbf{D}\boldsymbol{\xi} = \mathbf{M}\boldsymbol{\xi}\boldsymbol{\xi}^T$  – дисперсионная (корреляционная) матрица шума наблюдения,  $\mathbf{I}_N$  – единичная матрица размера  $N \times N$ .

Классическая модель линейной регрессии *при гауссовом шуме наблюдения* имеет вид:

$$\mathbf{Y} = \mathbf{F}\boldsymbol{\theta} + \boldsymbol{\xi}, \quad \boldsymbol{\xi} \sim \mathbf{N}(\mathbf{0}, \sigma^2 \mathbf{I}_N), \quad (2.6)$$

где запись  $\sim \mathbf{N}(\mathbf{a}, \mathbf{R})$  обозначает нормальный случайный вектор с математическим ожиданием  $\mathbf{a}$  и корреляционной матрицей  $\mathbf{R}$ .

### 2.3 Метод наименьших квадратов

*МНК-оценкой (англ. Least squares estimation) параметров* линейной регрессии  $\hat{\boldsymbol{\theta}} = (\hat{\theta}_1 \ \hat{\theta}_2 \ \dots \ \hat{\theta}_m)^T$  называется оценка параметров  $\boldsymbol{\theta} = (\theta_1 \ \theta_2 \ \dots \ \theta_m)^T$ , минимизирующая сумму квадратов шумовых невязок:

$$R_0^2 = \sum_{i=1}^N \left[ y^i - (f_1^i \theta_1 + f_2^i \theta_2 + f_3^i \theta_3 + \dots + f_m^i \theta_m) \right]^2 \rightarrow \min_{\theta_1, \theta_2, \dots, \theta_m}. \quad (2.7)$$

Величина  $R_0^2$  называется остаточной суммой квадратов невязок.

В матричном виде:  $R_0^2 = (\mathbf{Y} - \mathbf{F}\boldsymbol{\theta})^T (\mathbf{Y} - \mathbf{F}\boldsymbol{\theta}) \rightarrow \min_{\boldsymbol{\theta}}$ , то есть

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} (\mathbf{Y} - \mathbf{F}\boldsymbol{\theta})^T (\mathbf{Y} - \mathbf{F}\boldsymbol{\theta}).$$

Из условия минимума суммы квадратов невязок (2.7)

$$\frac{\partial R_0^2}{\partial \theta} = \frac{\partial}{\partial \theta} \left[ (\mathbf{Y} - \mathbf{F}\theta)^T (\mathbf{Y} - \mathbf{F}\theta) \right] = \mathbf{F}^T (\mathbf{Y} - \mathbf{F}\theta) = \mathbf{0}$$

получаем систему нормальных уравнений относительно МНК-оценки  $\hat{\theta}$ :

$$\mathbf{F}^T \mathbf{F} \hat{\theta} = \mathbf{F}^T \mathbf{Y}. \quad (2.8)$$

Подробный вывод системы нормальных уравнений приведён в [2].

Здесь и в дальнейшем считаем, что система нормальных уравнений (2.8) является невырожденной, то есть  $\text{rank } \mathbf{F} = m$  и обратная матрица  $(\mathbf{F}^T \mathbf{F})^{-1}$  существует. Тогда МНК-оценка параметров линейной регрессии  $\hat{\theta}$  может быть определена из следующего соотношения:

$$\hat{\theta} = (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{Y}. \quad (2.9)$$

Если система нормальных уравнений (2.8) является вырожденной, то получить МНК-оценку можно, воспользовавшись аппаратом псевдообратных матриц и наложением дополнительных условий [2].

## 2.4 Примеры использования метода наименьших квадратов

### *Пример 1. Линейная регрессия на прямую.*

Рассмотрим простейшую линейную регрессию на прямую, показанную на рис. 2.2. Наблюдению доступно  $N$  точек на плоскости  $(x_i, y_i)$ . Уравнение наблюдения:

$$y_i = ax_i + b + \xi_i.$$

Вектор неизвестных параметров:  $\boldsymbol{\theta} = \begin{pmatrix} a \\ b \end{pmatrix}$ .

$$\begin{pmatrix} Y_1 \\ Y_2 \\ \dots \\ Y_N \end{pmatrix} = \begin{pmatrix} x_1 & 1 \\ x_2 & 1 \\ \dots & \\ x_N & 1 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} + \begin{pmatrix} \xi_1 \\ \xi_2 \\ \dots \\ \xi_N \end{pmatrix} \Rightarrow \mathbf{Y} = \mathbf{F}\boldsymbol{\theta} + \boldsymbol{\xi}, \text{ матрица}$$

$$\text{эксперимента: } \mathbf{F} = \begin{pmatrix} x_1 & 1 \\ x_2 & 1 \\ \dots & \\ x_N & 1 \end{pmatrix},$$

$$\mathbf{F}^T \mathbf{F} = \begin{pmatrix} \sum_{i=1}^N x_i^2 & \sum_{i=1}^N x_i \\ \sum_{i=1}^N x_i & N \end{pmatrix} = N \begin{pmatrix} \overline{x^2} & \bar{x} \\ \bar{x} & 1 \end{pmatrix},$$

$$(\mathbf{F}^T \mathbf{F})^{-1} = \frac{1}{N(\overline{x^2} - \bar{x}^2)} \begin{pmatrix} 1 - \bar{x} & \\ -\bar{x} & \overline{x^2} \end{pmatrix},$$

$$\mathbf{F}^T \mathbf{Y} = \begin{pmatrix} \sum_{i=1}^N x_i Y_i \\ \sum_{i=1}^N Y_i \end{pmatrix} = N \begin{pmatrix} \overline{xy} \\ \bar{y} \end{pmatrix},$$

$$\hat{\boldsymbol{\theta}} = (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{Y} = \frac{1}{\overline{x^2} - \bar{x}^2} \begin{pmatrix} 1 - \bar{x} & \\ -\bar{x} & \overline{x^2} \end{pmatrix} \begin{pmatrix} \overline{xy} \\ \bar{y} \end{pmatrix}.$$

С использованием обозначений п.1.1 получаем:

$$\hat{a} = \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{\overline{x^2} - \bar{x}^2} = r \frac{s_y}{s_x}, \quad \hat{b} = \frac{\overline{x^2} \cdot \bar{y} - \bar{x} \cdot \overline{xy}}{\overline{x^2} - \bar{x}^2} = \frac{1}{\overline{x^2} - \bar{x}^2} \begin{pmatrix} \overline{xy} - \bar{x} \cdot \bar{y} \\ \overline{x^2} \cdot \bar{y} - \bar{x} \cdot \overline{xy} \end{pmatrix}.$$

**Пример 2. Квадратичная аппроксимация.**

Построить МНК-оценки параметров  $a, b, c$  квадратичной функции  $y = ax^2 + bx + c$  по наблюдениям в пяти точках  $x \in \{-2, -1, 0, 1, 2\}$ .

**Решение.**

$$\begin{pmatrix} Y_1 \\ Y_2 \\ Y_3 \\ Y_4 \\ Y_5 \end{pmatrix} = \begin{pmatrix} 4 & -2 & 1 \\ 1 & -1 & 1 \\ 0 & 0 & 1 \\ 1 & 1 & 1 \\ 4 & 2 & 1 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \end{pmatrix} + \begin{pmatrix} \xi_1 \\ \xi_2 \\ \xi_3 \\ \xi_4 \\ \xi_5 \end{pmatrix} \Rightarrow \mathbf{Y} = \mathbf{F}\boldsymbol{\theta} + \boldsymbol{\xi},$$

$$\mathbf{F} = \begin{pmatrix} 4 & -2 & 1 \\ 1 & -1 & 1 \\ 0 & 0 & 1 \\ 1 & 1 & 1 \\ 4 & 2 & 1 \end{pmatrix},$$

$$\mathbf{F}^T \mathbf{F} = \begin{pmatrix} 34 & 0 & 10 \\ 0 & 10 & 0 \\ 10 & 0 & 5 \end{pmatrix},$$

$$(\mathbf{F}^T \mathbf{F})^{-1} = \frac{1}{700} \begin{pmatrix} 50 & 0 & -100 \\ 0 & 70 & 0 \\ -100 & 0 & 340 \end{pmatrix} = \frac{1}{70} \begin{pmatrix} 5 & 0 & -10 \\ 0 & 7 & 0 \\ -10 & 0 & 34 \end{pmatrix},$$

$$\mathbf{F}^T \mathbf{Y} = \begin{pmatrix} 4 & 1 & 0 & 1 & 4 \\ -2 & -1 & 0 & 1 & 2 \\ 1 & 1 & 1 & 1 & 1 \end{pmatrix} \cdot \begin{pmatrix} Y_1 \\ Y_2 \\ Y_3 \\ Y_4 \\ Y_5 \end{pmatrix} = \begin{pmatrix} 4Y_1 + Y_2 + Y_4 + 4Y_5 \\ -2Y_1 - Y_2 + Y_4 + 2Y_5 \\ Y_1 + Y_2 + Y_3 + Y_4 + Y_5 \end{pmatrix},$$

$$\hat{\boldsymbol{\theta}} = (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{Y} = \frac{1}{70} \begin{pmatrix} 5 & 0 & -10 \\ 0 & 7 & 0 \\ -10 & 0 & 34 \end{pmatrix} \begin{pmatrix} 4y_1 + y_2 + y_4 + 4y_5 \\ -2y_1 - y_2 + y_4 + 2y_5 \\ y_1 + y_2 + y_3 + y_4 + y_5 \end{pmatrix},$$

$$\hat{\boldsymbol{\theta}} = \begin{pmatrix} \frac{1}{70} [5(4y_1 + y_2 + y_4 + 4y_5) - 10(y_1 + y_2 + y_3 + y_4 + y_5)] \\ \frac{1}{10} (-2y_1 - y_2 + y_4 + 2y_5) \\ \frac{1}{70} [-10(4y_1 + y_2 + y_4 + 4y_5) + 34(y_1 + y_2 + y_3 + y_4 + y_5)] \end{pmatrix},$$

$$\hat{a} = \frac{4y_1 - y_2 - 2y_3 - y_4 + 4y_5}{14}, \quad \hat{b} = \frac{-2y_1 - y_2 + y_4 + 2y_5}{10},$$

$$\hat{c} = \frac{-3y_1 + 12y_2 + 17y_3 + 12y_4 - 3y_5}{35}.$$

**Пример 3. Билинейная аппроксимация по четырём точкам.**

Построить МНК-оценки параметров  $a, b, c$  симметричной билинейной функции двух переменных  $z(x, y) = axy + b(x + y) + c$  по наблюдениям в четырёх точках:

$$z_1 = z(0; 0), \quad z_2 = z(0; 1), \quad z_3 = z(1; 0), \quad z_4 = z(1; 1).$$

**Решение.**

$$\boldsymbol{\theta} = \begin{pmatrix} a \\ b \\ c \end{pmatrix}, \quad \begin{pmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \end{pmatrix} = \begin{pmatrix} 0 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 1 & 1 \\ 1 & 2 & 1 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \end{pmatrix} + \begin{pmatrix} \xi_1 \\ \xi_2 \\ \xi_3 \\ \xi_4 \end{pmatrix} \Rightarrow \mathbf{Z} = \mathbf{F}\boldsymbol{\theta} + \boldsymbol{\xi},$$

$$\mathbf{F} = \begin{pmatrix} 0 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 1 & 1 \\ 1 & 2 & 1 \end{pmatrix}, \quad \mathbf{F}^T = \begin{pmatrix} 0 & 0 & 0 & 1 \\ 0 & 1 & 1 & 2 \\ 1 & 1 & 1 & 1 \end{pmatrix}, \quad \mathbf{F}^T \mathbf{F} = \begin{pmatrix} 1 & 2 & 1 \\ 2 & 6 & 4 \\ 1 & 4 & 4 \end{pmatrix},$$



$$\begin{aligned}
 (\mathbf{F}^T \mathbf{F})^{-1} &= \frac{1}{2} \begin{pmatrix} 8 & -4 & 2 \\ -4 & 3 & -2 \\ 2 & -2 & 2 \end{pmatrix}, \quad \mathbf{F}^T \mathbf{Y} = \begin{pmatrix} z_4 \\ z_2 + z_3 + 2z_4 \\ z_1 + z_2 + z_3 + z_4 \end{pmatrix}, \\
 \hat{\boldsymbol{\theta}} &= (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{Y} = \frac{1}{2} \begin{pmatrix} 8 & -4 & 2 \\ -4 & 3 & -2 \\ 2 & -2 & 2 \end{pmatrix} \cdot \begin{pmatrix} z_4 \\ z_2 + z_3 + 2z_4 \\ z_1 + z_2 + z_3 + z_4 \end{pmatrix} = \\
 &= \begin{pmatrix} z_1 - z_2 - z_3 + z_4 \\ \frac{1}{2}(z_2 + z_3 - 2z_1) \\ z_1 \end{pmatrix}.
 \end{aligned}$$

$$\text{Ответ: } \hat{a} = \frac{z_1 - z_2 - z_3 + z_4}{2}, \quad \hat{b} = \frac{z_2 + z_3}{2} - z_1, \quad \hat{c} = z_1.$$

**Пример 4. Билинейная аппроксимация по девяти точкам.**

Построить МНК-оценки параметров  $a, b, c, d$  билинейной функции двух переменных  $z(x, y) = axy + bx + cy + d$  по наблюдениям в девяти точках на сетке  $3 \times 3$ :

$$\begin{aligned}
 z_1 &= z(1; 1), \quad z_2 = z(1; 0), \quad z_3 = z(1; -1), \\
 z_4 &= z(0; 1), \quad z_5 = z(0; 0), \quad z_6 = z(0; -1), \\
 z_7 &= z(-1; 1), \quad z_8 = z(-1; 0), \quad z_9 = z(-1; -1).
 \end{aligned}$$

**Решение.**

$$\begin{pmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \\ z_5 \\ z_6 \\ z_7 \\ z_8 \\ z_9 \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ -1 & 1 & -1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & -1 & 1 \\ -1 & -1 & 1 & 1 \\ 0 & -1 & 0 & 1 \\ 1 & -1 & -1 & 1 \end{pmatrix} \cdot \begin{pmatrix} a \\ b \\ c \\ d \end{pmatrix} + \begin{pmatrix} \xi_1 \\ \xi_2 \\ \xi_3 \\ \xi_4 \\ \xi_5 \\ \xi_6 \\ \xi_7 \\ \xi_8 \\ \xi_9 \end{pmatrix} \Rightarrow \mathbf{Y} = \mathbf{F}\boldsymbol{\theta} + \boldsymbol{\xi},$$

$$\mathbf{F} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ -1 & 1 & -1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & -1 & 1 \\ -1 & -1 & 1 & 1 \\ 0 & -1 & 0 & 1 \\ 1 & -1 & -1 & 1 \end{pmatrix},$$

$$\mathbf{F}^T \mathbf{F} = \begin{pmatrix} 4 & 0 & 0 & 0 \\ 0 & 6 & 0 & 0 \\ 0 & 0 & 6 & 0 \\ 0 & 0 & 0 & 8 \end{pmatrix}, \quad (\mathbf{F}^T \mathbf{F})^{-1} = \begin{pmatrix} 1/4 & 0 & 0 & 0 \\ 0 & 1/6 & 0 & 0 \\ 0 & 0 & 1/6 & 0 \\ 0 & 0 & 0 & 1/8 \end{pmatrix},$$

$$\mathbf{F}^T \mathbf{Y} = \begin{pmatrix} 1 & 0 & -1 & 0 & 0 & 0 & -1 & 0 & 1 \\ 1 & 1 & 1 & 0 & 0 & 0 & -1 & -1 & -1 \\ 1 & 0 & -1 & 1 & 0 & -1 & 1 & 0 & -1 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{pmatrix} \cdot \begin{pmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \\ z_5 \\ z_6 \\ z_7 \\ z_8 \\ z_9 \end{pmatrix} =$$

$$= \begin{pmatrix} z_1 - z_3 - z_7 + z_9 \\ z_1 + z_2 + z_3 - z_7 - z_8 - z_9 \\ z_1 - z_3 + z_4 - z_6 + z_7 - z_9 \\ z_1 + z_2 + z_3 + z_4 + z_5 + z_6 + z_7 + z_8 + z_9 \end{pmatrix},$$

$$\hat{a} = \frac{1}{4}(z_1 - z_3 - z_7 + z_9), \quad \hat{c} = \frac{1}{6}(z_1 - z_3 + z_4 - z_6 + z_7 - z_9),$$

$$\hat{b} = \frac{1}{6}(z_1 + z_2 + z_3 - z_7 - z_8 - z_9),$$

$$\hat{d} = \frac{1}{8}(z_1 + z_2 + z_3 + z_4 + z_5 + z_6 + z_7 + z_8 + z_9).$$

## 2.5 Свойства МНК-оценки

Здесь и далее будем рассматривать классическую модель наблюдения (2.5). В соответствии с этой моделью вектор параметров  $\boldsymbol{\theta}$  считается неизвестной неслучайной (детерминированной) величиной.

Оценка некоторого вектора параметров называется несмещённой, если её математическое ожидание равно истинному значению вектора параметров. Докажем, что МНК-оценка  $\hat{\boldsymbol{\theta}}$  является *несмещённой* оценкой вектора параметров, то есть:

$$\mathbf{M}\hat{\boldsymbol{\theta}} = \boldsymbol{\theta}. \quad (2.10)$$

$$\begin{aligned} \mathbf{M}\hat{\boldsymbol{\theta}} &= \mathbf{M}(\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{Y} = \mathbf{M}(\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T (\mathbf{F}\boldsymbol{\theta} + \boldsymbol{\xi}) = \\ &= (\mathbf{F}^T \mathbf{F})^{-1} (\mathbf{F}^T \mathbf{F}) \mathbf{M}\boldsymbol{\theta} + (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{M}\boldsymbol{\xi} = \boldsymbol{\theta} + \mathbf{0} = \boldsymbol{\theta}. \end{aligned}$$

Докажем, что дисперсионная матрица МНК-оценки параметров (то есть ковариационная матрица ошибки оценивания  $\boldsymbol{\varepsilon} = \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}$ ) равна  $\sigma^2 (\mathbf{F}^T \mathbf{F})^{-1}$ , где  $\sigma^2$  – дисперсия шума наблюдения,  $\mathbf{F}$  – матрица эксперимента.

$$\begin{aligned}
\mathbf{D}\boldsymbol{\varepsilon} &= \mathbf{D}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) = \mathbf{D}\hat{\boldsymbol{\theta}} = \mathbf{D}\left[(\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T\mathbf{Y}\right] = \\
&= \mathbf{D}\left[(\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T(\mathbf{F}\boldsymbol{\theta} + \boldsymbol{\xi})\right] = \mathbf{D}\left[(\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T\boldsymbol{\xi}\right] = \\
&= \left((\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T\right)\mathbf{D}\boldsymbol{\xi}\left((\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T\right)^T = \\
&= \left((\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T\right)\sigma^2\mathbf{I}_N\left((\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T\right)^T = \sigma^2(\mathbf{F}^T\mathbf{F})^{-1}.
\end{aligned}$$

Индивидуальные дисперсии оценок параметров (дисперсии ошибок оценивания) определяются диагональными элементами матрицы:

$$\mathbf{D}\hat{\theta}_i = \mathbf{D}(\hat{\theta}_i - \theta_i) = \sigma^2 \left[ (\mathbf{F}^T\mathbf{F})^{-1} \right]_{ii}, \quad i = 1, 2, 3, \dots, m.$$

МНК-оценка параметров является *эффективной* оценкой (то есть обладающей минимальной дисперсией ошибки) в классе линейных несмещённых оценок, о чём утверждается в следующей теореме Гаусса–Маркова.

## 2.6 Теорема Гаусса–Маркова

*Линейной оценкой* вектора параметров  $\boldsymbol{\theta}$  называется оценка  $\hat{\boldsymbol{\theta}} = \mathbf{C}\mathbf{Y}$ , где  $\mathbf{C}$  – произвольная матрица размера  $M \times N$ . *Линейной несмещённой оценкой* вектора параметров  $\boldsymbol{\theta}$  называется такая линейная оценка  $\hat{\boldsymbol{\theta}}$ , для которой  $\mathbf{M}\hat{\boldsymbol{\theta}} = \boldsymbol{\theta}$  (несмещённость МНК-оценки доказана выше в п.2.5). *Наилучшей линейной несмещённой оценкой (НЛН-оценкой)* вектора параметров  $\boldsymbol{\theta}$  называется такая оценка  $\hat{\boldsymbol{\theta}}$ , для кото-

рой  $\mathbf{D}\hat{\boldsymbol{\theta}} \leq \mathbf{D}\hat{\boldsymbol{\theta}}$ , где  $\hat{\boldsymbol{\theta}}$  – произвольная линейная несмещённая оценка вектора параметров  $\boldsymbol{\theta}$ .

Далее будем использовать обозначения  $\mathbf{A} < \mathbf{B}$ ,  $\mathbf{B} > \mathbf{A}$ ,  $\mathbf{B} - \mathbf{A} > \mathbf{0}$  ( $\mathbf{A} \leq \mathbf{B}$ ,  $\mathbf{B} \geq \mathbf{A}$ ,  $\mathbf{B} - \mathbf{A} \geq \mathbf{0}$ ) для указания факта положительной (неотрицательной) определённости матрицы  $\mathbf{B} - \mathbf{A}$ .

*Теорема Гаусса–Маркова.* Для классической модели наблюдения (2.5) МНК-оценка  $\hat{\boldsymbol{\theta}} = (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{Y}$  является НЛН-оценкой  $\hat{\boldsymbol{\theta}}$  вектора параметров  $\boldsymbol{\theta}$ .

*Доказательство:* Пусть  $\hat{\boldsymbol{\theta}} = \mathbf{C}\mathbf{Y}$  – произвольная линейная оценка вектора параметров  $\boldsymbol{\theta}$ . Тогда  $\mathbf{M}\hat{\boldsymbol{\theta}} = \mathbf{M}\mathbf{C}\mathbf{Y} = \mathbf{M}\mathbf{C}(\mathbf{F}\boldsymbol{\theta} + \boldsymbol{\xi}) = \mathbf{C}\mathbf{F}\boldsymbol{\theta}$ . Условие несмещённости оценки  $\hat{\boldsymbol{\theta}}$ :  $\mathbf{C}\mathbf{F}\boldsymbol{\theta} = \boldsymbol{\theta}$ . Отсюда следует, что для несмещённости любой линейной оценки необходимо и достаточно, чтобы выполнялось условие:  $\mathbf{C}\mathbf{F} = \mathbf{I}_m$ .

Если  $\mathbf{C} = \mathbf{C}_0 = (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T$ , то получаем МНК-оценку  $\hat{\boldsymbol{\theta}} = \mathbf{C}_0 \mathbf{Y} = (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{Y}$  с дисперсией ошибки  $\sigma^2 (\mathbf{F}^T \mathbf{F})^{-1}$  (см. п.2.5). Обозначив  $\mathbf{A} = \mathbf{C} - \mathbf{C}_0$ , заметим, что  $\mathbf{A}\mathbf{C}_0^T = \mathbf{0}$ :

$$\begin{aligned} \mathbf{A}\mathbf{C}_0^T &= \mathbf{A} \left( (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \right)^T = (\mathbf{C} - \mathbf{C}_0) \mathbf{F} (\mathbf{F}^T \mathbf{F})^{-1} = (\mathbf{C}\mathbf{F} - \mathbf{C}_0 \mathbf{F}) (\mathbf{F}^T \mathbf{F})^{-1} = \\ &= (\mathbf{I}_m - \mathbf{I}_m) (\mathbf{F}^T \mathbf{F})^{-1} = \mathbf{0}. \end{aligned}$$

Аналогично  $\mathbf{C}_0 \mathbf{A}^T = \mathbf{0}$ .

$$\begin{aligned}
\mathbf{D}\hat{\boldsymbol{\theta}} &= \mathbf{D}(\mathbf{C}\mathbf{Y}) = \mathbf{C} \cdot \mathbf{D}\mathbf{Y} \cdot \mathbf{C}^T = \mathbf{C} \cdot \mathbf{D}(\mathbf{F}\boldsymbol{\theta} + \boldsymbol{\xi}) \cdot \mathbf{C}^T = \mathbf{C} \cdot \mathbf{D}\boldsymbol{\xi} \cdot \mathbf{C}^T = \\
&= \mathbf{C} \cdot \sigma^2 \mathbf{I}_N \cdot \mathbf{C}^T = \sigma^2 \mathbf{C}\mathbf{C}^T = \sigma^2 (\mathbf{C}_0 + \mathbf{A})(\mathbf{C}_0 + \mathbf{A})^T = \\
&= \sigma^2 (\mathbf{C}_0 \mathbf{C}_0^T + \mathbf{C}_0 \mathbf{A}^T + \mathbf{A}\mathbf{C}_0^T + \mathbf{A}\mathbf{A}^T) = \\
&= \sigma^2 \left( (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \left( (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \right)^T + \mathbf{A}\mathbf{A}^T \right) = \\
&= \sigma^2 \left( (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{F} (\mathbf{F}^T \mathbf{F})^{-1} + \mathbf{A}\mathbf{A}^T \right) = \\
&= \sigma^2 (\mathbf{F}^T \mathbf{F})^{-1} + \sigma^2 \mathbf{A}\mathbf{A}^T = \mathbf{D}\hat{\boldsymbol{\theta}}^{MНК} + \sigma^2 \mathbf{A}\mathbf{A}^T \geq \mathbf{D}\hat{\boldsymbol{\theta}}^{MНК},
\end{aligned}$$

так как  $\mathbf{D}\hat{\boldsymbol{\theta}}^{MНК} = \sigma^2 (\mathbf{F}^T \mathbf{F})^{-1}$ ,  $\mathbf{A}\mathbf{C}_0^T = \mathbf{C}_0 \mathbf{A}^T = \mathbf{0}$ , а матрица  $\mathbf{A}\mathbf{A}^T$  неотрицательно определена.

Мы доказали, что произвольная линейная несмещённая оценка имеет не меньшую дисперсию, чем оценка МНК, что означает эффективность оценки МНК.

## 2.7 Оценивание дисперсии шума наблюдения

Здесь будем рассматривать классическую модель линейной регрессии *при гауссовом шуме наблюдения* (2.6). Докажем, что статистика

$$s^2 = \frac{R_0^2}{N - m}$$

является несмещённой оценкой дисперсии шума наблюдения  $\sigma^2$ , то есть  $\mathbf{M}s^2 = \sigma^2$ . Здесь обозначено, как и ранее:  $N$  – число опытов (наблюдений),  $m$  – число оцениваемых параметров,  $R_0^2$  – остаточная сумма квадратов ошибок:

$$R_0^2 = (\mathbf{Y} - \mathbf{F}\hat{\boldsymbol{\theta}})^T (\mathbf{Y} - \mathbf{F}\hat{\boldsymbol{\theta}}) = \sum_{i=1}^N \left[ y^i - (f_1^i \hat{\theta}_1 + f_2^i \hat{\theta}_2 + f_3^i \hat{\theta}_3 + \dots + f_m^i \hat{\theta}_m) \right]^2.$$

Доказательство основано на том факте, что статистика  $R_0^2/\sigma^2$  имеет распределение хи-квадрат с  $(N-m)$  степенями свободы и её математическое ожидание равно  $(N-m)$ . Доказательство распределения статистики  $R_0^2/\sigma^2$  по закону хи-квадрат приведено в [2].

Находим

$$\mathbf{M}s^2 = \frac{1}{N-m} \mathbf{M}R_0^2 = \frac{\sigma^2}{N-m} \mathbf{M} \left( \frac{R_0^2}{\sigma^2} \right) = \frac{\sigma^2}{N-m} \mathbf{M} \chi_{N-m}^2 = \sigma^2$$

(см. приложение А4, формула (А7)).

## 2.8 Интервальные оценки параметров линейной регрессионной модели и проверка гипотез о коэффициентах регрессии

Для классической модели с гауссовским шумом наблюдения (2.6) вектор оценки  $\hat{\boldsymbol{\theta}}$  имеет многомерное нормальное распределение с математическим ожиданием  $\mathbf{M}\hat{\boldsymbol{\theta}} = \boldsymbol{\theta}$  и дисперсионной матрицей  $\mathbf{D}\hat{\boldsymbol{\theta}} = \sigma^2 (\mathbf{F}^T \mathbf{F})^{-1}$  (доказано выше в п.2.5):  $\hat{\boldsymbol{\theta}} \sim \mathbf{N}(\boldsymbol{\theta}, \sigma^2 (\mathbf{F}^T \mathbf{F})^{-1})$ . Вектор  $\hat{\boldsymbol{\theta}}$  имеет многомерное нормальное распределение, так как  $\hat{\boldsymbol{\theta}} = (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{Y} = (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T (\mathbf{F}\boldsymbol{\theta} + \boldsymbol{\xi})$  – линейное преобразование нормального случайного вектора  $\boldsymbol{\xi}$ . Поэтому индивидуальные оценки параметров регрессии имеют одномерные нормальные распределения:

$$\hat{\theta}_i \sim N\left(\theta_i, \sigma^2 \left[ (\mathbf{F}^T \mathbf{F})^{-1} \right]_{ii}\right), \quad i = 1, 2, \dots, m,$$

то есть статистика  $\frac{\hat{\theta}_i - \theta_i}{\sigma \sqrt{\left[ (\mathbf{F}^T \mathbf{F})^{-1} \right]_{ii}}}$  имеет стандартное нормальное распределение  $N(0,1)$ .

Докажем, что статистика  $t = \frac{\hat{\theta}_i - \theta_i}{s \sqrt{\left[ (\mathbf{F}^T \mathbf{F})^{-1} \right]_{ii}}}$ ,  $i = 1, 2, \dots, m$ ,

где  $s^2 = \frac{R_0^2}{N-m}$ , имеет распределение Стьюдента с  $(N-m)$  степенями свободы:  $t \sim t(N-m)$ .

*Доказательство.* Так как статистика  $R_0^2/\sigma^2$  имеет хи-квадрат распределение с  $(N-m)$  степенями свободы, то в соответствии с определением (А.9) получаем:

$$\frac{\hat{\theta}_i - \theta_i}{s \sqrt{\left[ (\mathbf{F}^T \mathbf{F})^{-1} \right]_{ii}}} = \frac{\frac{\hat{\theta}_i - \theta_i}{\sigma} \left[ (\mathbf{F}^T \mathbf{F})^{-1} \right]_{ii}^{-1/2}}{\sqrt{\frac{1}{N-m} \frac{R_0^2}{\sigma^2}}} \sim \frac{N(0;1)}{\sqrt{\frac{\chi_{N-m}^2}{N-m}}} = t(N-m), \quad (2.11)$$

то есть статистика  $t = \frac{\hat{\theta}_i - \theta_i}{s \sqrt{\left[ (\mathbf{F}^T \mathbf{F})^{-1} \right]_{ii}}}$  имеет распределение

Стьюдента с  $(N-m)$  степенями свободы, что и требовалось доказать. При доказательстве учитывалось, что случайные величины

$\frac{\hat{\theta}_i - \theta_i}{\sigma \sqrt{\left[ (\mathbf{F}^T \mathbf{F})^{-1} \right]_{ii}}}$  и  $R_0^2$  являются независимыми.



*Доказательство:* Сначала докажем, что случайные векторы невязки  $(\mathbf{Y} - \mathbf{F}\hat{\boldsymbol{\theta}})$  и оценки  $\hat{\boldsymbol{\theta}}$  не коррелированы.

$$\begin{aligned}
\mathbf{M}\left[(\mathbf{Y} - \mathbf{F}\hat{\boldsymbol{\theta}})\hat{\boldsymbol{\theta}}^T\right] &= \mathbf{M}\left[\left(\mathbf{Y} - \mathbf{F}(\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T\mathbf{Y}\right)\left((\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T\mathbf{Y}\right)^T\right] = \\
&= \mathbf{M}\left[\left(\mathbf{I}_N - \mathbf{F}(\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T\right)\mathbf{Y}\mathbf{Y}^T\mathbf{F}(\mathbf{F}^T\mathbf{F})^{-1}\right] = \\
&= \mathbf{M}\left[\left(\mathbf{I}_N - \mathbf{F}(\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T\right)(\mathbf{F}\boldsymbol{\theta} + \boldsymbol{\xi})(\mathbf{F}\boldsymbol{\theta} + \boldsymbol{\xi})^T\mathbf{F}(\mathbf{F}^T\mathbf{F})^{-1}\right] = \\
&= \left(\mathbf{I}_N - \mathbf{F}(\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T\right)\mathbf{M}(\boldsymbol{\xi}\boldsymbol{\xi}^T)\mathbf{F}(\mathbf{F}^T\mathbf{F})^{-1} = \\
&= \left(\mathbf{I}_N - \mathbf{F}(\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T\right)\mathbf{I}_N\mathbf{F}(\mathbf{F}^T\mathbf{F})^{-1} = \left(\mathbf{I}_N - \mathbf{F}(\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T\right)\mathbf{F}(\mathbf{F}^T\mathbf{F})^{-1} = \\
&= \mathbf{F}(\mathbf{F}^T\mathbf{F})^{-1} - \mathbf{F}(\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T\mathbf{F}(\mathbf{F}^T\mathbf{F})^{-1} = \mathbf{0},
\end{aligned}$$

$$\begin{aligned}
\mathbf{M}(\mathbf{Y} - \mathbf{F}\hat{\boldsymbol{\theta}}) &= \mathbf{M}(\mathbf{Y} - \mathbf{F}(\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T\mathbf{Y}) = \left(\mathbf{I}_N - \mathbf{F}(\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T\right)\mathbf{M}\mathbf{Y} = \\
&= \left(\mathbf{I}_N - \mathbf{F}(\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T\right)\mathbf{M}(\mathbf{F}\boldsymbol{\theta} + \boldsymbol{\xi}) = \left(\mathbf{I}_N - \mathbf{F}(\mathbf{F}^T\mathbf{F})^{-1}\mathbf{F}^T\right)\mathbf{F}\boldsymbol{\theta} = \mathbf{F}\boldsymbol{\theta} - \mathbf{F}\boldsymbol{\theta} = \mathbf{0}, \\
\mathbf{R}(\mathbf{Y} - \mathbf{F}\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\theta}}) &= \mathbf{M}\left[(\mathbf{Y} - \mathbf{F}\hat{\boldsymbol{\theta}})\hat{\boldsymbol{\theta}}^T\right] - \mathbf{M}(\mathbf{Y} - \mathbf{F}\hat{\boldsymbol{\theta}})\mathbf{M}\hat{\boldsymbol{\theta}}^T = \mathbf{0} - \mathbf{0} \cdot \mathbf{M}\hat{\boldsymbol{\theta}}^T = \mathbf{0}.
\end{aligned}$$

Так как случайные векторы  $\mathbf{Y} - \mathbf{F}\hat{\boldsymbol{\theta}}$  и  $\hat{\boldsymbol{\theta}}$  являются нормально распределёнными, то из их некоррелированности следует их независимость. Тогда также независимыми будут их функции:

$$R_0^2 = (\mathbf{Y} - \mathbf{F}\hat{\boldsymbol{\theta}})^T (\mathbf{Y} - \mathbf{F}\hat{\boldsymbol{\theta}}) \text{ и } \frac{\hat{\theta}_i - \theta_i}{\sigma \sqrt{\left[(\mathbf{F}^T\mathbf{F})^{-1}\right]_{ii}}}.$$

Из (2.11) следует, что соотношение для определения индивидуальных интервальных оценок параметров  $\theta_1, \theta_2, \dots, \theta_m$  имеет вид:

$$\Pr\left(\left|\hat{\theta}_i - \theta_i\right| < t_{N-m}^{(1+p)/2} s \sqrt{\left[(\mathbf{F}^T\mathbf{F})^{-1}\right]_{ii}}\right) = p,$$

$$\Pr\left(\hat{\theta}_i - t_{N-m}^{(1+p)/2} s \sqrt{\left[(\mathbf{F}^T \mathbf{F})^{-1}\right]_{ii}} < \theta_i < \hat{\theta}_i + t_{N-m}^{(1+p)/2} s \sqrt{\left[(\mathbf{F}^T \mathbf{F})^{-1}\right]_{ii}}\right) = p,$$

$$\left(\hat{\theta}_i - t_{N-m}^{(1+p)/2} s \sqrt{\left[(\mathbf{F}^T \mathbf{F})^{-1}\right]_{ii}}, \hat{\theta}_i + t_{N-m}^{(1+p)/2} s \sqrt{\left[(\mathbf{F}^T \mathbf{F})^{-1}\right]_{ii}}\right) - \text{довери-}$$

тельный интервал для параметра  $\theta_i$ ;  $p$  – доверительная веро-

ятность;  $s^2 = \frac{R_0^2}{N-m}$ ,  $t_{N-m}^{(1+p)/2}$  – квантиль на уровне  $(1+p)/2$

$t$ -распределения Стьюдента с  $(N-m)$  степенями свободы.

Для проверки статистической значимости параметра регрессии  $\theta_i$  (то есть его отличия от нуля) также воспользуемся  $t$ -статистикой Стьюдента. При выполнении нулевой гипотезы

$$(H_0: \theta_i = 0) \text{ статистика } t = \frac{\hat{\theta}_i}{s \sqrt{\left[(\mathbf{F}^T \mathbf{F})^{-1}\right]_{ii}}} \text{ имеет распределение}$$

Стьюдента с  $(N-m)$  степенями свободы.

Задавшись уровнем значимости  $\alpha$ , определим границы критической области из соотношения

$$\Pr\left(\left|\hat{\theta}_i\right| < t_{N-m}^{1-\alpha/2} s \sqrt{\left[(\mathbf{F}^T \mathbf{F})^{-1}\right]_{ii}}\right) = 1 - \alpha,$$

где  $s^2 = \frac{R_0^2}{N-m}$ ,  $t_{N-m}^{1-\alpha/2}$  – квантиль на уровне  $(1-\alpha/2)$

$t$ -распределения Стьюдента с  $(N-m)$  степенями свободы.

### 3 ДИСПЕРСИОННЫЙ АНАЛИЗ

*Дисперсионным анализом* называется статистический метод выявления влияния отдельных факторов на результат эксперимента путём исследования *значимости* различий в средних значениях. В литературе, в частности, при описании математических пакетов часто встречается обозначение *ANOVA* (от английского *ANalysis Of VAriance* – дисперсионный анализ). Фундаментальная концепция дисперсионного анализа предложена английским статистиком Рональдом Фишером в 1920 году. Изобретённая им *SS*-технология статистического анализа (от английского *Sum of Squares* – сумма квадратов) является достаточно удобной и наглядной. Эта технология основана на математической статистике с использованием вероятностных распределений *хи-квадрат* и *Фишера–Снедекора* (приложение А). Мы будем использовать полученные ранее результаты для линейного регрессионного анализа при гауссовском шуме наблюдения (см. главу 2).

#### 3.1 Вероятностные распределения хи-квадрат и Фишера

Если независимые случайные величины  $\xi_1, \xi_2, \dots, \xi_k$  имеют нормальное распределение с нулевым математическим ожиданием и единичной дисперсией, то случайная величина

$\chi_k^2 = \sum_{i=1}^k \xi_i^2$  имеет *хи-квадрат*-распределение с  $k$  степенями

свободы:

$$\sum_{i=1}^k \xi_i^2 \sim \chi^2(k). \quad (3.1)$$

Если нормальный случайный вектор  $\boldsymbol{\eta} = (\eta_1 \eta_2 \eta_3 \dots \eta_m)^T$  имеет нулевое математическое ожидание и дисперсионную матрицу  $\mathbf{D}_\eta$ , то случайная величина  $\chi_m^2 = \boldsymbol{\eta}^T \mathbf{D}_\eta^{-1} \boldsymbol{\eta}$  имеет *хи-квадрат*-распределение с  $m$  степенями свободы:

$$\chi_m^2 = \boldsymbol{\eta}^T \mathbf{D}_\eta^{-1} \boldsymbol{\eta} \sim \chi^2(m). \quad (3.2)$$

Если независимые случайные величины  $\chi_k^2$  и  $\chi_l^2$  имеют  $\chi^2$ -распределение соответственно с  $k$  и  $l$  степенями свободы, то случайная величина  $F_{k,l} = \frac{\chi_k^2/k}{\chi_l^2/l}$  имеет распределение Фишера ( $F$ -распределение, распределение Фишера–Снедекора) с  $k$  и  $l$  степенями свободы:

$$\frac{\chi_k^2/k}{\chi_l^2/l} \sim F(k, l). \quad (3.3)$$

### 3.2 Статистика Фишера

Статистика

$$F_{m, N-m} = \frac{1}{ms^2} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})^T \mathbf{F}^T \mathbf{F} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \sim F(m, N-m) \quad (3.4)$$

имеет распределение Фишера с  $m$  и  $(N-m)$  степенями свободы.

*Доказательство.* Из соотношения  $\frac{R_0^2}{\sigma^2} \sim \chi^2(N-m)$  (см.

п.2.7) следует, что величина  $\frac{(N-m)s^2}{\sigma^2}$  имеет хи-квадрат рас-

пределение с  $(N-m)$  степенями свободы. Из гауссовского распределения вектора оценки  $\hat{\boldsymbol{\theta}} \sim N\left(\boldsymbol{\theta}, \sigma^2 (\mathbf{F}^T \mathbf{F})^{-1}\right)$  следует,

что случайный вектор  $\boldsymbol{\eta} = \frac{\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}}{\sigma}$  является гауссовским с нуле-

вым математическим ожиданием и дисперсионной матрицей

$\mathbf{D}_{\boldsymbol{\eta}} = (\mathbf{F}^T \mathbf{F})^{-1}$ . Из (3.2) следует, что  $\boldsymbol{\eta}^T \left[ (\mathbf{F}^T \mathbf{F})^{-1} \right]^{-1} \boldsymbol{\eta} \sim \chi^2(m)$ , то

есть величина  $\left( \frac{\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}}{\sigma} \right)^T \mathbf{F}^T \mathbf{F} \left( \frac{\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}}{\sigma} \right)$  имеет распределение хи-

квадрат с  $m$  степенями свободы. Так как остаточная сумма

квадратов  $R_0^2$  и вектор оценки  $\hat{\boldsymbol{\theta}}$  вероятностно независимы

(доказывается в [2]), то также будут независимыми и функции

от них:  $\frac{(N-m)s^2}{\sigma^2}$  и  $\left( \frac{\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}}{\sigma} \right)^T \mathbf{F}^T \mathbf{F} \left( \frac{\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}}{\sigma} \right)$ . В соответствии с

(3.3) отношение двух независимых случайных величин, имею-

щих распределения хи-квадрат и нормированных на их число

степеней свободы, подчиняется закону распределения Фишера:

$$\frac{\left( \frac{\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}}{\sigma} \right)^T \mathbf{F}^T \mathbf{F} \left( \frac{\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}}{\sigma} \right) / m}{\frac{(N-m)s^2}{\sigma^2} / (N-m)} \sim F(m, N-m),$$

то есть  $\frac{1}{ms^2}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})^T \mathbf{F}^T \mathbf{F} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \sim F(m, N - m)$ .

Статистика Фишера (3.4) используется для проверки гипотезы о значимости линейной регрессионной модели в целом относительно нулевой гипотезы  $H_0: \boldsymbol{\theta} = \mathbf{0}$  (то есть  $H_0: \theta_1 = \theta_2 = \theta_3 = \dots = \theta_m = 0$ ). Регрессионная модель считается значимой ( $\boldsymbol{\theta} \neq \mathbf{0}$ ) на уровне значимости  $\alpha$ , если рассчитанное значение статистики

$$F = \frac{1}{ms^2} \hat{\boldsymbol{\theta}}^T \mathbf{F}^T \mathbf{F} \hat{\boldsymbol{\theta}}$$

превышает критическое значение  $F_{m, N-m}^{1-\alpha}$  – квантиля на уровне  $(1-\alpha)$  распределения Фишера с  $m$  и  $(N-m)$  степенями свободы.

### 3.3 МНК-оценка линейной векторной функции коэффициентов регрессии

Линейная векторная функция параметров регрессии задается как

$$\boldsymbol{\tau} = \mathbf{T}\boldsymbol{\theta},$$

где  $\mathbf{T}$  – произвольная матрица полного ранга линейного преобразования параметров, имеющая размер  $k \times m$ ,  $1 \leq k \leq m$  ( $\text{rank } \mathbf{T} = k$ ).

Тогда НЛН-оценка вектора  $\boldsymbol{\tau}$  может быть определена через МНК-оценку параметров  $\hat{\boldsymbol{\theta}}$  (2.9):

$$\hat{\boldsymbol{\tau}} = \mathbf{T}\hat{\boldsymbol{\theta}} = \mathbf{T}(\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{Y}. \quad (3.5)$$

Оценка  $\hat{\boldsymbol{\tau}}$  (3.5) является несмещённой ( $\mathbf{M}\hat{\boldsymbol{\tau}} = \boldsymbol{\tau}$ ), а её дисперсионная матрица (дисперсионная матрица ошибки оценивания) равна

$$\mathbf{D}(\hat{\boldsymbol{\tau}} - \boldsymbol{\tau}) = \mathbf{D}\hat{\boldsymbol{\tau}} = \sigma^2 \mathbf{T}(\mathbf{F}^T \mathbf{F})^{-1} \mathbf{T}^T.$$

Основное назначение матрицы  $\mathbf{T}$  – выделять статистически значимые компоненты регрессионной модели при проверке гипотез. Например, при  $m = 5$ ,  $k = 3$  матрица

$$\mathbf{T} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

выделяет из первоначального набора из пяти параметров  $\boldsymbol{\theta} = (\theta_1 \ \theta_2 \ \theta_3 \ \theta_4 \ \theta_5)^T$  набор из трёх параметров  $\boldsymbol{\tau} = (\theta_1 \ \theta_2 \ \theta_4)^T$  для анализа.

Другое назначение – приведение гипотезы о проверке параметров к нулевой гипотезе. Например, если требуется проверить гипотезу о равенстве четырёх параметров регрессии

$$H_0 : \theta_1 = \theta_2 = \theta_3 = \theta_4, \text{ то есть } \theta_1 - \theta_2 = 0, \theta_1 - \theta_3 = 0, \theta_1 - \theta_4 = 0,$$

то эквивалентная нулевая гипотеза  $H'_0 : \boldsymbol{\tau} = \mathbf{0}$  получается при использовании линейного преобразования параметров:

$$\boldsymbol{\tau} = \mathbf{T}\boldsymbol{\theta} = \begin{bmatrix} 1 & -1 & 0 & 0 \\ 1 & 0 & -1 & 0 \\ 1 & 0 & 0 & -1 \end{bmatrix} \begin{pmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \end{pmatrix} = \mathbf{0}, \text{ где } \mathbf{T} = \begin{bmatrix} 1 & -1 & 0 & 0 \\ 1 & 0 & -1 & 0 \\ 1 & 0 & 0 & -1 \end{bmatrix}.$$

Ниже без доказательства приведены вероятностные распределения и основные свойства оценки  $\hat{\boldsymbol{\tau}}$ . Доказательства аналогичны приведённым выше.

- Оценка  $\hat{\boldsymbol{\tau}} = \mathbf{T}(\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{Y}$  является НЛН-оценкой  $\boldsymbol{\tau}$ .

- Вектор оценки  $\hat{\boldsymbol{\tau}}$  имеет многомерное нормальное распределение с математическим ожиданием  $\mathbf{M}\hat{\boldsymbol{\tau}} = \boldsymbol{\tau}$  и дисперсионной матрицей  $\mathbf{D}\hat{\boldsymbol{\tau}} = \sigma^2 \mathbf{T}(\mathbf{F}^T \mathbf{F})^{-1} \mathbf{T}^T$ :

$$\hat{\boldsymbol{\tau}} \sim \mathbf{N}\left(\boldsymbol{\tau}, \sigma^2 \mathbf{T}(\mathbf{F}^T \mathbf{F})^{-1} \mathbf{T}^T\right).$$

- Остаточная сумма квадратов  $R_0^2$  и вектор оценки  $\hat{\boldsymbol{\tau}}$  вероятностно независимы.
- Статистика

$$F = \frac{1}{ks^2} (\hat{\boldsymbol{\tau}} - \boldsymbol{\tau})^T \left[ \mathbf{T}(\mathbf{F}^T \mathbf{F})^{-1} \mathbf{T}^T \right]^{-1} (\hat{\boldsymbol{\tau}} - \boldsymbol{\tau}), \quad (3.6)$$

имеет распределение Фишера с  $k$  и  $(N - m)$  степенями свободы:  $F \sim F(k, N - m)$ , где  $k$  – число строк в матрице  $\mathbf{T}$ .

### 3.4 Однофакторный двухуровневый дисперсионный анализ

Для начала рассмотрим простейший случай дисперсионного анализа, относящийся к исследованию значимости влияния какого-либо одного фактора на выходной фактор. На этом примере рассмотрим основные понятия и закономерности дисперсионного анализа.

*Однофакторный двухуровневый дисперсионный анализ* при одинаковом числе наблюдений решает простейшую задачу исследования значимости различия между средними значениями.

Рассмотрим наблюдения двух независимых выборок  $\{y_1^1, y_1^2, y_1^3, \dots, y_1^n\}$  и  $\{y_2^1, y_2^2, y_2^3, \dots, y_2^n\}$  из двух нормальных генеральных совокупностей  $Y_1$  и  $Y_2$  с одинаковыми дисперсиями  $\sigma^2$ , но различными математическими ожиданиями  $a$  и  $b$ :



$$y_1^i = a + \xi_1^i, \quad i = \overline{1, n},$$

$$y_2^i = b + \xi_2^i, \quad i = \overline{1, n},$$

где  $\xi_1^i$  и  $\xi_2^i$  можно интерпретировать как шум наблюдения с нулевым математическим ожиданием и дисперсией  $\sigma^2$ . Эти уравнения можно представить как результат наблюдения в классической линейной регрессионной модели (2.6)  $\mathbf{Y} = \mathbf{F}\boldsymbol{\theta} + \boldsymbol{\xi}$  при  $N = 2n$ ,  $m = 2$ ,

$$\boldsymbol{\theta} = (a \ b)^T, \quad \mathbf{Y} = (y_1^1 \ y_1^2 \ y_1^3 \ \dots \ y_1^n \ y_2^1 \ y_2^2 \ y_2^3 \ \dots \ y_2^n)^T,$$

$$\boldsymbol{\xi} = (\xi_1^1 \ \xi_1^2 \ \xi_1^3 \ \dots \ \xi_1^n \ \xi_2^1 \ \xi_2^2 \ \xi_2^3 \ \dots \ \xi_2^n)^T,$$

$$\mathbf{F} = \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}, \quad \mathbf{D}\boldsymbol{\xi} = \sigma^2 \mathbf{I}_{2n}.$$

Тогда

$$\mathbf{F}^T \mathbf{F} = n \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = n \mathbf{I}_2, \quad (\mathbf{F}^T \mathbf{F})^{-1} = \frac{1}{n} \mathbf{I}_2, \quad \mathbf{F}^T \mathbf{Y} = \left( \sum_{i=1}^n y_1^i \quad \sum_{i=1}^n y_2^i \right)^T,$$

$$\hat{\boldsymbol{\theta}} = (\hat{a} \ \hat{b})^T = (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{Y} = (\bar{y}_1 \ \bar{y}_2)^T, \quad \text{то есть } \hat{a} = \bar{y}_1, \quad \hat{b} = \bar{y}_2,$$

где обозначено:  $\bar{y}_1 = \frac{1}{n} \sum_{i=1}^n y_1^i$ ,  $\bar{y}_2 = \frac{1}{n} \sum_{i=1}^n y_2^i$ .

Для проверки гипотезы  $H_0 : a = b$  используется статистика Фишера (3.6). При линейном преобразовании вектора параметров

$$\mathbf{T} = [\mathbf{1} \ -\mathbf{1}], \quad k = 1$$

получаем, что гипотеза  $H_0 : a = b$  эквивалентна гипотезе  $H'_0 : \boldsymbol{\tau} = \mathbf{0}$ ,

$$\left[ \mathbf{T}(\mathbf{F}^T \mathbf{F})^{-1} \mathbf{T}^T \right]^{-1} = \frac{n}{2},$$

$$\hat{\boldsymbol{\tau}} = \mathbf{T} \hat{\boldsymbol{\theta}} = \mathbf{T}(\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{Y} = \bar{y}_1 - \bar{y}_2.$$

Остаточная сумма квадратов ошибок:

$$\begin{aligned} R_0^2 &= \sum_{i=1}^n (y_1^i - \bar{y}_1)^2 + \sum_{i=1}^n (y_2^i - \bar{y}_2)^2 = \\ &= \sum_{i=1}^n (y_1^i)^2 - 2\bar{y}_1 \sum_{i=1}^n y_1^i + n\bar{y}_1^2 + \sum_{i=1}^n (y_2^i)^2 - 2\bar{y}_2 \sum_{i=1}^n y_2^i + n\bar{y}_2^2 = \\ &= n(\bar{y}_1^2 - \bar{y}_1^2) + n(\bar{y}_2^2 - \bar{y}_2^2). \end{aligned}$$

Несмещённая оценка дисперсии шума наблюдения:

$$s^2 = \frac{n(\bar{y}_1^2 - \bar{y}_1^2) + n(\bar{y}_2^2 - \bar{y}_2^2)}{2n - 2}.$$

Статистика Фишера (3.6):

$$\begin{aligned} F &= \frac{1}{ks^2} \hat{\boldsymbol{\tau}}^T \left[ \mathbf{T}(\mathbf{F}^T \mathbf{F})^{-1} \mathbf{T}^T \right]^{-1} \hat{\boldsymbol{\tau}} = \frac{\frac{n}{2}(\bar{y}_1 - \bar{y}_2)^2}{\frac{n(\bar{y}_1^2 - \bar{y}_1^2) + n(\bar{y}_2^2 - \bar{y}_2^2)}{2n - 2}} = \\ &= \frac{(n-1)(\bar{y}_1 - \bar{y}_2)^2}{(\bar{y}_1^2 - \bar{y}_1^2) + (\bar{y}_2^2 - \bar{y}_2^2)} \sim F(1, 2n-2). \end{aligned} \quad (3.7)$$

Обозначим через  $F_{1,2n-2}^{1-\alpha}$  квантиль на уровне  $(1-\alpha)$  распределения Фишера с  $(1, 2n-2)$  степенями свободы. Если  $F_{1,2n-2} > F_{1,2n-2}^{1-\alpha}$ , то гипотезу  $H_0$  о равенстве математических ожиданий отвергаем на уровне значимости  $\alpha$ , то есть считаем, что различие математических ожиданий существенно, значимо ( $a \neq b$ ). Если  $F_{1,2n-2} \leq F_{1,2n-2}^{1-\alpha}$ , то считаем, что эксперименталь-

ные данные не противоречат гипотезе  $H_0 : a = b$  (различие математических ожиданий несущественно и объясняется наличием шума наблюдения).

Найдём **общую** изменчивость результирующего фактора  $y$  (сумму квадратов отклонений значений результирующего фактора от его среднего значения  $\bar{y}$  по всем наблюдениям), которая обозначается как  $SS_{общ}$  (от английского *Sum of Squares*). Иначе  $SS_{общ}$  обозначается как  $R_1^2$ . В общем виде в рамках классической линейной регрессионной модели с  $N = 2n$  наблюдениями:

$$\begin{aligned} SS_{общ} = R_1^2 &= \sum_{i=1}^N (y^i - \bar{y})^2 = \sum_{i=1}^N (y^i)^2 - 2 \sum_{i=1}^N y^i \bar{y} + \sum_{i=1}^N \bar{y}^2 = \\ &= N \overline{y^2} - 2N\bar{y}^2 + N\bar{y}^2 = N \overline{y^2} - N\bar{y}^2 = n \overline{y_1^2} + n \overline{y_2^2} - \frac{n}{2} (\bar{y}_1 + \bar{y}_2)^2. \end{aligned}$$

Изменчивость результирующего фактора  $y$  в результате воздействия **шума наблюдения** (обозначается  $SS_{ош}$  – сумма квадратов ошибок, или  $R_0^2$ ) описывается известной нам остаточной суммой квадратов ошибок

$$\begin{aligned} SS_{ош} = R_0^2 &= \sum_{i=1}^n (y_1^i - \bar{y}_1)^2 + \sum_{i=1}^n (y_2^i - \bar{y}_2)^2 = \\ &= \sum_{i=1}^n (y_1^i)^2 - 2\bar{y}_1 \sum_{i=1}^n y_1^i + \sum_{i=1}^n \bar{y}_1^2 + \sum_{i=1}^n (y_2^i)^2 - 2\bar{y}_2 \sum_{i=1}^n y_2^i + \sum_{i=1}^n \bar{y}_2^2 = \\ &= n \left( \overline{y_1^2} - \bar{y}_1^2 \right) + n \left( \overline{y_2^2} - \bar{y}_2^2 \right). \end{aligned}$$

Изменчивость результирующего фактора  $y$  в результате различных уровней **входного фактора** (обозначается  $SS_{факт}$  – сумма квадратов отклонений значений выходного фактора, со-

ответствующих различным значениям входного фактора, от общего среднего):

$$SS_{\text{факт}} = \sum_{i=1}^n (\bar{y}_1 - \bar{y})^2 + \sum_{i=1}^n (\bar{y}_2 - \bar{y})^2.$$

Для нашего случая однофакторного двухуровневого дисперсионного анализа можно получить более простое выражение для  $SS_{\text{факт}}$ :

$$\begin{aligned} SS_{\text{факт}} &= n(\bar{y}_1 - \bar{y})^2 + n(\bar{y}_2 - \bar{y})^2 = \\ &= n\left(\bar{y}_1 - \frac{\bar{y}_1 + \bar{y}_2}{2}\right)^2 + n\left(\bar{y}_2 - \frac{\bar{y}_1 + \bar{y}_2}{2}\right)^2 = \\ &= n\left(\frac{\bar{y}_1 - \bar{y}_2}{2}\right)^2 + n\left(\frac{\bar{y}_2 - \bar{y}_1}{2}\right)^2 = \frac{n}{2}(\bar{y}_1 - \bar{y}_2)^2. \end{aligned}$$

Докажем «закон сохранения суммы квадратов»

$$SS_{\text{общ}} = SS_{\text{факт}} + SS_{\text{ош}} \quad (3.8)$$

на примере нашего частного случая:

$$\begin{aligned} SS_{\text{общ}} - SS_{\text{ош}} &= \sum_{i=1}^n (y_1^i - \bar{y})^2 + \sum_{i=1}^n (y_2^i - \bar{y})^2 - \sum_{i=1}^n (y_1^i - \bar{y}_1)^2 - \sum_{i=1}^n (y_2^i - \bar{y}_2)^2 = \\ &= \cancel{\sum_{i=1}^n (y_1^i)^2} - 2\bar{y} \sum_{i=1}^n y_1^i + n\bar{y}^2 + \cancel{\sum_{i=1}^n (y_2^i)^2} - 2\bar{y} \sum_{i=1}^n y_2^i + n\bar{y}^2 - \\ &- \cancel{\sum_{i=1}^n (y_1^i)^2} + 2\bar{y}_1 \sum_{i=1}^n y_1^i - n\bar{y}_1^2 - \cancel{\sum_{i=1}^n (y_2^i)^2} + 2\bar{y}_2 \sum_{i=1}^n y_2^i - n\bar{y}_2^2 = \\ &= -2n\bar{y}\bar{y}_1 + n\bar{y}^2 - 2n\bar{y}\bar{y}_2 + n\bar{y}^2 + n\bar{y}_1^2 + n\bar{y}_2^2 = n(\bar{y}_1 - \bar{y})^2 + n(\bar{y}_2 - \bar{y})^2 = \\ &= SS_{\text{факт}}. \end{aligned}$$

На самом деле вычисление статистики Фишера на основе соотношения (3.6) является частным случаем гораздо более

глобального соотношения, которое будет рассмотрено ниже для других моделей наблюдения. Общий подход к проверке гипотезы о значимости влияния входных факторов на результат основывается на вычислении статистики Фишера в виде:

$$F_{K,N-m} = \frac{SS_{\text{факт}}/K}{SS_{\text{ош}}/(N-m)} \sim F(K, N-m), \quad (3.9)$$

где  $N$  – общее число наблюдений,  
 $m$  – число оцениваемых параметров,  
 $K$  – число связей (ограничений), наложенных на входные факторы проверяемой гипотезой.

В нашем случае  $N=2n$ ,  $m=2$ ,  $K=1$ .

$$\begin{aligned} \frac{SS_{\text{факт}}/K}{SS_{\text{ош}}/(N-m)} &= \frac{\frac{n}{2}(\bar{y}_1 - \bar{y}_2)^2 / 1}{\left( n(\bar{y}_1^2 - \bar{y}_1^2) + n(\bar{y}_2^2 - \bar{y}_2^2) \right) / (2n-2)} = \\ &= \frac{(n-1)(\bar{y}_1 - \bar{y}_2)^2}{\left( \bar{y}_1^2 - \bar{y}_1^2 \right) + \left( \bar{y}_2^2 - \bar{y}_2^2 \right)}. \end{aligned}$$

Видим, что результат совпадает с (3.7). Важным является то, что расчёт статистики Фишера на основе общего подхода является более простым и наглядным.

### 3.5 Однофакторный многоуровневый дисперсионный анализ

В табл. 3.1 приведены данные для однофакторного многоуровневого дисперсионного анализа. Уравнения линейной регрессии в скалярном виде имеют вид:

$$y_i^k = \theta_i + \zeta_i^k, \quad i = \overline{1, m}, \quad k = \overline{1, n}, \quad (3.10)$$

где  $\theta_1, \theta_2, \theta_3, \dots, \theta_m$  – влияния значений входного фактора различных уровней на результат,  $m$  – число уровней входного фактора,  $n$  – число наблюдений на каждом уровне входного фактора (всего  $N = mn$  наблюдений),  $\xi_i^k$  – некоррелированный гауссовский шум наблюдения с нулевым математическим ожиданием и дисперсией  $\sigma^2$ .

Таблица 3.1 – Основная таблица дисперсионного анализа

Номер уровня входного фактора	Значение входного фактора	Влияние уровня входного фактора на результат	Значения выходного фактора					Среднее значение
			1	2	3	...	n	
1	$\mathbf{x}_1$	$\theta_1$	$y_1^1$	$y_1^2$	$y_1^3$	...	$y_1^n$	$\bar{y}_1$
2	$\mathbf{x}_2$	$\theta_2$	$y_2^1$	$y_2^2$	$y_2^3$	...	$y_2^n$	$\bar{y}_2$
...	...	...	...	...	...	...	...	...
$m$	$\mathbf{x}_m$	$\theta_m$	$y_m^1$	$y_m^2$	$y_m^3$	...	$y_m^n$	$\bar{y}_m$
Общее среднее: $\bar{y} = \frac{1}{mn} \sum_{k=1}^m \sum_{i=1}^n y_k^i$								$\bar{y}$

В традиционной матричной записи уравнения (3.10) имеют вид:

$$\mathbf{Y} = \mathbf{F}\boldsymbol{\theta} + \boldsymbol{\xi},$$

где  $\boldsymbol{\theta} = (\theta_1 \ \theta_2 \ \theta_3 \ \dots \ \theta_m)^T$  – вектор параметров,

$\mathbf{Y} = (y_1^1 \ y_1^2 \ \dots \ y_1^n \ y_2^1 \ y_2^2 \ \dots \ y_2^n \ \dots \ y_m^1 \ \dots \ y_m^n)^T$  – вектор наблюдений,

$\xi = (\xi_1^1 \ \xi_1^2 \ \dots \ \xi_1^n \ \xi_2^1 \ \xi_2^2 \ \dots \ \xi_2^n \ \dots \ \xi_m^1 \ \dots \ \xi_m^n)^T$  – вектор шума наблюдения, который считается гауссовским с нулевым математическим ожиданием и дисперсионной матрицей  $\mathbf{D}\xi = \sigma^2 \mathbf{I}_N$ ,  $N = mn$ ,

$$\mathbf{F}^T = \begin{bmatrix}
 \mathbf{1} & \mathbf{1} & \dots & \mathbf{1} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \dots & \dots & \dots & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} \\
 \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{1} & \mathbf{1} & \dots & \mathbf{1} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \dots & \dots & \dots & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} \\
 \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{1} & \mathbf{1} & \dots & \mathbf{1} & \dots & \dots & \dots & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} \\
 \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\
 \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \dots & \dots & \dots & \mathbf{1} & \mathbf{1} & \dots & \mathbf{1}
 \end{bmatrix} \begin{matrix} 1 \\ 2 \\ 3 \\ \dots \\ m \end{matrix}$$

$1 \quad 2 \quad \dots \quad n \qquad \qquad \qquad mn$

Тогда

$$\mathbf{F}^T \mathbf{F} = \begin{bmatrix}
 n & 0 & \dots & 0 \\
 0 & n & \dots & 0 \\
 \dots & \dots & \dots & \dots \\
 0 & 0 & \dots & n
 \end{bmatrix} = n \mathbf{I}_m,$$

$$(\mathbf{F}^T \mathbf{F})^{-1} = \begin{bmatrix}
 1/n & 0 & \dots & 0 \\
 0 & 1/n & \dots & 0 \\
 \dots & \dots & \dots & \dots \\
 0 & 0 & \dots & 1/n
 \end{bmatrix} = \frac{1}{n} \mathbf{I}_m,$$

$$\mathbf{F}^T \mathbf{Y} = \begin{pmatrix} \sum_{i=1}^n y_1^i \\ \sum_{i=1}^n y_2^i \\ \dots \\ \sum_{i=1}^n y_m^i \end{pmatrix} = n \begin{pmatrix} \bar{y}_1 \\ \bar{y}_2 \\ \dots \\ \bar{y}_m \end{pmatrix},$$

$$\widehat{\boldsymbol{\theta}} = (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{Y} = \frac{1}{n} \mathbf{I}_m n \begin{pmatrix} \bar{y}_1 \\ \bar{y}_2 \\ \dots \\ \bar{y}_m \end{pmatrix} = \begin{pmatrix} \bar{y}_1 \\ \bar{y}_2 \\ \dots \\ \bar{y}_m \end{pmatrix}.$$

МНК-оценки параметров:  $\widehat{\theta}_k = \bar{y}_k = \frac{1}{n} \sum_{i=1}^n y_k^i$ ,  $k = \overline{1, m}$ .

Для проверки гипотезы  $H_0: \theta_1 = \theta_2 = \theta_3 = \dots = \theta_m$  (нулевая гипотеза о незначимом влиянии различных уровней входного фактора на результат) воспользуемся статистикой Фишера

$$(3.6): F = \frac{1}{kS^2} (\widehat{\boldsymbol{\tau}} - \boldsymbol{\tau})^T \left[ \mathbf{T} (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{T}^T \right]^{-1} (\widehat{\boldsymbol{\tau}} - \boldsymbol{\tau}) \sim F(k, N - m) \text{ при}$$

следующих значениях параметров:  $\boldsymbol{\tau} = \mathbf{0}$ ,  $k = m - 1$ . Матрицу  $\mathbf{T}$  определим из условия гипотезы  $\mathbf{T}\boldsymbol{\theta} = \mathbf{0}$ :

$$\mathbf{T} = \begin{array}{cccccccc|c} \mathbf{1} & -\mathbf{1} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} & \mathbf{0} & 1 \\ \mathbf{1} & \mathbf{0} & -\mathbf{1} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} & \mathbf{0} & 2 \\ \mathbf{1} & \mathbf{0} & \mathbf{0} & -\mathbf{1} & \dots & \mathbf{0} & \mathbf{0} & \mathbf{0} & 3 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & -\mathbf{1} & \mathbf{0} & m-2 \\ \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} & -\mathbf{1} & m-1 \\ \hline 1 & 2 & 3 & & & & & m & \end{array}.$$

Явный вид статистики (3.6):

$$F_{m-1, m(n-1)} = \frac{mn(n-1) \sum_{k=1}^m (\bar{y}_k - \bar{y})^2}{(m-1) \sum_{k=1}^m \sum_{i=1}^n (y_k^i - \bar{y}_k)^2} \sim F(m-1, nm-m). \quad (3.11)$$



Доказательство проведём с использованием классической модели линейной регрессии при гауссовом шуме наблюдения, затем полученный результат сравним с *SS*-технологией Фишера.

### 1. Регрессионная модель.

$$\begin{aligned} \mathbf{T}(\mathbf{F}^T \mathbf{F})^{-1} \mathbf{T}^T &= \begin{bmatrix} 1 & -1 & 0 & 0 & \dots & 0 & 0 & 0 \\ 1 & 0 & -1 & 0 & \dots & 0 & 0 & 0 \\ 1 & 0 & 0 & -1 & \dots & 0 & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 0 & 0 & 0 & \dots & 0 & -1 & 0 \\ 1 & 0 & 0 & 0 & \dots & 0 & 0 & -1 \end{bmatrix}_{(m-1) \times m} \cdot \frac{1}{n} \cdot \mathbf{I}_m = \begin{bmatrix} 1 & 1 & 1 & 1 & \dots & 1 & 1 & 1 \\ -1 & 0 & 0 & 0 & \dots & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 & \dots & 0 & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & 0 & \dots & 0 & -1 & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 & 0 & -1 \end{bmatrix}_{m \times (m-1)} = \\ &= \frac{1}{n} \begin{bmatrix} 2 & 1 & 1 & 1 & \dots & 1 & 1 & 1 \\ 1 & 2 & 1 & 1 & \dots & 1 & 1 & 1 \\ 1 & 1 & 2 & 1 & \dots & 1 & 1 & 1 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 1 & 1 & 1 & \dots & 1 & 2 & 1 \\ 1 & 1 & 1 & 1 & \dots & 1 & 1 & 2 \end{bmatrix}_{(m-1) \times (m-1)} = \frac{1}{n} (\mathbf{I}_{m-1} + \mathbf{E}_{m-1}), \end{aligned}$$

ГДЕ

$$\mathbf{E}_{m-1} = \begin{bmatrix} 1 & 1 & 1 & 1 & \dots & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & \dots & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & \dots & 1 & 1 & 1 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 1 & 1 & 1 & \dots & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & \dots & 1 & 1 & 1 \end{bmatrix}_{(m-1) \times (m-1)},$$

$$\mathbf{I}_{m-1} = \begin{bmatrix} 1 & 0 & 0 & 0 & \dots & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & \dots & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & \dots & 0 & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & 0 & \dots & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 & 0 & 1 \end{bmatrix}_{(m-1) \times (m-1)}.$$

$$\left[ \mathbf{T}(\mathbf{F}^T \mathbf{F})^{-1} \mathbf{T}^T \right]^{-1} = \left[ \frac{1}{n} (\mathbf{I}_{m-1} + \mathbf{E}_{m-1}) \right]^{-1} = n \left( \mathbf{I}_{m-1} - \frac{1}{m} \mathbf{E}_{m-1} \right), \text{ так как}$$

$$\begin{aligned}
& (\mathbf{I}_{m-1} + \mathbf{E}_{m-1}) \left( \mathbf{I}_{m-1} - \frac{1}{m} \mathbf{E}_{m-1} \right) = \mathbf{I}_{m-1} + \mathbf{E}_{m-1} - \frac{1}{m} \mathbf{E}_{m-1} - \frac{1}{m} \mathbf{E}_{m-1}^2 = \\
& = \mathbf{I}_{m-1} + \frac{m-1}{m} \mathbf{E}_{m-1} - \frac{m-1}{m} \mathbf{E}_{m-1} = \mathbf{I}_{m-1}.
\end{aligned}$$

(Здесь использовано матричное тождество для матриц из единиц:  $\mathbf{E}_m \cdot \mathbf{E}_m = m\mathbf{E}_m$ ).

$$\hat{\boldsymbol{\tau}} = \mathbf{T}\hat{\boldsymbol{\theta}} = \begin{bmatrix} 1 & -1 & 0 & 0 & \dots & 0 & 0 & 0 & 0 \\ 1 & 0 & -1 & 0 & \dots & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & -1 & \dots & 0 & 0 & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 0 & 0 & 0 & \dots & 0 & -1 & 0 & 0 \\ 1 & 0 & 0 & 0 & \dots & 0 & 0 & -1 & 0 \end{bmatrix} \begin{bmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \\ \hat{\theta}_3 \\ \dots \\ \hat{\theta}_m \end{bmatrix} = \begin{bmatrix} \bar{y}_1 - \bar{y}_2 \\ \bar{y}_1 - \bar{y}_3 \\ \bar{y}_1 - \bar{y}_4 \\ \dots \\ \bar{y}_1 - \bar{y}_m \end{bmatrix},$$

$$\begin{aligned}
\hat{\boldsymbol{\tau}}^T \left[ \mathbf{T}(\mathbf{F}^T \mathbf{F})^{-1} \mathbf{T}^T \right]^{-1} \hat{\boldsymbol{\tau}} &= \hat{\boldsymbol{\tau}}^T n \left( \mathbf{I}_{m-1} - \frac{1}{m} \mathbf{E}_{m-1} \right) \hat{\boldsymbol{\tau}} = n \hat{\boldsymbol{\tau}}^T \hat{\boldsymbol{\tau}} - \frac{n}{m} \hat{\boldsymbol{\tau}}^T \mathbf{E}_{m-1} \hat{\boldsymbol{\tau}} = \\
&= n \sum_{k=2}^m (\bar{y}_1 - \bar{y}_k)^2 - \frac{n}{m} \sum_{k=2}^m \sum_{i=2}^m (\bar{y}_1 - \bar{y}_k)(\bar{y}_1 - \bar{y}_i) = n \sum_{k=1}^m (\bar{y}_1 - \bar{y}_k)^2 - \frac{n}{m} \sum_{k=1}^m \sum_{i=1}^m (\bar{y}_1 - \bar{y}_k), \\
(\bar{y}_1 - \bar{y}_i) &= \frac{nm\bar{y}_1^2}{m} - \frac{2nm\bar{y}_1\bar{y}}{m} + n \sum_{k=1}^m \bar{y}_k^2 - \frac{nm\bar{y}_1^2}{m} + \frac{2nm\bar{y}_1\bar{y}}{m} - nm\bar{y}^2 = n \sum_{k=1}^m \bar{y}_k^2 - nm\bar{y}^2 = \\
&= n \sum_{k=1}^m \bar{y}_k^2 - 2nm\bar{y}^2 + nm\bar{y}^2 = n \sum_{k=1}^m \bar{y}_k^2 - 2n\bar{y} \sum_{k=1}^m \bar{y}_k + n \sum_{k=1}^m \bar{y}^2 = n \sum_{k=1}^m (\bar{y}_k - \bar{y})^2,
\end{aligned}$$

$$ks^2 = (m-1) \frac{R_0^2}{N-m} = \frac{m-1}{N-m} \sum_{k=1}^m \sum_{i=1}^n (y_k^i - \hat{\theta}_k)^2 = \frac{m-1}{nm-m} \sum_{k=1}^m \sum_{i=1}^n (y_k^i - \bar{y}_k)^2,$$

$$\begin{aligned}
F &= \frac{1}{ks^2} \hat{\boldsymbol{\tau}}^T \left[ \mathbf{T}(\mathbf{F}^T \mathbf{F})^{-1} \mathbf{T}^T \right]^{-1} \hat{\boldsymbol{\tau}} = \frac{n \sum_{k=1}^m (\bar{y}_k - \bar{y})^2}{\frac{m-1}{nm-m} \sum_{k=1}^m \sum_{i=1}^n (y_k^i - \bar{y}_k)^2} = \\
&= \frac{mn(n-1) \sum_{k=1}^m (\bar{y}_k - \bar{y})^2}{(m-1) \sum_{k=1}^m \sum_{i=1}^n (y_k^i - \bar{y}_k)^2}.
\end{aligned}$$

## 2. *SS-технология.*

$$F = \frac{SS_{\text{факт}} / (m-1)}{SS_{\text{ош}} / (N-m)},$$

$$SS_{\text{ош}} = R_0^2 = \sum_{k=1}^m \sum_{i=1}^n (y_k^i - \bar{y}_k)^2, \quad SS_{\text{факт}} = \sum_{k=1}^m \sum_{i=1}^n (\bar{y}_k - \bar{y})^2 = n \sum_{k=1}^m (\bar{y}_k - \bar{y})^2,$$

$$F = \frac{SS_{\text{факт}} / (m-1)}{SS_{\text{ош}} / (nm-m)} = \frac{n \sum_{k=1}^m (\bar{y}_k - \bar{y})^2 / (m-1)}{\sum_{k=1}^m \sum_{i=1}^n (y_k^i - \bar{y}_k)^2 / (nm-m)} =$$

$$= \frac{mn(n-1) \sum_{k=1}^m (\bar{y}_k - \bar{y})^2}{(m-1) \sum_{k=1}^m \sum_{i=1}^n (y_k^i - \bar{y}_k)^2} \sim F(m-1, nm-m).$$

В основе дисперсионного анализа лежит разделение дисперсии на части или компоненты. В обозначениях, принятых в дисперсионном анализе (*SS* – от английского *Sum of Squares*):

$SS_{\text{общ}} = R_1^2$	компонента, описывающая общую дисперсию (изменчивость) выходного фактора;
$SS_{\text{ош}} = R_0^2$	компонента, описывающая влияние шума наблюдения на выходной фактор;
$SS_{\text{факт}} = R_1^2 - R_0^2$	компонента, описывающая влияние входного фактора (входных факторов) на выходной фактор.

Доказательство зависимости между указанными компонентами

$$SS_{\text{общ}} = SS_{\text{факт}} + SS_{\text{ош}}$$

аналогично (3.8) («закон сохранения суммы квадратов»).

### 3.6 Двухфакторный дисперсионный анализ при произвольном числе наблюдений

В этом случае уравнения линейной регрессии в скалярном виде имеют вид:

$$y_{ij}^k = \theta_{ij} + \xi_{ij}^k, \quad i = \overline{1, m_1}, \quad j = \overline{1, m_2}, \quad k = \overline{1, n_{ij}},$$

где  $\theta_{ij}$  – совместное влияние значения первого входного фактора с номером уровня  $i$  и второго входного фактора с номером уровня  $j$  на результат;  $m_1$  – число уровней первого входного фактора;  $m_2$  – число уровней второго входного фактора;  $n_{ij}$  – число наблюдений выходного фактора при значениях первого входного фактора с номером уровня  $i$  и второго входного фактора с номером уровня  $j$  (всего  $m = m_1 m_2$  неизвестных параметров и  $N = \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} n_{ij}$  наблюдений);  $y_{ij}^k$  – значение выходного фактора в  $k$ -м наблюдении при сочетании воздействия первого входного фактора с номером уровня  $i$  и второго входного фактора с номером уровня  $j$ ;  $\xi_{ij}^k$  – некоррелированный гауссовский шум наблюдения с нулевым математическим ожиданием и дисперсией  $\sigma^2$ . Будем считать, что число наблюдений  $n_{ij}$  при каждом сочетании уровней двух входных факторов не является одинаковым (полагаем, что  $\forall(i, j) n_{ij} \neq 0$ ).

В табл. 3.2 обозначено:  $n_i = \sum_{j=1}^{m_2} n_{ij}$ ,  $n_j = \sum_{i=1}^{m_1} n_{ij}$ ,  $n = \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} n_{ij}$ .

**Таблица 3.2 – Таблица средних для двухфакторного  
многоруовнего дисперсионного анализа при неодинаковом числе  
наблюдений сочетаний двух факторов**

Уровни первого фактора	Уровни второго фактора									Общие средние	
	1		2		3		...	$m_2$			
1	$\bar{y}_{11}$	$n_{11}$	$\bar{y}_{12}$	$n_{12}$	$\bar{y}_{13}$	$n_{13}$	...	$\bar{y}_{1m_2}$	$n_{1m_2}$	$\bar{y}_1$	$n_1$
2	$\bar{y}_{21}$	$n_{21}$	$\bar{y}_{22}$	$n_{22}$	$\bar{y}_{23}$	$n_{23}$	...	$\bar{y}_{2m_2}$	$n_{2m_2}$	$\bar{y}_2$	$n_2$
...	...	...	...	...	...	...	...	...	...	...	...
$m_1$	$\bar{y}_{m_1 1}$	$n_{m_1 1}$	$\bar{y}_{m_1 2}$	$n_{m_1 2}$	$\bar{y}_{m_1 3}$	$n_{m_1 3}$	...	$\bar{y}_{m_1 m_2}$	$n_{m_1 m_2}$	$\bar{y}_{m_1}$	$n_{m_1}$
Общие средние	$\bar{y}_1$	$n_1$	$\bar{y}_2$	$n_2$	$\bar{y}_3$	$n_3$	...	$\bar{y}_{m_2}$	$n_{m_2}$	$\bar{y}$	$n$

Параметры статистики Фишера:

$$N = n = \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} n_{ij} \text{ – общее число наблюдений;}$$

$m = m_1 m_2$  – число оцениваемых параметров

$$\theta_{11}, \theta_{12}, \dots, \theta_{1m_2}, \theta_{21}, \theta_{22}, \dots, \theta_{2m_2}, \dots, \theta_{m_1 1}, \theta_{m_1 2}, \dots, \theta_{m_1 m_2};$$

$K = m - 1 = m_1 m_2 - 1$  – число связей, накладываемых проверяемой гипотезой

$$H_0 : \theta_{11} = \theta_{12} = \dots = \theta_{1m_2} = \theta_{21} = \theta_{22} = \dots = \theta_{2m_2} = \dots = \theta_{m_1 1} = \theta_{m_1 2} = \dots = \theta_{m_1 m_2} .$$

Средние значения:

$$\bar{y}_{ij} = \frac{1}{n_{ij}} \sum_{k=1}^{n_{ij}} y_{ij}^k \text{ – среднее значение наблюдаемого (выходного,}$$

результатирующего) фактора при воздействии

первого фактора с номером уровня  $i$ ,  $i = \overline{1, m_1}$   
и второго фактора с номером уровня  $j$ ,  
 $j = \overline{1, m_2}$ ;

$\bar{y}_i = \frac{1}{n_i} \sum_{j=1}^{m_2} n_{ij} \bar{y}_{ij}$  – среднее значение наблюдаемого фактора при  
воздействии первого фактора с номером  
уровня  $i$ ,  $i = \overline{1, m_1}$ ;

$\bar{y}_j = \frac{1}{n_j} \sum_{i=1}^{m_1} n_{ij} \bar{y}_{ij}$  – среднее значение наблюдаемого фактора при  
воздействии второго фактора с номером  
уровня  $j$ ,  $j = \overline{1, m_2}$ ;

$\bar{y} = \frac{1}{n} \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} \sum_{k=1}^{n_{ij}} y_{ij}^k = \frac{1}{n} \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} n_{ij} \bar{y}_{ij} = \frac{1}{n} \sum_{j=1}^{m_2} n_j \bar{y}_j = \frac{1}{n} \sum_{i=1}^{m_1} n_i \bar{y}_i$  – гло-  
бальное среднее значение наблюдаемого  
фактора;

$SS_{\text{общ}} = \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} \sum_{k=1}^{n_{ij}} (y_{ij}^k - \bar{y})^2$  – общая изменчивость выходного  
фактора.

Заметим, что МНК-оценка неизвестных параметров в рам-  
ках используемой модели наблюдения всегда равна среднему  
значению выходного фактора при заданном сочетании входных  
факторов:

$$\hat{\theta}_{ij} = \bar{y}_{ij}, \quad i = \overline{1, m_1}, \quad j = \overline{1, m_2}.$$

Поэтому сумма квадратов ошибок наблюдения равна

$$SS_{\text{ош}} = R_0^2 = \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} \sum_{k=1}^{n_{ij}} (y_{ij}^k - \hat{\theta}_{ij})^2 = \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} \sum_{k=1}^{n_{ij}} (y_{ij}^k - \bar{y}_{ij})^2.$$

Статистика Фишера:

$$F = \frac{(SS_{\text{общ}} - SS_{\text{ош}})/(m-1)}{SS_{\text{ош}}/(N-m)} \sim F(m-1, N-m).$$

### 3.7 Общий подход к дисперсионному анализу

Повторим ещё раз  $SS$ -технологии Фишера применительно к самому общему случаю дисперсионного анализа для произвольного числа входных факторов. Изменчивость (вариабельность) наблюдаемого (результатирующего) фактора определяется: 1) изменчивостью шума наблюдения; 2) изменчивостью значений входных факторов. Проверка гипотезы о значимости влияния входных факторов на результат основывается на вычислении статистики Фишера:

$$F = \frac{SS_{\text{факт}}/K}{SS_{\text{ош}}/(N-m)} = \frac{(SS_{\text{общ}} - SS_{\text{ош}})/K}{SS_{\text{ош}}/(N-m)} \sim F(K, N-m). \quad (3.12)$$

$SS_{\text{общ}}$  описывает общую изменчивость (дисперсию, вариабельность) выходного фактора и вычисляется как сумма квадратов отклонений наблюдаемого выходного фактора от его среднего значения:

$$SS_{\text{общ}} = R_1^2 = \sum_i (y^i - \bar{y})^2, \quad \bar{y} = \frac{1}{N} \sum_i y^i, \quad \mathbf{i} = \{(i_1, i_2, i_3 \dots)\},$$

где  $\bar{y}$  – среднее значение выходного фактора (по всем наблюдениям  $y^i$ );  $N$  – общее число наблюдений (опытов);  $i$  – значение вектора индексов, определяющее все используемые сочетания значений входных факторов и номера опытов для каждого сочетания значений входных факторов.

Компонента  $SS_{ou}$  описывает влияние шума наблюдения и других случайных факторов на выходной фактор и рассчитывается как остаточная сумма квадратов значений невязки:

$$SS_{ou} = R_0^2 = \sum_{\mathbf{k}} \sum_{\mathbf{i}} (y_{\mathbf{k}}^{\mathbf{i}} - \bar{y}_{\mathbf{k}})^2,$$

где  $\bar{y}_{\mathbf{k}} = \frac{1}{n_{\mathbf{k}}} \sum_{\mathbf{i}} y_{\mathbf{k}}^{\mathbf{i}}$  – среднее значение выходного фактора в группе  $\mathbf{k}$ ;  $\mathbf{k} = \{(k_1, k_2, k_3 \dots)\}$  – вектор индексов групп входных факторов.

Компонента  $SS_{факт}$  описывает влияние значений исследуемых входных факторов на результат наблюдения (выходной фактор):

$$SS_{факт} = \sum_{\mathbf{k}} n_{\mathbf{k}} (\bar{y}_{\mathbf{k}} - \bar{y})^2.$$

Здесь  $n_{\mathbf{k}}$  – число наблюдений в группе  $\mathbf{k}$ ,  $\bar{y}_{\mathbf{k}}$  – среднее значение выходного фактора в группе  $\mathbf{k}$ . Эту факторную компоненту общей изменчивости можно также определить как

$$SS_{факт} = SS_{общ} - SS_{ou}.$$

Через  $K$  обозначено число связей (ограничений), наложенных на входные факторы проверяемой гипотезой.

Все рассмотренные выше модели дисперсионного анализа могут быть сведены к общей модели дисперсионного анализа. В табл. 3.3 показано, как это может быть выполнено. Результаты использования общей модели (3.12) совпадают с ранее полученными выражениями для расчёта статистики Фишера.



Таблица 3.3 – Частные модели дисперсионного анализа

$$F = \frac{SS_{\text{факт}}/K}{SS_{\text{ост}}/(N-m)} = \frac{(SS_{\text{общ}} - SS_{\text{ост}})/K}{SS_{\text{ост}}/(N-m)} \sim F(K, N-m)$$

Модель дисперсионного анализа	Число опытов, $N$	Число параметров, $m$	Связи входных факторов, накладываемые проверяемой гипотезой; число связей $K$	Сумма квадратов ошибок $SS_{\text{ост}}$	Общая изменчивость $SS_{\text{общ}}$
Однофакторная двухуровневая (п.п.3.2, 3.4)	$N = 2n$	$m = 2$	$a = b, K = 1$	$n \sum_{i=1}^n (y_i^j - \bar{y}_1)^2 + n \sum_{i=1}^n (y_i^j - \bar{y}_2)^2$	$\sum_{i=1}^n (y_i^j - \bar{y})^2 + \sum_{i=1}^n (y_i^j - \bar{y})^2$
Однофакторная многоуровневая (п.3.5)	$N = mn$	$m$	$\theta_1 = \theta_2 = \theta_3 = \dots = \theta_m, K = m - 1$	$\sum_{i=1}^m \sum_{k=1}^n (y_i^k - \bar{y}_i)^2$	$\sum_{i=1}^m \sum_{k=1}^n (y_i^k - \bar{y})^2$
Двухфакторная многоуровневая (п.3.6)	$N = \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} n_{ij}$	$m = m_1 m_2$	$\theta_{11} = \theta_{12} = \dots = \theta_{1m_2} = \theta_{21} = \theta_{22} = \dots = \theta_{2m_2} = \dots = \theta_{m_1 m_2} = \theta_{m_1 1} = \theta_{m_1 2} = \dots = \theta_{m_1 m_2}, K = m_1 m_2 - 1$	$\sum_{k=1}^n \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} (y_{ij}^k - \bar{y}_{ij})^2$	$\sum_{k=1}^n \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} (y_{ij}^k - \bar{y})^2$
Общая (п.3.7)	$N = \sum_k n_k$	$m$	$K = m - 1$	$\sum_k \sum_i (y_k^i - \bar{y}_k)^2$	$\sum_k n_k (\bar{y}_k - \bar{y})^2$

Для проведения дисперсионного анализа не всегда нужно записывать регрессионную модель наблюдения и строить МНК-оценки параметров. Достаточно использовать общий подход, определяемый уравнением

### 3.8 Пример двухфакторного дисперсионного анализа

*Исследование потребления топлива автомобилем.* Исследователь хочет выяснить, оказывают ли тип потребляемого бензина и тип автомобиля влияние на расход топлива. Для этого будут использованы два типа бензина – обычный и высокооктановый, и для каждой группы будут использованы два типа автомобилей – с двумя ведущими колесами и с четырьмя. Для каждой группы будут использованы по два автомобиля, всего восемь. (Американский галлон  $\approx 3,785$  литра). Исходные данные заданы табл. 3.4.

Таблица 3.4 – **Исходные данные для исследования потребления топлива автомобилем**

<i>Топливо</i> [«92», «98»]	<i>Привод</i> [«2WD», «4WD»]	<i>Расход топлива</i> [миль/галлон]
92	2WD	<b>26.7</b>
92	2WD	<b>25.2</b>
92	4WD	<b>28.6</b>
92	4WD	<b>29.3</b>
98	2WD	<b>32.3</b>
98	2WD	<b>32.8</b>
98	4WD	<b>26.1</b>
98	4WD	<b>24.2</b>

Проверяем гипотезу  $H_0$ , заключающуюся в том, что тип топлива и тип трансмиссии автомобиля не оказывают эффекта

на потребление бензина. Альтернативной является гипотеза  $H_1$ , заключающаяся в том, что тип топлива **или** тип трансмиссии автомобиля оказывают влияние на потребление бензина. Таблица средних представлена в табл. 3.5.

**Таблица 3.5 – Таблица средних для двухфакторного дисперсионного анализа при одинаковом числе наблюдений сочетаний двух факторов**

Уровни первого фактора	Уровни второго фактора		Общие средние
	2WD	4WD	
«92»	25.95	28.95	27.45
«98»	32.55	25.15	28.85
<i>Общие средние</i>	29.25	27.05	<b>28.15</b>

Общая изменчивость (вариабельность) выходного фактора:

$$SS_{\text{общ}} = R_1^2 = \sum_i (y^i - \bar{y})^2 = 70.98.$$

Остаточная сумма квадратов, характеризующая влияние шума наблюдения:

$$SS_{\text{ош}} = R_0^2 = \sum_k \sum_i (y_k^i - \bar{y}_k)^2 = \\ = (26.7 - 25.95)^2 + (25.2 - 25.95)^2 + \dots + (24.2 - 27.05)^2 = \mathbf{3.3}.$$

Изменчивость выходного фактора из-за влияния значений входных факторов:

$$SS_{\text{факт}} = \sum_k n_k (\bar{y}_k - \bar{y})^2 = \\ = 2(25.95 - 28.15)^2 + 2(28.95 - 28.15)^2 + 2(32.55 - 28.15)^2 + \\ + 2(27.05 - 28.15)^2 = \mathbf{67.68}.$$

Или так:  $SS_{факт} = SS_{общ} - SS_{ош} = 70.98 - 3.3 = 67.68$ ;

$$F = \frac{SS_{факт} / K}{SS_{ош} / (N - m)} = \frac{67.68 / 3}{3.3 / (8 - 4)} = 27.34545455 \sim$$

$$\sim F(K, N - m) = F(3, 4).$$

$K$  – число связей (ограничений), наложенных на входные факторы проверяемой гипотезой.

$$N = 8; m = 4; K = 3.$$

Проведём исследование гипотезы  $H_0$  при различных значениях уровня значимости  $\alpha$  (табл. 36).

Таблица 3.6 – Анализ результатов дисперсионного анализа

Уровень значимости $\alpha$	Квантиль распределения Фишера $F_{3,4}^{1-\alpha}$	Сравнение $F$ и $F_{3,4}^{1-\alpha}$	Анализ выполнения гипотезы $H_0$
0.05	6.591382	$F > F_{3,4}^{0.95}$	Гипотеза $H_0$ не подтверждается. Тип топлива <b>или</b> тип трансмиссии автомобиля оказывают влияние на потребление бензина. Вероятность ошибки равна 0.05
0.01	16.69437	$F > F_{3,4}^{0.99}$	Гипотеза $H_0$ не подтверждается. Тип топлива <b>или</b> тип трансмиссии автомобиля оказывают влияние на потребление бензина. Вероятность ошибки равна 0.01

Окончание табл. 3.6

0.005	24.25912	$F > F_{3,4}^{0.995}$	Гипотеза $H_0$ не подтверждается. Тип топлива <b>или</b> тип трансмиссии автомобиля оказывают влияние на потребление бензина. Вероятность ошибки равна 0.005
0.001	56.17719	$F < F_{3,4}^{0.999}$	Данные не противоречат гипотезе $H_0$ . Гипотеза $H_0$ принимается (тип топлива и тип трансмиссии автомобиля не оказывают эффекта на потребление бензина). Вероятность ошибки не установлена
<b>0.003989453</b>	<b>27.34545455</b>	<b><math>F = F_{3,4}^{0.996010547}</math></b>	Граничное значение уровня значимости. Гипотеза $H_0$ не подтверждается. Тип топлива <b>или</b> тип трансмиссии автомобиля оказывают влияние на потребление бензина. Вероятность ошибки равна 0.004

## 4 МЕТОДЫ И АЛГОРИТМЫ КЛАСТЕРИЗАЦИИ ДАННЫХ

### 4.1 Задача кластеризации

Кластер (англ. *cluster*) переводится как «*группа*», «*рой*», «*пачка*», «*скопление*», «*сгусток*», «*связка*», «*гроздь*». Употребляется этот термин в экономике, компьютерной инженерии, музыке, астрономии и означает множество объектов функционально схожих между собой и собранных в одну связку. Применительно к (интеллектуальному) анализу данных под кластерным анализом (*Data clustering*) понимается задача разбиения заданной выборки объектов на непересекающиеся подмножества, называемые кластерами, так чтобы каждый кластер состоял из схожих объектов, а объекты разных кластеров существенно отличались.

*Классификация*, являясь наиболее простой задачей *Data Mining*, относится к стратегии «обучение с учителем», для ее решения обучающая выборка должна содержать значения как входных переменных, так и выходных (целевых) переменных. *Кластеризация*, напротив, является задачей *Data Mining*, относящейся к стратегии «обучение без учителя», т.е. не требует наличия значения целевых переменных в обучающей выборке. Таким образом, синонимами термина «**кластеризация**» являются «*автоматическая классификация*», «*обучение без учителя*» и «*таксономия*» (*taxonomy*).

*Таксоно́мия* (от древнегреческого *τάξις* – *строй, порядок* и *νόμος* – *закон*) – учение о принципах и практике *классификации* и

*систематизации*. Термин «таксономия» впервые был предложен в 1813 году Огюстеном Декандром, занимавшимся классификацией растений, и изначально применялся только в биологии. Позже этот термин стал использоваться для обозначения общей теории классификации и систематизации в других областях знаний. Классическим примером таксономии на основе сходства является биномиальная номенклатура (деление на роды / виды живых существ), предложенная Карлом Линнеем в середине XVIII века. В современном представлении биологическая иерархия имеет около 30 уровней, 7 из них считаются основными: царство, тип, класс, отряд, семейство, род, вид. Аналогичные систематизации строятся во многих областях знания, чтобы упорядочить информацию о большом количестве объектов.

Кластеризация предназначена для разбиения совокупности объектов на однородные группы (кластеры или классы). Обычно данные представляют собой выборки точек в признаковом пространстве. Наибольшее применение кластеризация первоначально получила в таких науках как биология, антропология, психология. Кластерный анализ полезен, когда нужно классифицировать большое количество информации. Так, в медицине используется кластеризация заболеваний, лечения заболеваний или их симптомов, а также таксономия пациентов, препаратов и т.д. В археологии устанавливаются таксономии каменных сооружений и древних объектов и т.д. В маркетинге это может быть задача сегментации конкурентов и потребителей. В менеджменте примером задачи кластеризации будет разбиение персонала на различные группы, классификация потребителей и поставщиков, анализ поведения клиента-потребителя, выявление схожих производственных ситуаций, при которых возникает брак. В медицине – классификация симптомов. В социологии задача кластеризации – разбиение респондентов на однородные группы.

Алгоритмы кластеризации являются в большой степени *эвристическими*. Эвристический алгоритм – это алгоритм решения задачи, правильность которого для всех возможных случаев не доказана, но про который известно, что он даёт достаточно хорошее решение в большинстве случаев. В действительности может быть даже известно, что эвристический алгоритм формально неверен. Его всё равно можно применять, если при этом он даёт неверный результат только в отдельных, достаточно редких и хорошо выделяемых случаях или же даёт неточный, но всё же приемлемый результат. Проще говоря, эвристика – это не полностью математически обоснованный (или даже «не совсем корректный»), но при этом практически полезный алгоритм. Важно понимать, что эвристика, в отличие от корректного алгоритма решения задачи, обладает следующими особенностями:

- Она не гарантирует нахождение лучшего решения.
- Она не гарантирует нахождение решения, даже если оно заведомо существует (возможен «пропуск цели»).
- Она может дать неверное решение в некоторых случаях.

### **Формальная постановка задачи кластеризации**

Пусть  $\mathbf{X} = \{x^1, x^2, x^3, \dots, x^N\}$  – множество  $N$  объектов, заданных в  $n$ -мерном векторном признаковом пространстве:

$$x^k = (x_1^k \ x_2^k \ x_3^k \ \dots \ x_n^k)^T, \quad k = \overline{1, N},$$

$\mathbf{Y} = \{1, 2, 3, \dots\}$  – множество номеров (имён, меток) кластеров.

Задана функция расстояния между объектами  $\rho(x, x')$ .

Наиболее часто используется евклидова метрика:

$$\rho_2(x, x') = \sqrt{\sum_{m=1}^n (x_m - x'_m)^2}.$$



Используются также

- манхэттенское расстояние (расстояние городских кварталов)

$$\rho_1(x, x') = \sum_{j=1}^n |x_j - x'_j|;$$

- максимальное различие по координатам

$$\rho_\infty(x, x') = \max_{j \in [1:n]} |x_j - x'_j|.$$

Все эти расстояния являются частным случаем расстояния Минковского:

$$\rho_p(x, y) = \left( \sum_{j=1}^n |x_j - y_j|^p \right)^{1/p}.$$

Имеется конечная выборка объектов  $\mathbf{X} = \{x^1, x^2, \dots, x^N\}$ . Требуется разбить эту выборку на  $K$  непересекающихся подмножеств  $S_k, k = \overline{1, K}$ ;  $\mathbf{X} = \bigcup_{k=1}^K S_k$ , называемых кластерами, так чтобы каждый кластер состоял из объектов, близких по метрике  $\rho$ , а объекты разных кластеров существенно отличались. При этом каждому объекту  $x^i \in \mathbf{X}$  приписывается номер кластера  $y_i \in \mathbf{Y}$ . Алгоритм кластеризации – это функция  $\mathbf{X} \rightarrow \mathbf{Y}$ , которая любому объекту  $x \in \mathbf{X}$  ставит в соответствие номер кластера  $y \in \mathbf{Y}$ . Множество  $\mathbf{Y}$  в некоторых случаях известно заранее, однако чаще ставится задача определить оптимальное число кластеров с точки зрения того или иного критерия качества кластеризации.

Кластеризация (обучение без учителя) отличается от классификации (обучения с учителем) тем, что метки исходных объектов  $y_i$  изначально не заданы и даже может быть неизвестно само множество  $\mathbf{Y}$ , т.е. число объектов в кластере.

**Центром кластера**  $S_k$  (центроидом) называется геометрический центр точек кластера  $k$  в евклидовом пространстве:

$$X_k = \frac{1}{|S_k|} \sum_{x^i \in S_k} x^i, \text{ где } |S_k| - \text{число точек в кластере } k,$$

$k = 1, 2, 3, \dots, K$ ;  $K$  – число кластеров.

**Дисперсия кластера** – это мера рассеяния точек в пространстве относительно центра кластера:

$$D_k = \frac{1}{|S_k|} \sum_{x^i \in S_k} \rho^2(x^i, X_k).$$

**Радиус кластера** – это тоже мера рассеяния точек в пространстве относительно центра кластера – максимальное расстояние до центра кластера:

$$R_k = \max_{x^i \in S_k} \rho(x^i, X_k).$$

Цели кластеризации могут быть различными в зависимости от особенностей конкретной прикладной задачи:

- Понять структуру множества объектов  $X$ , разбив его на группы схожих объектов. Упростить дальнейшую обработку данных и принятия решений, работая с каждым кластером по отдельности (стратегия «разделяй и властвуй»).
- Сократить объём хранимых данных в случае сверхбольшой выборки  $X$ , оставив по одному наиболее типичному представителю от каждого кластера.
- Выделить нетипичные объекты, которые не подходят ни к одному из кластеров. Эту задачу называют одноклассовой классификацией, обнаружением нетипичности или новизны (novelty detection).

В первом случае число кластеров стараются сделать поменьше. Во втором случае важнее обеспечить высокую степень сходства объектов внутри каждого кластера, а кластеров может быть сколько угодно. В третьем случае наибольший интерес представляют отдельные объекты, не вписывающиеся ни в один из кластеров.

Во всех этих случаях может применяться иерархическая кластеризация, когда крупные кластеры дробятся на более мелкие, те, в свою очередь, дробятся ещё мельче и т. д. Такие задачи называются задачами таксономии (taxonomy). Результатом таксономии является не простое разбиение множества объектов на кластеры, а древообразная иерархическая структура. Вместо номера кластера объект характеризуется перечислением всех кластеров, которым он принадлежит: от крупного к мелкому. Таксономии строятся во многих областях знания, чтобы упорядочить информацию о большом количестве объектов. Мы будем рассматривать алгоритмы иерархической кластеризации, позволяющие автоматизировать процесс построения таксономий.

Решение задачи кластеризации принципиально неоднозначно по следующим причинам:

1. Не существует однозначно наилучшего критерия качества кластеризации. Известен целый ряд эвристических критериев, а также ряд алгоритмов, не имеющих чётко выраженного критерия, но осуществляющих достаточно разумную кластеризацию «по построению». Все они могут давать разные результаты. Следовательно, для определения качества кластеризации требуется эксперт предметной области, который бы мог оценить осмысленность выделения кластеров.

2. Число кластеров, как правило, неизвестно заранее и устанавливается в соответствии с некоторым субъективным критерием. Это справедливо только для методов дискриминации, так

как в методах кластеризации выделение кластеров идёт за счёт формализованного подхода на основе мер близости.

3. Результат кластеризации существенно зависит от метрики, выбор которой, как правило, также субъективен и определяется экспертом. Но стоит отметить, что есть ряд рекомендаций к выбору мер близости для различных задач.

## 4.2 Метод *k*-средних (*k*-means)

К настоящему времени разработано более сотни различных алгоритмов кластеризации, в результате применения которых получаются неодинаковые результаты, что объясняется особенностью работы того или иного алгоритма, ориентированного на решение *конкретной* задачи.

Наиболее популярным является метод *k-средних*. Название «*k-means*» было впервые использовано James MacQueen в 1967 году. Этот алгоритм относится к классу эвристических EM-алгоритмов (англ. *Expectation-maximization*), используемых в математической статистике для нахождения оценок максимального правдоподобия параметров вероятностных моделей. Каждая итерация алгоритма состоит из двух шагов. На E-шаге (*expectation*) вычисляется ожидаемое значение функции правдоподобия. На M-шаге (*maximization*) вычисляется оценка максимального правдоподобия, увеличивая ожидаемое правдоподобие, вычисляемое на E-шаге. Затем это значение используется для E-шага на следующей итерации. Алгоритм выполняется до сходимости.

Основная идея алгоритма *k-средних* заключается в минимизации суммарного квадратичного отклонения точек кластеров от центров этих кластеров, то есть

$$J = \sum_{k=1}^K \sum_{x_j \in S_k} \rho(x^j, X_k) \rightarrow \min_{S_1, S_2, \dots, S_K},$$

где  $K$  – число кластеров (считается известным);  $S_k$  – полученные кластеры;  $X_k$  – центр масс (центроид) векторов из кластера  $S_k$ .

Алгоритм заключается в том, что на каждой итерации заново вычисляется центр масс  $X'_k$  для каждого кластера, полученного на предыдущем шаге, затем векторы разбиваются на кластеры вновь в соответствии с тем, какой из новых центров оказался ближе по выбранной метрике:

$$x^i \in S'_k, \text{ если } \rho(x^i, X'_k) = \min_{k=1, K}.$$

На рис. 4.1 показан пример работы алгоритма k-means.

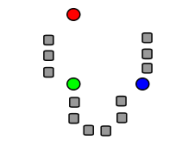
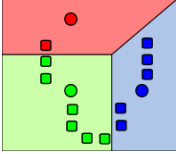
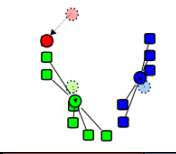
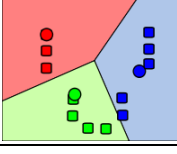
Шаг 1		Начальное расположение 10 объектов и случайно выбранные центры кластеров (красная, зелёная и синяя точки)
Шаг 2		Точки, отнесённые к кластерам по критерию близости к центру кластера по евклидовому расстоянию
Шаг 3		Вычисление новых центров кластеров
Шаг 4		Переход к шагу 2 с новым расположением центров кластеров или остановка, если нет изменений в кластерах

Рис. 4.1. Пример алгоритма k-means для трёх кластеров  
(<https://ru.wikipedia.org/wiki/K-means>)

Алгоритм завершается, когда на какой-то итерации не происходит изменения центра масс кластеров. Число итераций конечно, так как конечно множество возможных разбиений конечного множества векторов, а заикливания не происходит, так как  $J$  не увеличивается. Доказано (*David Arthur, Sergei Vassilvitskii*; 2006), что на некоторых классах множеств вычислительная сложность алгоритма (число итераций) равна  $2^{\Omega(\sqrt{n})}$ , где  $\Omega(\cdot)$  – полиномиальная сложность.

### **Алгоритм кластеризация k-средних**

#### ***Начало***

*Сформировать начальное приближение центров всех  $K$  кластеров, взяв наиболее удалённые друг от друга объекты выборки  $X_k^{(0)}$ ,  $k \in Y = \{1, 2, 3, \dots, K\}$ .*

#### ***Repeat***

*Отнести каждый объект к ближайшему центру, разбивая пространство признаков и формируя новые кластеры  $S_k^{(j)}$ ,  $k \in \{1, 2, 3, \dots, K\}$ :*

$$\forall x_i, i = \overline{1, N} : k = \arg \min_{k \in \overline{1, K}} \rho(x^i, X_k^{(j-1)}), \quad S_k^{(j)} = S_k^{(j-1)} \cup x^i.$$

*Вычислить новое положение центров:*

$$X_k^{(j)} = \frac{1}{|S_k^{(j)}|} \sum_{x_i \in S_k^{(j)}} x_i.$$

*Перейти к следующей итерации:*

$$j := j + 1$$

#### ***Until***

*Пока состав кластеров  $S_k^{(j)}$  не перестанет изменяться.*

**Недостатки алгоритма *k-means*:**

- Не гарантируется достижение глобального минимума суммарного квадратичного отклонения  $J$ , а только одного из локальных минимумов.
- Результат зависит от начального выбора центров кластеров  $\{X_k^{(0)}\}$ , их оптимальный выбор неизвестен.
- Число кластеров  $K$  надо знать заранее.

Пример 1. Дан набор из 8 точек на плоскости:

	A	B	C	D	E	F	G	H
	$x^1$	$x^2$	$x^3$	$x^4$	$x^5$	$x^6$	$x^7$	$x^8$
$x_1$	1	4	5	6	1	5	1	2
$x_2$	3	3	3	3	2	2	1	1

Выполнить кластеризацию точек по алгоритму  $k$ -means: число кластеров равно 2, метрика – манхеттенская.

$K=2$ . Произвольные начальные центры кластеров:

$$X_1^{(0)} = \begin{pmatrix} 5 \\ 2 \end{pmatrix}; \quad X_2^{(0)} = \begin{pmatrix} 4 \\ 3 \end{pmatrix}.$$

Начальные расстояния до центров кластеров и отнесение в кластер:

	A	B	C	D	E	F	G	H
	$x^1$	$x^2$	$x^3$	$x^4$	$x^5$	$x^6$	$x^7$	$x^8$
$x_1$	1	4	5	6	1	5	1	2
$x_2$	3	3	3	3	2	2	1	1
$\rho_1 \left( x^i, \begin{pmatrix} 5 \\ 2 \end{pmatrix} \right)$	5	2	1	2	4	0	5	4
$\rho_1 \left( x^i, \begin{pmatrix} 4 \\ 3 \end{pmatrix} \right)$	3	0	1	2	4	2	5	4
Кластер	2	2	1	1	1	1	1	1

При одинаковых расстояниях до центров кластеров принадлежность кластеру вектора признаков назначается произвольно (для определённости выбираем первый кластер).

Центры кластеров на следующем шаге:

$$X_1^{(1)} = \begin{pmatrix} 10/3 \\ 2 \end{pmatrix}; \quad X_2^{(1)} = \begin{pmatrix} 5/2 \\ 3 \end{pmatrix}.$$

Новые расстояния до центров и отнесение в кластер:

	A	B	C	D	E	F	G	H
	$x^1$	$x^2$	$x^3$	$x^4$	$x^5$	$x^6$	$x^7$	$x^8$
$x_1$	<b>1</b>	<b>4</b>	<b>5</b>	<b>6</b>	<b>1</b>	<b>5</b>	<b>1</b>	<b>2</b>
$x_2$	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>2</b>	<b>2</b>	<b>1</b>	<b>1</b>
$\rho_1 \left( x^i, \begin{pmatrix} 10/3 \\ 2 \end{pmatrix} \right)$	10/3	5/3	8/3	11/3	7/3	5/3	10/3	7/3
$\rho_1 \left( x^i, \begin{pmatrix} 5/2 \\ 3 \end{pmatrix} \right)$	3/2	3/2	5/2	7/2	5/2	7/2	7/2	5/2
<i>Кластер</i>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>

Центры кластеров на следующем шаге:

$$X_1^{(2)} = \begin{pmatrix} 9/4 \\ 3/2 \end{pmatrix}; \quad X_2^{(2)} = \begin{pmatrix} 4 \\ 3 \end{pmatrix}.$$

Новые расстояния до центров и отнесение в кластер:

	A	B	C	D	E	F	G	H
	$x^1$	$x^2$	$x^3$	$x^4$	$x^5$	$x^6$	$x^7$	$x^8$
$x_1$	<b>1</b>	<b>4</b>	<b>5</b>	<b>6</b>	<b>1</b>	<b>5</b>	<b>1</b>	<b>2</b>
$x_2$	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>2</b>	<b>2</b>	<b>1</b>	<b>1</b>
$\rho_1 \left( x^i, \begin{pmatrix} 9/4 \\ 3/2 \end{pmatrix} \right)$	11/4	13/4	17/4	21/4	7/4	13/4	7/4	3/4
$\rho_1 \left( x^i, \begin{pmatrix} 4 \\ 3 \end{pmatrix} \right)$	3	0	1	2	4	2	5	4
<i>Кластер</i>	<b>1</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>1</b>	<b>2</b>	<b>1</b>	<b>1</b>



Центры кластеров на следующем шаге:

$$X_1^{(3)} = \begin{pmatrix} 5/4 \\ 7/4 \end{pmatrix}; \quad X_2^{(3)} = \begin{pmatrix} 5 \\ 11/4 \end{pmatrix}.$$

Новые расстояния до центров и отнесение в кластер:

	A	B	C	D	E	F	G	H
	$x^1$	$x^2$	$x^3$	$x^4$	$x^5$	$x^6$	$x^7$	$x^8$
$x_1$	<b>1</b>	<b>4</b>	<b>5</b>	<b>6</b>	<b>1</b>	<b>5</b>	<b>1</b>	<b>2</b>
$x_2$	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>2</b>	<b>2</b>	<b>1</b>	<b>1</b>
$\rho_1 \left( x^i, \begin{pmatrix} 5/4 \\ 7/4 \end{pmatrix} \right)$	3/2	4	5	6	1/2	4	1	3/2
$\rho_1 \left( x^i, \begin{pmatrix} 5 \\ 11/4 \end{pmatrix} \right)$	17/4	5/4	1/4	5/4	19/4	3/4	23/4	19/4
<i>Кластер</i>	<b>1</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>1</b>	<b>2</b>	<b>1</b>	<b>1</b>

Видим, что состав кластеров на последней итерации не изменился. Поэтому итерации заканчиваются.

*Ответ:*  $S_1 = \{A, E, G, H\}$ ;  $S_2 = \{B, C, D, F\}$ .

### *Алгоритм кластеризации k-means++*

Это расширение направлено на оптимальный, а не случайный **выбор начальных значений** центров кластеров (David Arthur, Sergei Vassilvitskii; 2007). Описание алгоритма:

1 Выбрать первый центроид случайным образом (среди всех точек).

2 Для каждой точки найти значение квадрата расстояния до ближайшего центроида (из тех, которые уже выбраны)  $dx^2$ .

3 Выбрать из этих точек следующий центроид так, чтобы вероятность выбора точки была пропорциональна вычисленному для неё квадрату расстояния. Это можно сделать следу-

ющим образом. На шаге 2 нужно параллельно с расчётом  $dx^2$  подсчитывать сумму  $S = \sum dx^2$ . После накопления суммы взять случайное значение  $R$  из равномерно распределённого на отрезке  $[0; S]$ , по которому остаётся только определить, какой точке это соответствует. Для этого нужно снова начать подсчитывать сумму  $S = \sum dx^2$  до тех пор, пока сумма не превысит  $R$ . Текущая точка берётся в качестве центроида.

4 При выборе каждого следующего центроида специально следить за тем, чтобы он не совпал с одной из уже выбранных в качестве центроидов точек, не нужно, так как вероятность повторного выбора некоторой точки равна 0.

5 Повторять шаги 2 и 3 до тех пор, пока не будут найдены все необходимые центроиды.

Далее выполняется основной алгоритм *k-means*.

### 4.3 Алгоритм кластеризации Isodata (ИСОМАД)

**ИСОМАД** (Итеративный самоорганизующийся метод анализа данных, **Isodata** – Iterative Self-Organizing Data Analysis Techniques) аналогичен методу **k-means**, однако обладает более широким набором параметров и вспомогательных эвристических процедур. Алгоритм эвристический, результат работы во многом зависит от исходных данных и заданных начальных значениях параметров.

Исходные данные: набор векторов  $\{x_1, x_2, x_3, \dots, x_N\}$ . Параметры алгоритма:  $K$  – необходимое число кластеров;  $Q_N$  – параметр, с которым сравнивается количество выборочных об-

разов, вошедших в кластер;  $Q_S$  – параметр, характеризующий среднеквадратичное отклонение;  $Q_C$  – параметр, характеризующий компактность;  $L$  – максимальное количество пар центров кластеров, которые можно объединить;  $I$  – допустимое число циклов итерации.

Подробное описание алгоритма *Isodata* приведено в книге [1].

#### 4.4 Иерархическая кластеризация

Алгоритмы **иерархической кластеризации** называются также графовыми алгоритмами кластеризации. **Дендрограмма** – дерево (граф без циклов), построенное по матрице мер близости. Дендрограмма позволяет изобразить взаимные связи между объектами из заданного множества. Для создания дендрограммы требуется матрица сходства, которая определяет уровень сходства между любой парой объектов.

Для **агломеративных методов** (*AGglomerative NESTing, AGNES*) – методы слияния – проводится рекурсивное поочерёдное слияние близких кластеров. При слиянии отдельные точки-объекты также считаются кластерами. На каждом шаге выбирается для слияния пара наиболее схожих кластеров. На последующих шагах объединение продолжается до тех пор, пока все объекты не будут составлять один кластер.

**Дивизивные методы** (*DIVisive ANALysis, DIANA*) – методы деления – являются логической противоположностью агломеративным методам. В начале работы алгоритма все объекты принадлежат одному кластеру, который на последующих шагах делится на меньшие кластеры и в результате образуется последовательность расщепляющих групп.

Различные алгоритмы отличаются выбором метрики измерения расстояний (см. п.4.1) и критерия схожести кластеров (приведены ниже).

### **Методы измерения расстояний между кластерами**

**Метод ближнего соседа.** Здесь расстояние между двумя кластерами определяется расстоянием между двумя наиболее близкими объектами (ближайшими соседями) в различных кластерах.

**Метод наиболее удаленных соседей.** Здесь расстояния между кластерами определяются наибольшим расстоянием между любыми двумя объектами в различных кластерах (т.е. «наиболее удаленными соседями»).

**Метод Варда (Ward).** В качестве расстояния между кластерами берется прирост суммы квадратов расстояний объектов до центров кластеров, получаемый в результате их объединения. В отличие от других методов кластерного анализа для оценки расстояний между кластерами здесь используются методы дисперсионного анализа. На каждом шаге алгоритма объединяются такие два кластера, которые приводят к минимальному увеличению целевой функции, т.е. внутригрупповой суммы квадратов.

**Метод невзвешенного попарного среднего.** В качестве расстояния между двумя кластерами берется среднее расстояние между всеми парами объектов в них.

**Метод взвешенного попарного среднего.** Этот метод похож на метод невзвешенного попарного среднего, разница состоит лишь в том, что здесь в качестве весового коэффициента используется размер кластера (число объектов, содержащихся в кластере).

**Невзвешенный центроидный метод.** В качестве расстояния между двумя кластерами в этом методе берется расстояние между их геометрическими центрами.

**Взвешенный центроидный метод.** Для учета разницы между размерами кластеров (числе объектов в них) используются веса.

**Пример.** Необходимо провести кластеризацию пяти предприятий, каждое из которых характеризуется тремя переменными:  $x_1$  – среднегодовая стоимость основных производственных фондов, млрд. руб.;  $x_2$  – материальные затраты на 1 руб. произведенной продукции;  $x_3$  – объем произведенной продукции, млрд. руб. Значения переменных приведены в табл. 4.1.

Таблица 4.1 – **Финансовые показатели  
производственных предприятий**

Номер предприятия	$x_1$	$x_2$	$x_3$
1	120,0	94,0	164,0
2	85,0	75,2	92,0
3	145,0	81,0	120,0
4	78,0	76,8	86,0
5	70,0	75,9	104,0
Среднее значение, $\bar{x}_k$	99,6	80,6	113,2
Среднее квадратическое отклонение, $S_k$	28,4	10,9	27,9

Нормировка исходных данных:

$$z_k^i = \frac{x_k^i - \bar{x}_k}{s_k},$$

$$Z = \begin{pmatrix} 0.718 & 1.229 & 1.821 \\ -0.514 & -2.238 & -0.760 \\ 1.514 & 0.037 & 0.244 \\ -0.760 & -0.349 & -0.975 \\ -1,042 & -0.431 & -0.330 \end{pmatrix}.$$

Классификацию проведем при помощи иерархического **агломеративного** метода. Для построения матрицы расстояний воспользуемся квадратом евклидового расстояния. Тогда, например, квадрат расстояния между первым и вторым объектами:

$$d^2(1,2) = (0.718 - (-0.514))^2 + (1.229 - (-2.238))^2 + (1.821 - (-0.760))^2 = 20.20,$$

квадрат расстояния между первым и третьим объектами:

$$d^2(1,3) = (0.718 - 1.514)^2 + (1.229 - 0.037)^2 + (1.821 - 0.244)^2 = 4.54 \text{ и т.д.}$$

В результате получаем первоначальную матрицу расстояний, характеризующую расстояния между отдельными объектами, каждый из которых изначально является отдельным кластером:

	1	2	3	4	5
$D_0^2 =$	0	20.20	4.54	12.49	10.48
	2	0	10.30	3.68	3.73
	3		0	6.81	7.08
	4			0	<b>0.50</b>
	5				0

Как видно по элементам матрицы  $D_0^2$ , наиболее близкими являются объекты 4 и 5:  $d^2(4, 5) = 0.50$ . Объединим их в один кластер и присвоим ему метку 45. Пересчитаем расстояния всех оставшихся объектов (кластеров) до кластера 45. В матрице  $D_1^2$  расстояния между кластерами определяются по алгоритму «дальнего соседа».

Расстояние между кластером 1 и кластером 45:

$$d^2(1, 45) = \max(d^2(1, 4), d^2(1, 5)) = \max(12.49, 10.48) = 12.49.$$

Расстояние между кластером 2 и кластером 45:

$$d^2(2, 45) = \max(d^2(2, 4), d^2(2, 5)) = \max(3.68, 3.73) = 3.73.$$

Расстояние между кластером 3 и кластером 45:

$$d^2(3, 45) = \max(d^2(3, 4), d^2(3, 5)) = \max(6.81, 7.08) = 7.08.$$

Получим новую матрицу расстояний:

	1	2	3	45
$D_1^2 =$	1	0	20.20	4.54
	2	0	10.30	<b>3.73</b>
	3	10.30	0	7.08
	45			0

В матрице  $D_1^2$  опять находим самые близкие кластеры. Это будут кластеры 2 и 45,  $d^2(2, 45) = 3.73$ . Следовательно, на этом шаге объединяем кластеры 2 и 45; получим новый кластер, состоящий из кластеров 2 и 45 (объекты 2, 4, 5). Присвоим ему метку 245. Теперь имеем три кластера: 1, 245, 3. Пересчитываем расстояния  $d^2(1, 245)$ ,  $d^2(3, 245)$  и получаем матрицу  $D_2$ :

$$\begin{aligned}
 d^2(1, 245) &= \max(d^2(1, 2), d^2(1, 45)) = \\
 &= \max(20.2, 12.49) = 20.2, \\
 d^2(3, 245) &= \max(d^2(3, 2), d^2(3, 45)) = \\
 &= \max(10.3, 7.08) = 10.3.
 \end{aligned}$$

$D_2^2 =$	1	0	245	3
	1		20.20	<b>4.54</b>
	245		0	10.30
	3			0

На следующем шаге объединяем кластеры 1 и 3 ( $d^2(1, 3) = 4.54$ ) в один кластер и присваиваем ему номер 13. Теперь имеем только два кластера: 13, 245.

$$\begin{aligned}
 d^2(13, 245) &= \max(d^2(1, 245), d^2(3, 245)) = \\
 &= \max(20.20, 10.3) = 20.20
 \end{aligned}$$

И, наконец, на последнем шаге объединяем кластеры 13 и 245 на квадрате расстояния 20.20.

		13	245
$D_3^2 =$	13	0	<b>20.20</b>
	245		0

Представим результаты классификации в виде дендрограммы (рис. 4.2). Дендрограмма свидетельствует о том, что кластер  $S_2$  более однороден по составу входящих объектов, так как в нем объединение происходит при меньших расстояниях, чем в кластере  $S_1$ .



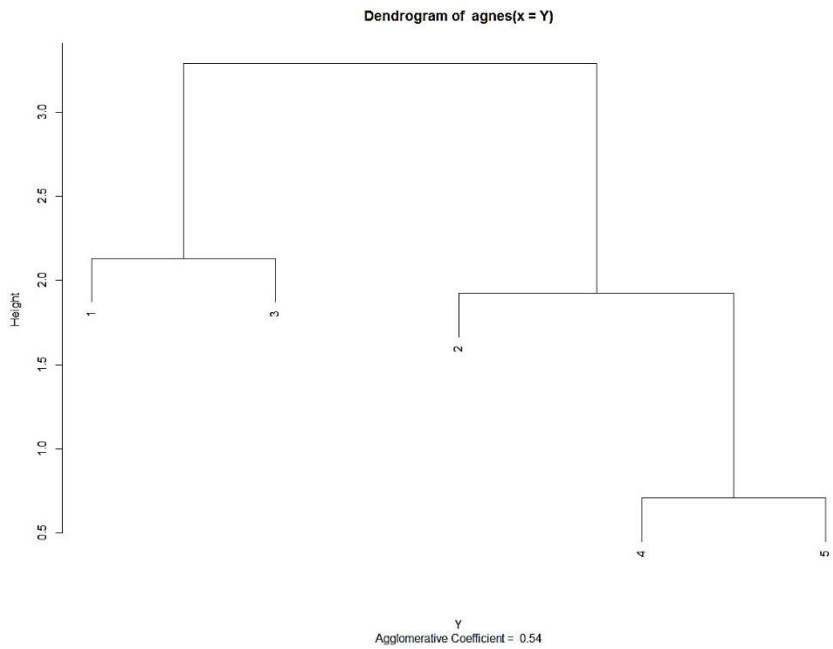


Рис. 4.2. Дендрограмма кластеризации пяти  
производственных предприятий

## 5 РАСПОЗНАВАНИЕ ОБРАЗОВ

Задачей распознавания образов является отнесение наблюдаемого объекта (вектора признаков)  $\mathbf{x}$  к одному из классов  $\{\omega_1, \omega_2, \omega_3, \omega_4, \dots\}$ . То есть требуется разбить всё множество  $X$  возможных значений признаков объектов на непересекающиеся области  $G_1, G_2, G_3, G_4, \dots$  из условия минимизации вероятности ошибки. Будем считать, что условная вероятность принадлежности объекта каждому из классов (априорная вероятность) известна или её можно определить/измерить на основе известной модели наблюдения.

### 5.1 Статистическая проверка гипотез

Простейшей задачей распознавания образов является задача статистической проверки гипотез из математической статистики. Рассмотрим **байесовское решающее правило**, минимизирующее вероятность ошибки решения, для случая двух альтернатив. Рассмотрим модель наблюдения, в которой признаки объекта являются непрерывными случайными величинами, которые описываются безусловной и условной плотностями вероятностей. Пусть  $x$  – вектор наблюдений признаков (объект); задача заключается в том, чтобы определить его принадлежность классу  $\omega_1$  (гипотеза  $H_1$ ) или  $\omega_2$  (гипотеза  $H_2$ ). Решающее правило основано на максимизации **апостериорной**

**вероятности**  $\Pr\{\omega_i|x\}$ , вычисляемой по формуле Байеса (Томас Байес (англ. *Reverend Thomas Bayes*) – английский математик XVIII века), которая для случая непрерывных признаков имеет вид:

$$\Pr\{\omega_i|x\} = \frac{f(x|\omega_i)\mathbf{p}_i}{f(x)}, i = \overline{1,2}, \quad (5.1)$$

$$\begin{aligned} \omega_i &= \arg \max_{\omega_i} \Pr\{\omega_i|x\} = \arg \max_{\omega_i} \frac{f(x|\omega_i)\mathbf{p}_i}{f(x)} = \\ &= \arg \max_{\omega_1, \omega_2} \left\{ \frac{f(x|\omega_1)\mathbf{p}_1}{f(x)}, \frac{f(x|\omega_2)\mathbf{p}_2}{f(x)} \right\}, \end{aligned}$$

где  $\mathbf{p}_1 = \Pr\{\omega_1\}$ ,  $\mathbf{p}_2 = \Pr\{\omega_2\}$  – **априорные** вероятности классов объектов (гипотез);

$f(x) = f(x|\omega_1)\mathbf{p}_1 + f(x|\omega_2)\mathbf{p}_2$  – безусловная плотность вероятности,  $\mathbf{p}_1 + \mathbf{p}_2 = 1$ ;

$f(x|\omega_1)$ ,  $f(x|\omega_2)$  – условные плотности вероятности.

Так как решающее правило (5.1) не зависит от значения знаменателя  $f(x)$ , его можно записать в виде

$$f(x|\omega_1)\mathbf{p}_1 > f(x|\omega_2)\mathbf{p}_2 \rightarrow x \in \omega_1,$$

$$f(x|\omega_1)\mathbf{p}_1 < f(x|\omega_2)\mathbf{p}_2 \rightarrow x \in \omega_2.$$

С использованием отношения правдоподобия  $l(x)$  решающее правило принимает вид:

$$\begin{aligned}
 l(x) &= \frac{f(x|\omega_1)}{f(x|\omega_2)} > h \rightarrow x \in \omega_1, \\
 l(x) &= \frac{f(x|\omega_1)}{f(x|\omega_2)} < h \rightarrow x \in \omega_2.
 \end{aligned}
 \tag{5.2}$$

где  $h = \mathbf{p}_2/\mathbf{p}_1$  – пороговое значение.

Уравнение (5.2) называется **байесовским критерием**, минимизирующим ошибку решения. Вероятность ошибки классификации:

$$\begin{aligned}
 \varepsilon &= \Pr\{ош\} = \mathbf{p}_1 \Pr\{ош|\omega_1\} + \mathbf{p}_2 \Pr\{ош|\omega_2\} = \\
 &= \mathbf{p}_1 \Pr\{x \in G_2|\omega_1\} + \mathbf{p}_2 \Pr\{x \in G_1|\omega_2\} = \\
 &= \mathbf{p}_1 \int_{G_2} f(x|\omega_1) dx + \mathbf{p}_2 \int_{G_1} f(x|\omega_2) dx = \mathbf{p}_1 \varepsilon_1 + \mathbf{p}_2 \varepsilon_2.
 \end{aligned}$$

Величины

$$\varepsilon_1 = \int_{G_2} f(x|\omega_1) dx \quad \text{и} \quad \varepsilon_2 = \int_{G_1} f(x|\omega_2) dx$$

называются вероятностями ошибок *первого рода* и *второго рода*.

В радиолокации эти ошибки называются вероятностью *пропуска цели* и вероятностью *ложной тревоги*.

В медицинской диагностике используются названия *чувствительность* ( $1 - \varepsilon_1$ ) и *специфичность* ( $1 - \varepsilon_2$ ).

Докажем, что байесовский классификатор обеспечивает минимальную вероятность ошибки. Воспользуемся тем фактом, что события  $x \in G_1$  и  $x \in G_2$  образуют полную группу:

$$\begin{aligned}
& \int_{G_2} f(x|\omega_1) dx + \int_{G_1} f(x|\omega_1) dx = 1. \\
\varepsilon &= \mathbf{p}_1 \int_{G_2} f(x|\omega_1) dx + \mathbf{p}_2 \int_{G_1} f(x|\omega_2) dx = \\
&= \mathbf{p}_1 \left( 1 - \int_{G_1} f(x|\omega_1) dx \right) + \mathbf{p}_2 \int_{G_1} f(x|\omega_2) dx = \\
&= \mathbf{p}_1 + \int_{G_1} [\mathbf{p}_2 f(x|\omega_2) - \mathbf{p}_1 f(x|\omega_1)] dx \rightarrow \min_{G_1}.
\end{aligned}$$

Из последнего соотношения видно, что минимум вероятности ошибки достигается, когда в область  $G_1$  включаются точки  $x$ , для которых значение подынтегральной функции отрицательно, то есть  $\mathbf{p}_2 f(x|\omega_2) - \mathbf{p}_1 f(x|\omega_1) < 0$  или  $\frac{f(x|\omega_1)}{f(x|\omega_2)} > \frac{\mathbf{p}_2}{\mathbf{p}_1} \rightarrow x \in G_1$ , что совпадает с правилом Байеса (5.2).

**Пример 5.1.** Построить байесовское решающее правило для двух классов независимых гауссовских наблюдений случайной величины  $X$  с плотностями вероятностей

$$\begin{aligned}
f(x|\omega_1) &= \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}(x-1)^2\right], \\
f(x|\omega_2) &= \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}(x-3)^2\right]
\end{aligned}$$

и априорными вероятностями классов  $\mathbf{p}_1 = 1/3$ ,  $\mathbf{p}_2 = 2/3$ .

Отношение правдоподобия:

$$\begin{aligned}
 l(x) &= \frac{f(x|\omega_1)}{f(x|\omega_2)} = \frac{\frac{1}{\sqrt{2\pi}} \exp\left[-\frac{(x-1)^2}{2}\right]}{\frac{1}{\sqrt{2\pi}} \exp\left[-\frac{(x-3)^2}{2}\right]} = \\
 &= \exp\left[\frac{(x-3)^2 - (x-1)^2}{2}\right] = \exp(4-2x).
 \end{aligned}$$

Получаем следующее решающее правило:

$$\exp(4-2x) > \frac{\mathbf{p}_2}{\mathbf{p}_1} = 2 \rightarrow x \in \omega_1, \text{ то есть}$$

$$\text{если } x < 2 - \frac{\ln 2}{2}, \text{ то } x \in \omega_1;$$

$$\text{если } x > 2 - \frac{\ln 2}{2}, \text{ то } x \in \omega_2.$$

Вероятность ошибки классификации:

$$\begin{aligned}
 \varepsilon &= \mathbf{p}_1 \Pr\{\text{ошибка}|\omega_1\} + \mathbf{p}_2 \Pr\{\text{ошибка}|\omega_2\} = \\
 &= \mathbf{p}_1 \Pr\{X \in G_2|\omega_1\} + \mathbf{p}_2 \Pr\{X \in G_1|\omega_2\} = \\
 &= \frac{1}{3} \int_{2-\ln 2/2}^{+\infty} \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}(x-1)^2\right] dx + \frac{2}{3} \int_{-\infty}^{2-\ln 2/2} \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}(x-2)^2\right] dx = \\
 &= \frac{1}{3} \int_{1-\ln 2/2}^{+\infty} \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}x^2\right] dx + \frac{2}{3} \int_{-\infty}^{-\ln 2/2} \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}x^2\right] dx = \\
 &= \frac{1}{3}(1 - \Phi(1 - \ln 2/2)) + \frac{2}{3}\Phi(-\ln 2/2) \approx 0.4135779.
 \end{aligned}$$

Здесь  $\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z \exp\left(-\frac{x^2}{2}\right) dx$  – интегральная функция

стандартного нормального распределения (интеграл вероятности, функция Лапласа).

При гауссовских наблюдениях  $x$  удобнее использовать логарифм отношения правдоподобия. Тогда решающее правило (5.2) для случая двух классов может быть записано в виде

$$\begin{aligned}
 -\ln l(x) = -\ln \frac{f(x|\omega_1)}{f(x|\omega_2)} < \ln \frac{p_1}{p_2} &\rightarrow x \in \omega_1, \\
 -\ln l(x) = -\ln \frac{f(x|\omega_1)}{f(x|\omega_2)} > \ln \frac{p_1}{p_2} &\rightarrow x \in \omega_2.
 \end{aligned}
 \tag{5.3}$$

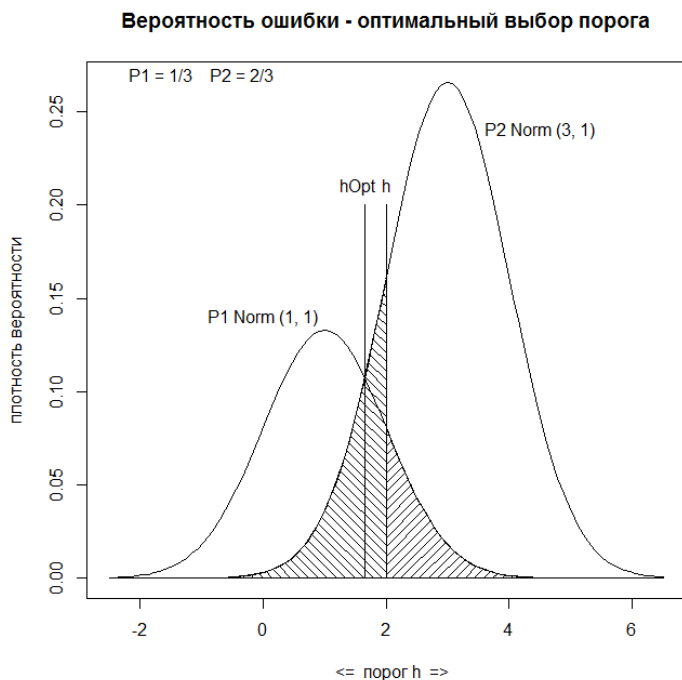


Рис. 5.1. Графическая иллюстрация оптимального выбора порога (к примеру 5.1)

**Пример 5.2.** Построить байесовское решающее правило (5.2) для двух классов независимых гауссовских наблюдений

двумерного случайного вектора  $\mathbf{x}$  с математическими ожиданиями  $\mathbf{a}_1 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$ ,  $\mathbf{a}_2 = \begin{pmatrix} 2 \\ 0 \end{pmatrix}$ , корреляционными матрицами  $\mathbf{R}_1 = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$ ,  $\mathbf{R}_2 = \begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix}$  и априорными вероятностями классов  $\mathbf{p}_1 = p$ ,  $\mathbf{p}_2 = 1 - p$ . Изобразить графически семейство разделяющих поверхностей для различных значений  $p$ .

$$f(\mathbf{x}|\boldsymbol{\omega}_1) = f(x_1, x_2|\boldsymbol{\omega}_1) = \frac{1}{2\pi\sqrt{|\mathbf{R}_1|}} \exp\left[-\frac{1}{2}(\mathbf{x}-\mathbf{a}_1)^T \mathbf{R}_1^{-1}(\mathbf{x}-\mathbf{a}_1)\right],$$

$$f(\mathbf{x}|\boldsymbol{\omega}_2) = f(x_1, x_2|\boldsymbol{\omega}_2) = \frac{1}{2\pi\sqrt{|\mathbf{R}_2|}} \exp\left[-\frac{1}{2}(\mathbf{x}-\mathbf{a}_2)^T \mathbf{R}_2^{-1}(\mathbf{x}-\mathbf{a}_2)\right].$$

Логарифм отношения правдоподобия:

$$\begin{aligned} -\ln l(\mathbf{x}) &= -\ln \frac{f(\mathbf{x}|\boldsymbol{\omega}_1)}{f(\mathbf{x}|\boldsymbol{\omega}_2)} = -\ln \frac{\frac{1}{2\pi\sqrt{|\mathbf{R}_1|}} \exp\left[-\frac{1}{2}(\mathbf{x}-\mathbf{a}_1)^T \mathbf{R}_1^{-1}(\mathbf{x}-\mathbf{a}_1)\right]}{\frac{1}{2\pi\sqrt{|\mathbf{R}_2|}} \exp\left[-\frac{1}{2}(\mathbf{x}-\mathbf{a}_2)^T \mathbf{R}_2^{-1}(\mathbf{x}-\mathbf{a}_2)\right]} = \\ &= -\ln \frac{|\mathbf{R}_2|}{|\mathbf{R}_1|} + \frac{1}{2}(\mathbf{x}-\mathbf{a}_1)^T \mathbf{R}_1^{-1}(\mathbf{x}-\mathbf{a}_1) - \frac{1}{2}(\mathbf{x}-\mathbf{a}_2)^T \mathbf{R}_2^{-1}(\mathbf{x}-\mathbf{a}_2) = \\ &= \frac{1}{2} \begin{pmatrix} x_1 \\ x_2 - 1 \end{pmatrix}^T \begin{pmatrix} 1 & 0 \\ 0 & \frac{1}{2} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 - 1 \end{pmatrix} - \frac{1}{2} \begin{pmatrix} x_1 - 2 \\ x_2 \end{pmatrix}^T \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 - 2 \\ x_2 \end{pmatrix} = \\ &= \frac{1}{2} x_1^2 + \frac{1}{4} (x_2 - 1)^2 - \frac{1}{4} (x_1 - 2)^2 - \frac{1}{2} x_2^2 = \frac{1}{4} (x_1 + 2)^2 - \frac{1}{4} (x_2 + 1)^2 - \frac{3}{4}. \end{aligned}$$

Решающее правило:



Если  $-\ln l(\mathbf{x}) < -\ln \frac{1-p}{p}$ , то  $x \in \omega_1$ ; иначе  $x \in \omega_2$ .

Уравнение разделяющей поверхности:

$$(x_1 + 2)^2 - (x_2 + 1)^2 = 3 - 4 \ln \frac{1-p}{p}.$$

Это есть гипербола с центром в точке  $(-2, -1)$  и асимптотами  $x_2 = x_1 + 1$  и  $x_2 = -x_1 - 3$ .

Центры классов (математические ожидания) обозначены на рис. 5.2 красной (класс 1) и синей (класс 2) точками. График разделяющей поверхности построен для  $p = 0.73$  ( $p_1 = 0.73$ ,  $p_2 = 0.27$ ). Заметим, что точки, лежащие левее левой ветви гиперболы, будут оптимально классифицироваться в класс 2 (апостериорная вероятность выше для второго класса, хотя эти точки лежат ближе к центру первого класса).

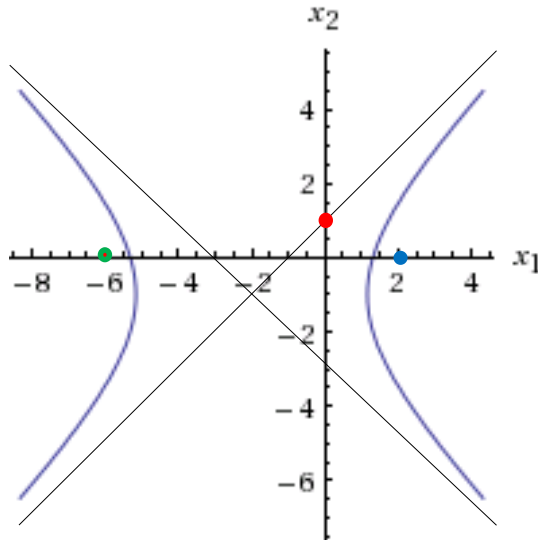


Рис. 5.2. Разделяющая поверхность двух классов гауссовых наблюдений с различными априорными вероятностями

Покажем, что для точки  $\mathbf{x} = \begin{pmatrix} -6 \\ 0 \end{pmatrix}$  (зелёный цвет) выполняется неравенство  $p_2 f(\mathbf{x}|\omega_2) > p_1 f(\mathbf{x}|\omega_1)$  (то есть точка классифицируется в класс 2):

$$p_2 f(\mathbf{x}|\omega_2) - p_1 f(\mathbf{x}|\omega_1) = p_2 \frac{1}{2\pi\sqrt{|\mathbf{R}_2|}} \exp\left[-\frac{1}{2}(\mathbf{x}-\mathbf{a}_2)^T \mathbf{R}_2^{-1}(\mathbf{x}-\mathbf{a}_2)\right] -$$

$$- p_1 \frac{1}{2\pi\sqrt{|\mathbf{R}_1|}} \exp\left[-\frac{1}{2}(\mathbf{x}-\mathbf{a}_1)^T \mathbf{R}_1^{-1}(\mathbf{x}-\mathbf{a}_1)\right].$$

$$\frac{1}{2}(\mathbf{x}-\mathbf{a}_1)^T \mathbf{R}_1^{-1}(\mathbf{x}-\mathbf{a}_1) = \frac{1}{2} \begin{pmatrix} -6 \\ -1 \end{pmatrix}^T \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}^{-1} \begin{pmatrix} -6 \\ -1 \end{pmatrix} = 18.25,$$

$$\frac{1}{2}(\mathbf{x}-\mathbf{a}_2)^T \mathbf{R}_2^{-1}(\mathbf{x}-\mathbf{a}_2) = \frac{1}{2} \begin{pmatrix} -8 \\ 0 \end{pmatrix}^T \begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix}^{-1} \begin{pmatrix} -8 \\ 0 \end{pmatrix} = 16,$$

$$p_2 f(\mathbf{x}|\omega_2) - p_1 f(\mathbf{x}|\omega_1) = \frac{1}{2\pi\sqrt{2}} 0.73e^{-16} - \frac{1}{2\pi\sqrt{2}} 0.27e^{-18.25} =$$

$$= \frac{1}{2\pi\sqrt{2}} (8.215 \cdot 10^{-8} - 3.203 \cdot 10^{-9}) > 0.$$

**Пример 5.3.** Найти разделяющие поверхности трёх классов **I**, **II**, **III** гауссовских наблюдений двумерного случайного вектора  $X$  с математическими ожиданиями  $\mathbf{a}_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ ,  $\mathbf{a}_2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$ ,

$\mathbf{a}_3 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ , некоррелированными компонентами  $\mathbf{R}_1 = \mathbf{R}_2 = \mathbf{R}_3 = \mathbf{R} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$

и априорными вероятностями классов  $p_1 = 1/2$ ,  $p_2 = 1/3$ ,  $p_3 = 1/6$ . Плотность вероятности каждого класса

$$f(\mathbf{x}|\omega_i) = f(x_1, x_2|\omega_i) = \frac{1}{2\pi\sqrt{|\mathbf{R}|}} \exp\left[-\frac{1}{2}(\mathbf{x}-\mathbf{a}_i)^T \mathbf{R}^{-1}(\mathbf{x}-\mathbf{a}_i)\right],$$

$i = 1, 2, 3$ .

Логарифм функции правдоподобия для разделения объектов первого и второго классов:

$$\begin{aligned}
 -\ln \frac{f(\mathbf{x}|\omega_1)}{f(\mathbf{x}|\omega_2)} &= -\ln \frac{\frac{1}{2\pi\sqrt{|\mathbf{R}|}} \exp\left[-\frac{1}{2}(\mathbf{x}-\mathbf{a}_1)^T \mathbf{R}^{-1}(\mathbf{x}-\mathbf{a}_1)\right]}{\frac{1}{2\pi\sqrt{|\mathbf{R}|}} \exp\left[-\frac{1}{2}(\mathbf{x}-\mathbf{a}_2)^T \mathbf{R}^{-1}(\mathbf{x}-\mathbf{a}_2)\right]} = \\
 &= \frac{1}{2}(\mathbf{x}-\mathbf{a}_1)^T \mathbf{R}^{-1}(\mathbf{x}-\mathbf{a}_1) - \frac{1}{2}(\mathbf{x}-\mathbf{a}_2)^T \mathbf{R}^{-1}(\mathbf{x}-\mathbf{a}_2) = \\
 &= \frac{1}{2}\|\mathbf{x}-\mathbf{a}_1\|^2 - \frac{1}{2}\|\mathbf{x}-\mathbf{a}_2\|^2 = \frac{1}{2}\left[(x_1-1)^2 + x_2^2 - x_1^2 - (x_2-1)^2\right] = x_2 - x_1.
 \end{aligned}$$

Разделяющая поверхность первого и второго классов находится из уравнения

$$-\ln \frac{f(\mathbf{x}|\omega_1)}{f(\mathbf{x}|\omega_2)} = \ln \frac{\mathbf{p}_1}{\mathbf{p}_2} \rightarrow x_2 - x_1 = \ln \frac{3}{2} \rightarrow x_2 = x_1 + \ln \frac{3}{2}.$$

Логарифм функции правдоподобия для разделения объектов первого и третьего классов:

$$\begin{aligned}
 -\ln \frac{f(\mathbf{x}|\omega_1)}{f(\mathbf{x}|\omega_3)} &= \frac{1}{2}\|\mathbf{x}-\mathbf{a}_1\|^2 - \frac{1}{2}\|\mathbf{x}-\mathbf{a}_3\|^2 = \\
 &= \frac{1}{2}\left[(x_1-1)^2 + x_2^2 - (x_1-1)^2 - (x_2-1)^2\right] = x_2 - \frac{1}{2}.
 \end{aligned}$$

Разделяющая поверхность первого и третьего классов находится из уравнения

$$-\ln \frac{f(\mathbf{x}|\omega_1)}{f(\mathbf{x}|\omega_3)} = \ln \frac{\mathbf{p}_1}{\mathbf{p}_3} \rightarrow x_2 - \frac{1}{2} = \ln 3 \rightarrow x_2 = \frac{1}{2} + \ln 3.$$

Логарифм функции правдоподобия для разделения объектов второго и третьего классов:

$$\begin{aligned}
 -\ln \frac{f(\mathbf{x}|\omega_2)}{f(\mathbf{x}|\omega_3)} &= \frac{1}{2} \|\mathbf{x} - \mathbf{a}_2\|^2 - \frac{1}{2} \|\mathbf{x} - \mathbf{a}_3\|^2 = \\
 &= \frac{1}{2} [x_1^2 + (x_2 - 1)^2 - (x_1 - 1)^2 - (x_2 - 1)^2] = x_1 - \frac{1}{2}.
 \end{aligned}$$

Разделяющая поверхность второго и третьего классов находится из уравнения

$$-\ln \frac{f(\mathbf{x}|\omega_2)}{f(\mathbf{x}|\omega_3)} = \ln \frac{\mathbf{p}_2}{\mathbf{p}_3} \rightarrow x_1 - \frac{1}{2} = \ln 2 \rightarrow x_1 = \frac{1}{2} + \ln 2.$$

На рис. 5.3 цветными точками обозначены центры (математические ожидания) классов. Видим, что из-за неодинаковой априорной вероятности классов происходит смещение области решений в сторону класса с низкой априорной вероятностью.

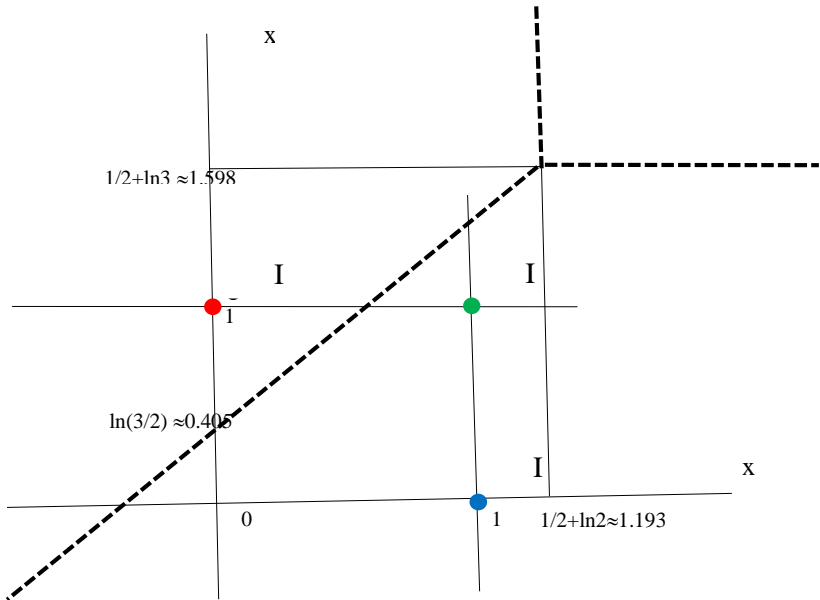


Рис. 5.3. Разделяющие прямые

Например, математическое ожидание третьего класса – точка  $\mathbf{a}_3$  – будет оптимально (с минимальной вероятностью ошибки) классифицирована как принадлежащая первому классу, так как

$$\begin{aligned} \mathbf{p}_1 f(\mathbf{x}|\boldsymbol{\omega}_1)|_{\mathbf{x}=\mathbf{a}_3} > \mathbf{p}_3 f(\mathbf{x}|\boldsymbol{\omega}_3)|_{\mathbf{x}=\mathbf{a}_3} &\Leftrightarrow \frac{1}{2} \frac{1}{2\pi\sqrt{|\mathbf{R}|}} \exp\left[-\frac{1}{2}\|\mathbf{a}_3-\mathbf{a}_1\|^2\right] > \\ > \frac{1}{6} \frac{1}{2\pi\sqrt{|\mathbf{R}|}} e^0 &\Leftrightarrow \frac{1}{\sqrt{e}} > \frac{1}{3}. \end{aligned}$$

## 5.2 Минимизация функции риска

Четыре типа штрафов:  $c_{ij}$ ,  $i=1,2$ ,  $j=1,2$  – штраф за решенные  $x \in \boldsymbol{\omega}_j$ , если в действительности  $x \in \boldsymbol{\omega}_i$ . Риск  $\mathbf{r}$  – математическое ожидание штрафа:

$$\begin{aligned} \mathbf{r} = c_{11}\mathbf{p}_1\Pr\{x \in G_1|\boldsymbol{\omega}_1\} + c_{21}\mathbf{p}_2\Pr\{x \in G_1|\boldsymbol{\omega}_2\} + c_{12}\mathbf{p}_1\Pr\{x \in G_2|\boldsymbol{\omega}_1\} + \\ + c_{22}\mathbf{p}_2\Pr\{x \in G_2|\boldsymbol{\omega}_2\}. \end{aligned}$$

Риск можно минимизировать за счёт оптимального выбора областей принятия решений  $G_1$ ,  $G_2$ :

$$\begin{aligned} \mathbf{r} = c_{11}\mathbf{p}_1 \int_{G_1} f(x|\boldsymbol{\omega}_1) dx + c_{21}\mathbf{p}_2 \int_{G_1} f(x|\boldsymbol{\omega}_2) dx + c_{12}\mathbf{p}_1 \int_{G_2} f(x|\boldsymbol{\omega}_1) dx + c_{22}\mathbf{p}_2 \int_{G_2} f(x|\boldsymbol{\omega}_2) dx = \\ = c_{11}\mathbf{p}_1 \int_{G_1} f(x|\boldsymbol{\omega}_1) dx + c_{21}\mathbf{p}_2 \int_{G_1} f(x|\boldsymbol{\omega}_2) dx + c_{12}\mathbf{p}_1 \left(1 - \int_{G_1} f(x|\boldsymbol{\omega}_1) dx\right) + c_{22}\mathbf{p}_2 \left(1 - \int_{G_1} f(x|\boldsymbol{\omega}_2) dx\right) = (5.4) \\ = c_{12}\mathbf{p}_1 + c_{22}\mathbf{p}_2 + \int_{G_1} [-(c_{12}-c_{11})\mathbf{p}_1 f(x|\boldsymbol{\omega}_1) + (c_{21}-c_{22})\mathbf{p}_2 f(x|\boldsymbol{\omega}_2)] dx \rightarrow \min_{G_1}. \end{aligned}$$

При выводе соотношения (5.4) использовался тот факт, что события  $x \in G_1$  и  $x \in G_2$  образуют полную группу:

$$\int_{G_2} f(x|\boldsymbol{\omega}_i) dx + \int_{G_1} f(x|\boldsymbol{\omega}_i) dx = 1.$$

Для случая  $c_{21} > c_{22}$ ,  $c_{12} > c_{11}$ , минимизируя подинтегральное выражение в (5.4) в каждой точке  $x$ , получаем по аналогии с (5.2) **байесовский критерий, минимизирующий риск:**

$$\begin{aligned}
 l(x) &= \frac{f(x|\omega_1)}{f(x|\omega_2)} > \frac{(c_{21} - c_{22})\mathbf{P}_2}{(c_{12} - c_{11})\mathbf{P}_1} \rightarrow x \in \omega_1, \\
 l(x) &= \frac{f(x|\omega_1)}{f(x|\omega_2)} \leq \frac{(c_{21} - c_{22})\mathbf{P}_2}{(c_{12} - c_{11})\mathbf{P}_1} \rightarrow x \in \omega_2.
 \end{aligned}
 \tag{5.5}$$

Сравнив (5.2) и (5.5), видим, что назначение штрафов эквивалентно изменению отношения априорных вероятностей (порога функции правдоподобия). В случае симметричной функции штрафа ( $c_{21} - c_{22} = c_{12} - c_{11}$ ) получаем байесовский критерий (5.2), минимизирующий вероятность ошибки.

### 5.3 Критерий Неймана-Пирсона

Одним из существенных недостатков байесовского правила обнаружения сигналов является большое количество априорной информации о потерях и вероятностях состоянии объекта, которая должна быть в распоряжении наблюдателя. Иногда указать априорные вероятности наличия цели и потери за счет ложной тревоги или пропуска цели оказывается весьма затруднительным. Поэтому в подобных задачах вместо байесовского критерия обычно используется критерий Неймана-Пирсона. Согласно этому критерию выбирается такое правило обнаружения, которое обеспечивает минимальную величину вероятности пропуска сигнала, то есть минимальную вероятность ошибки первого рода  $\varepsilon_1$  (максимальную вероятность правильного обнаружения), при условии, что вероятность ложной тре-

воги (вероятность ошибки второго рода  $\varepsilon_2$ ) не превышает заданной величины.

Минимизируется вероятность ошибки  $\varepsilon_1$  при условии, что вероятность ошибки  $\varepsilon_2$  равна заданной величине (например,  $\varepsilon_0$ ). Для определения решающего правила минимизируем выражение

$$\mathbf{r} = \varepsilon_1 + \lambda(\varepsilon_2 - \varepsilon_0),$$

где  $\lambda$  – множитель Лагранжа.

$$\begin{aligned} \mathbf{r} = \varepsilon_1 + \lambda(\varepsilon_2 - \varepsilon_0) &= \int_{G_2} f(\mathbf{x}|\omega_1) d\mathbf{x} + \lambda \left( \int_{G_1} f(\mathbf{x}|\omega_2) d\mathbf{x} - \varepsilon_0 \right) = \\ &= (1 - \lambda\varepsilon_0) + \int_{G_1} [\lambda f(\mathbf{x}|\omega_2) - f(\mathbf{x}|\omega_1)] d\mathbf{x} \rightarrow \min_{G_1}. \end{aligned}$$

Решающее правило:

$$\begin{aligned} l(\mathbf{x}) = \frac{f(\mathbf{x}|\omega_1)}{f(\mathbf{x}|\omega_2)} > \lambda &\rightarrow \mathbf{x} \in \omega_1, \\ l(\mathbf{x}) = \frac{f(\mathbf{x}|\omega_1)}{f(\mathbf{x}|\omega_2)} \leq \lambda &\rightarrow \mathbf{x} \in \omega_2. \end{aligned} \tag{5.6}$$

Из сравнения (5.6) и (5.2) можно заключить, что оптимальное, в смысле критерия Неймана-Пирсона, правило обнаружения отличается от байесовского лишь величиной порогового уровня, с которым производится сравнение отношения правдоподобия. Порог  $\lambda$  находится из уравнения  $\frac{\partial \mathbf{r}}{\partial \lambda} = 0$ , то есть

$$\int_{G_1(\lambda)} f(\mathbf{x}|\omega_2) d\mathbf{x} = \varepsilon_0 = \varepsilon_2.$$

**Пример 5.4.** Построить решающее правило Неймана-Пирсона (5.6) для двух классов независимых гауссовских наблюдений случайной величины  $X$  с условной плотностью вероятностей

$$f(x|\omega_1) = \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}(x-1)^2\right],$$

$$f(x|\omega_2) = \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}(x-3)^2\right].$$

Отношение правдоподобия:

$$\begin{aligned} l(x) &= \frac{f_x(x|\omega_1)}{f_x(x|\omega_2)} = \frac{\frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}(x-1)^2\right]}{\frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}(x-3)^2\right]} = \\ &= \exp\left[\frac{1}{2}(x-3)^2 - \frac{1}{2}(x-1)^2\right] = \exp(4-2x). \end{aligned}$$

Решающее правило:

$$x < h \rightarrow x \in \omega_1.$$

Порог  $h$  находится из уравнения, фиксирующего вероятность ошибки второго рода:

$$\int_{-\infty}^h f(x|\omega_2) dx = \varepsilon_2, \quad \frac{1}{\sqrt{2\pi}} \int_{-\infty}^h \exp\left[-\frac{1}{2}(x-3)^2\right] dx = \varepsilon_2,$$

$$\Phi(h-3) = \varepsilon_2, \quad h = \Phi^{-1}(\varepsilon_2) + 3.$$

Здесь, как и раньше, через  $\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z \exp\left(-\frac{x^2}{2}\right) dx$  обозначена интегральная функция стандартного нормального рас-



предела (интеграл вероятности, функция Лапласа). Тогда вероятность ошибки первого рода:

$$\begin{aligned} \varepsilon_1 &= \int_h^{+\infty} f(x|\omega_1) dx = \int_{\Phi^{-1}(\varepsilon_2)+3}^{+\infty} \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}(x-1)^2\right] dx = \\ &= \int_{\Phi^{-1}(\varepsilon_2)+2}^{+\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}u^2\right) du = 1 - \Phi(\Phi^{-1}(\varepsilon_2)+2), \end{aligned}$$

то есть:  $\varepsilon_1 = 1 - \Phi(\Phi^{-1}(\varepsilon_2) + 2)$ . (5.7)

Уравнение (5.7) задаёт зависимость вероятности ошибки первого рода от установленного значения вероятности ошибки второго рода. В табл. 5.1 приведены значения связанных по этому уравнению вероятностей ошибок первого и второго рода.

Таблица 5.1 – Соотношение между вероятностями ошибок первого и второго рода

$\varepsilon_2$	0.010	0.050	0.100	0.500	0.900	0.950
$h$	0.674	1.355	1.718	3.000	4.282	4.645
$\varepsilon_1$	0.628	0.361	0.236	0.023	0.001	0.000

На рис. 5.1 приведена графическая иллюстрация соотношения между вероятностями ошибок первого и второго рода.

**Пример 5.5.** Построить решающее правило Неймана-Пирсона (5.5) для двух классов независимых гауссовских наблюдений двумерного случайного вектора  $X$  с математическими ожиданиями  $\mathbf{a}_1 = (-1 \ 0)^T$ ,  $\mathbf{a}_2 = (1 \ 0)^T$ , корреляционными матрицами  $\mathbf{R}_1 = \mathbf{R}_2 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ .

Логарифм отношения правдоподобия:

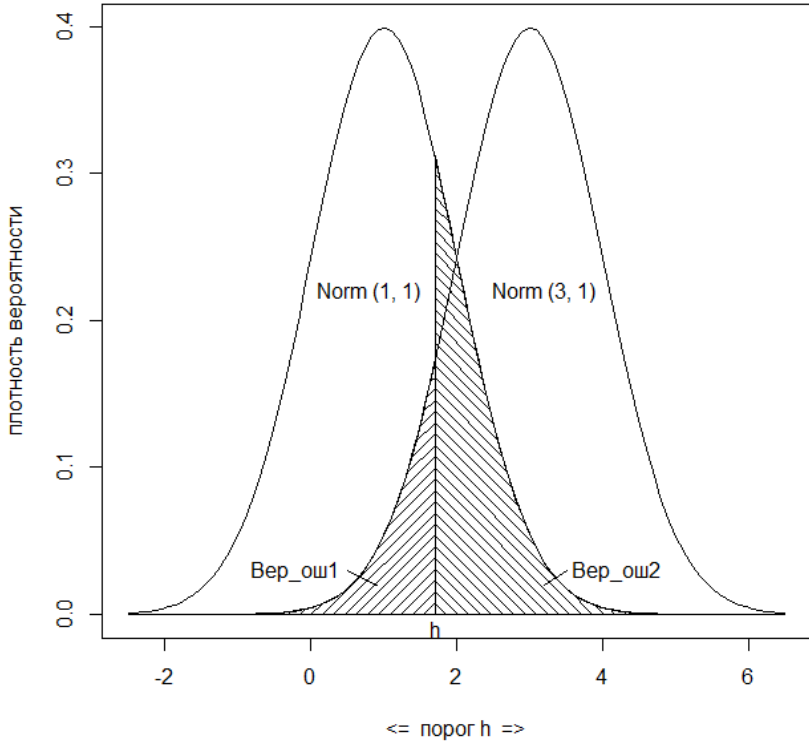


Рис. 5.4. Решающее правило Неймана–Пирсона.

Соотношение между вероятностями ошибок первого и второго рода

$$\begin{aligned}
 -\ln l(\mathbf{x}) &= -\ln \frac{f(\mathbf{x}|\omega_1)}{f(\mathbf{x}|\omega_2)} = -\ln \frac{\frac{1}{2\pi\sqrt{|\mathbf{R}_1|}} \exp\left[-\frac{1}{2}(\mathbf{x}-\mathbf{a}_1)^T \mathbf{R}_1^{-1}(\mathbf{x}-\mathbf{a}_1)\right]}{\frac{1}{2\pi\sqrt{|\mathbf{R}_2|}} \exp\left[-\frac{1}{2}(\mathbf{x}-\mathbf{a}_2)^T \mathbf{R}_2^{-1}(\mathbf{x}-\mathbf{a}_2)\right]} = \\
 &= -\frac{1}{2} \begin{pmatrix} x_1 - 1 \\ x_2 \end{pmatrix}^T \begin{pmatrix} x_1 - 1 \\ x_2 \end{pmatrix} + \frac{1}{2} \begin{pmatrix} x_1 + 1 \\ x_2 \end{pmatrix}^T \begin{pmatrix} x_1 + 1 \\ x_2 \end{pmatrix} = 2x_1 < \ln \lambda \rightarrow \mathbf{x} \in \omega_1.
 \end{aligned}$$

Порог  $\lambda$  находится из уравнения  $\int_{2x_1 > \ln \lambda} f(\mathbf{x}|\omega_1) d\mathbf{x} = \varepsilon_2$ , т.е.

$$\begin{aligned}\varepsilon_2 &= \frac{1}{2\pi} \int_{+\infty}^{+\infty} \left( \int_{\frac{\ln \lambda}{2}}^{+\infty} \exp \left[ -\frac{1}{2} \left( (x_1 + 1)^2 + x_2^2 \right) \right] dx_1 \right) dx_2 = \\ &= \frac{1}{\sqrt{2\pi}} \int_{+\infty}^{+\infty} \left( \frac{1}{\sqrt{2\pi}} \int_{\frac{\ln \lambda}{2} + 1}^{+\infty} \exp \left[ -\frac{1}{2} x_1^2 \right] dx_1 \right) \exp \left( -\frac{1}{2} x_2^2 \right) dx_2 = 1 - \Phi \left( \frac{\ln \lambda}{2} + 1 \right).\end{aligned}$$

Вероятность ошибки первого рода находится из уравнения

$$\varepsilon_1 = \int_{2x_1 < \ln \lambda} f_X(\mathbf{x} | \omega_2) d\mathbf{x}, \text{ т.е.}$$

$$\begin{aligned}\varepsilon_1 &= \frac{1}{2\pi} \int_{+\infty}^{+\infty} \left( \int_{-\infty}^{\frac{\ln \lambda}{2}} \exp \left[ -\frac{1}{2} \left( (x_1 - 1)^2 + x_2^2 \right) \right] dx_1 \right) = \\ &= \frac{1}{\sqrt{2\pi}} \int_{+\infty}^{+\infty} \left( \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\frac{\ln \lambda}{2} - 1} \exp \left( -\frac{1}{2} x_1^2 \right) dx_1 \right) \exp \left( -\frac{1}{2} x_2^2 \right) dx_2 = \Phi \left( \frac{\ln \lambda}{2} - 1 \right).\end{aligned}$$

Из последних двух соотношений  $\varepsilon_1 = \Phi \left( \frac{\ln \lambda}{2} - 1 \right)$ ,

$\varepsilon_2 = 1 - \Phi \left( \frac{\ln \lambda}{2} + 1 \right)$  получаем явную зависимость  $\varepsilon_1 = \varepsilon_1(\varepsilon_2)$ :

$$\varepsilon_1 = \Phi \left( \Phi^{-1}(1 - \varepsilon_2) - 2 \right).$$

В табл. 5.2 приведён ряд рассчитанных значений вероятностей ошибок первого и второго рода для некоторых значений порога  $\lambda$ .

Таблица 5.2 – Соотношение между вероятностями ошибок первого и второго рода

$\lambda$	8	4	2	1	1/2	1/4	1/8
$\varepsilon_1$	0.516	0.379	0.257	0.159	0.090	0.045	0.021
$\varepsilon_2$	0.021	0.045	0.090	0.159	0.257	0.379	0.516

## 5.4 ROC-кривая

В предыдущих разделах мы видели, что вероятности ошибок первого и второго рода изменяются в зависимости от выбранного порога, но жестко связаны друг с другом. Их зависимость традиционно представляется в виде так называемой *ROC-кривой* (*Receiver Operating Characteristic*, рабочая характеристика приёмника).

Эта кривая позволяет оценить качество бинарной классификации, отображает соотношение между долей объектов от общего количества носителей признака, верно классифицированных как несущих признак (*True Positive Rate*, *TPR*, называемой **чувствительностью** (*Sensitivity*) алгоритма классификации, в наших обозначениях  $TPR = 1 - \varepsilon_1$ ) и долей объектов, не несущих признака, ошибочно классифицированных как несущих признак (*False Positive Rate*, *FPR*, величина  $(1 - FPR)$  называется **специфичностью** (*Specificity*) алгоритма классификации, в наших обозначениях  $FPR = \varepsilon_2$ ) при варьировании порога решающего правила.

На рис. 5.5 приведена ROC-кривая для решающего правила Неймана-Пирсона для двух классов независимых гауссовских наблюдений двумерного случайного вектора (см. пример 5.4, уравнение (5.7)):  $TPR = \Phi(\Phi^{-1}(FPR) - 2)$ .

Площадь под ROC-кривой ( $AUC$  – *Area Under Curve*)

$$AUC = \int_0^1 TPR(FPR) dFPR$$
 позволяет оценить качество используемого алгоритма классификации при варьировании порога решающего правила.

зубаемого алгоритма классификации при варьировании порога решающего правила.

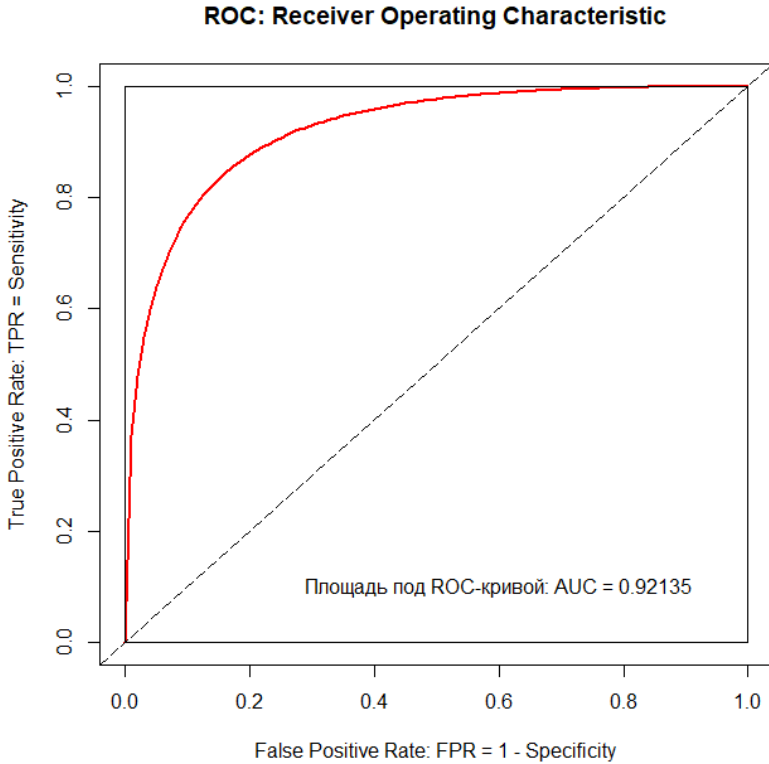


Рис. 5.5. ROC-кривая для ошибок классификации из примера 5.4

## 5.5 Наивный байесовский классификатор

*Наивный* байесовский классификатор основан на наивном предположении, что объекты описываются *независимыми признаками*.

Будем считать, что  $n$  вероятностно независимых признаков (не обязательно числовых)  $(x_1, x_2, \dots, x_n)$  представлены вектором  $\mathbf{x}$ . Вероятностная независимость означает, что многомерное распределение вероятностей равно произведению одномерных распределений:

$$P(\mathbf{x}) = P(x_1)P(x_2)\dots P(x_n).$$

В случае условных вероятностей, опуская обозначения признаков, получаем:

$$P(\mathbf{x}|\omega_i) = P(x_1|\omega_i)P(x_2|\omega_i)\dots P(x_n|\omega_i), i = \overline{1, m} - \text{номер класса.}$$

Запишем формулу Байеса для **апостериорной вероятности** класса в виде

$$P(\omega_i|\mathbf{x}) = \frac{P(\mathbf{x}|\omega_i)P(\omega_i)}{P(\mathbf{x})} = \frac{P(\mathbf{x}|\omega_i)P(\omega_i)}{P(\mathbf{x}|\omega_1)P(\omega_1) + \dots + P(\mathbf{x}|\omega_m)P(\omega_m)},$$

$$i = \overline{1, m},$$

где  $P(\omega_i)$  – **априорные** вероятности классов объектов;

$P(\mathbf{x}|\omega_i)$  – условные распределения вероятности признаков для различных классов объектов;

$P(\mathbf{x}) = P(\omega_1)P(\mathbf{x}|\omega_1) + \dots + P(\omega_m)P(\mathbf{x}|\omega_m)$  – распределение вероятностей признаков (формула полной вероятности).

Так как в знаменателе формулы Байеса стоит величина вероятности, не зависящая от класса объекта  $\mathbf{x}$ , то класс объекта, соответствующий максимуму апостериорной вероятности, определяется числителем формулы Байеса:

$$\omega_* = \arg \max_{\omega_1, \dots, \omega_m} \{P\{\omega_1\}P(\mathbf{x}|\omega_1), \dots, P\{\omega_m\}P(\mathbf{x}|\omega_m)\}.$$

Для случая двух классов получаем:

$$\omega_* = \arg \max_{\omega_1, \omega_2} \left\{ P(\omega_1)P(\mathbf{x}|\omega_1), P(\omega_2)P(\mathbf{x}|\omega_2) \right\}.$$

**Пример 5.6. Байесовская фильтрация спама.** Обозначим события:  $S$  – рассматриваемое сообщение является спамом (*Spam*, класс  $\omega_1$ ),  $\bar{S}$  – сообщение не является спамом (*non-spam* или *ham*, класс  $\omega_2$ ). Недавние статистические исследования показали, что на сегодняшний день вероятность любого сообщения быть спамом составляет по меньшей мере 80%:  $P\{S\} = 0.8$ ,  $P\{\bar{S}\} = 0.2$ . Однако большинство байесовских программ обнаружения спама делают предположение об отсутствии априорных предпочтений у сообщения быть спамом, а не «ham», и полагают, что у обеих альтернатив есть равные вероятности 50%:  $P\{S\} = P\{\bar{S}\} = 0.5$  (ещё одно наивное предположение). О фильтрах спама, которые используют эту гипотезу, говорят как о фильтрах «без предубеждений». Это означает, что у них нет никакого предубеждения относительно входящей электронной почты. Это предположение позволяет упростить общую формулу (считаем, что  $P\{S\} = P\{\bar{S}\} = 0.5$ ):

$$P\{S|\mathbf{x}\} = \frac{P(\mathbf{x}|S)P\{S\}}{P(\mathbf{x}|S)P\{S\} + P(\mathbf{x}|\bar{S})P\{\bar{S}\}}, \text{ то есть}$$

$$\Pr\{S|x_1, x_2, \dots, x_n\} = \frac{P(x_1|S)P(x_2|S)\dots P(x_n|S)}{P(x_1|S)P(x_2|S)\dots P(x_n|S) + P(x_1|\bar{S})P(x_2|\bar{S})\dots P(x_n|\bar{S})}.$$

Здесь  $x_1, x_2, \dots, x_n$  – слова в проверяемом на спам сообщении;  $P(x_k|S)$  – оценка вероятности наличия слова  $x_k$  в спамовом сообщении, определяемой относительной частотой сооб-

щений, содержащих слово  $x_k$  в сообщениях, идентифицированных как спам во время фазы обучения;  $P(x_k | \bar{S})$  – оценка вероятности наличия слова  $x_k$  в неспамовом сообщении, определяемой относительной частотой сообщений, содержащих слово  $x_k$  в сообщениях, идентифицированных как неспам во время фазы обучения.

$$\begin{aligned}
 P\{S|\mathbf{x}\} &= \frac{\frac{P(S|x_1)}{P(S)} \cdot \frac{P(S|x_2)}{P(S)} \cdot \dots \cdot \frac{P(S|x_n)}{P(S)}}{\frac{P(S|x_1)}{P(S)} \cdot \frac{P(S|x_2)}{P(S)} \cdot \dots \cdot \frac{P(S|x_n)}{P(S)} + \frac{P(\bar{S}|x_1)}{P(\bar{S})} \cdot \frac{P(\bar{S}|x_2)}{P(\bar{S})} \cdot \dots \cdot \frac{P(\bar{S}|x_n)}{P(\bar{S})}} = \\
 &= \frac{p_1 p_2 \dots p_n}{p_1 p_2 \dots p_n + (1-p_1)(1-p_2) \dots (1-p_n)},
 \end{aligned}$$

где  $p_k = \mathbf{P}(S|x_k)$  – условная вероятность того, что сообщение является спамом при условии, что оно содержит  $k$ -е слово ( $k = \overline{1, n}$ ) (тогда  $P(\bar{S}|x_k) = 1 - p_k$ ). То есть  $p_1$  – условная вероятность того, что сообщение  $S$  – спам при условии, что оно содержит слово  $x_1$ ;  $p_2$  – условная вероятность того, что сообщение  $S$  – спам при условии, что оно содержит слово  $x_2$  и т.д.



## 6 ЛИНЕЙНЫЕ КЛАССИФИКАТОРЫ

Байесовский классификатор, рассмотренный в предыдущем разделе, основан на вычислении отношения правдоподобия и является оптимальным в смысле минимизации риска или вероятности ошибки. Более простые методы классификации основаны на задании математического вида классификатора с последующим определением его оптимальных параметров. Наиболее простым является *линейный* или *кусочно-линейный* классификатор. Следует, однако, иметь в виду, что никакой линейный классификатор по качеству работы не превосходит классификатор, основанный на вычислении отношения правдоподобия.

### 6.1 Байесовский линейный классификатор

Для случайного вектора признаков  $\mathbf{x} \in R^n$  при двух классах гауссовских распределений с параметрами  $(\mathbf{M}_1, R_1)$  и  $(\mathbf{M}_2, R_2)$ , с одинаковыми ковариационными матрицами  $R_1 = R_2 = R$  и с заданными априорными вероятностями классов  $\mathbf{p}_1 = \mathbf{Pr}\{\omega_1\}$  и  $\mathbf{p}_2 = \mathbf{Pr}\{\omega_2\}$  и условными плотностями распределений  $f_1(\mathbf{x}) = f(\mathbf{x}|\omega_1)$ ,  $f_2(\mathbf{x}) = f(\mathbf{x}|\omega_2)$  плотность вероятности смеси распределений имеет вид:

$$f(\mathbf{x}) = \mathbf{p}_1 f_1(\mathbf{x}) + \mathbf{p}_2 f_2(\mathbf{x}) = \mathbf{p}_1 \frac{1}{(2\pi)^{n/2} |R|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \mathbf{M}_1)^T R^{-1}(\mathbf{x} - \mathbf{M}_1)\right] + \mathbf{p}_2 \frac{1}{(2\pi)^{n/2} |R|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \mathbf{M}_2)^T R^{-1}(\mathbf{x} - \mathbf{M}_2)\right].$$

Байесовское решающее правило имеет линейный относительно вектора наблюдений  $\mathbf{x}$  вид:

$$-\ln l(\mathbf{x}) = -\ln \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} = \frac{1}{2}(\mathbf{x} - \mathbf{M}_1)^T R^{-1}(\mathbf{x} - \mathbf{M}_1) - \frac{1}{2}(\mathbf{x} - \mathbf{M}_2)^T R^{-1}(\mathbf{x} - \mathbf{M}_2) < \ln \frac{\mathbf{p}_1}{\mathbf{p}_2} \rightarrow \omega_1,$$

то есть

$$(\mathbf{M}_2 - \mathbf{M}_1)^T R^{-1} \mathbf{x} + \frac{1}{2}(\mathbf{M}_1^T R^{-1} \mathbf{M}_1 - \mathbf{M}_2^T R^{-1} \mathbf{M}_2) < \ln(\mathbf{p}_1/\mathbf{p}_2) \rightarrow \omega_1.$$

Уравнение разделяющей поверхности:

$$(\mathbf{M}_2 - \mathbf{M}_1)^T R^{-1} \mathbf{x} = \ln(\mathbf{p}_1/\mathbf{p}_2) - \frac{1}{2}(\mathbf{M}_1^T R^{-1} \mathbf{M}_1 - \mathbf{M}_2^T R^{-1} \mathbf{M}_2).$$

Это есть гиперплоскость в пространстве  $\mathbf{R}^n$ .

Если компоненты вектора наблюдений  $X$  являются некоррелированными ( $R = \sigma^2 I_n$ ,  $\sigma^2$  – дисперсия белого шума наблюдения), то из

$$\frac{1}{2}(\mathbf{x} - \mathbf{M}_1)^T \sigma^{-2} I_n (\mathbf{x} - \mathbf{M}_1) - \frac{1}{2}(\mathbf{x} - \mathbf{M}_2)^T \sigma^{-2} I_n (\mathbf{x} - \mathbf{M}_2) < \ln(\mathbf{p}_1/\mathbf{p}_2) \rightarrow \omega_1$$

получаем:

$$\|\mathbf{x} - \mathbf{M}_1\|^2 - \|\mathbf{x} - \mathbf{M}_2\|^2 < 2\sigma^2 \ln(\mathbf{p}_1/\mathbf{p}_2) \rightarrow \omega_1,$$

где  $\|\cdot\|$  – евклидова норма вектора.

Уравнение разделяющей поверхности:

$$\|\mathbf{x} - \mathbf{M}_1\|^2 - \|\mathbf{x} - \mathbf{M}_2\|^2 = 2\sigma^2 \ln(\mathbf{p}_1/\mathbf{p}_2).$$

*Геометрическая интерпретация.* Это есть точка  $x = (\mathbf{M}_1 + \mathbf{M}_2)/2 + \sigma^2 \ln(\mathbf{p}_1/\mathbf{p}_2)$  в случае  $n=1$ , прямая в случае  $n=2$ , плоскость в случае  $n=3$  или гиперплоскость при  $n>3$ , перпендикулярная отрезку, соединяющему точки  $\mathbf{M}_1$  и  $\mathbf{M}_2$ . При  $\mathbf{p}_1 = \mathbf{p}_2 = \mathbf{1}/2$  разделяющая поверхность проходит через середину этого отрезка.

В случае, когда наблюдения не являются некоррелированными, т.е. когда  $R_X \neq \sigma^2 I_n$ , также можно использовать принцип минимизации расстояния. Для этого проведём декорреляцию вектора наблюдений с помощью некоторого линейного преобразования  $A: Y = AX$ :

$$R_Y = \mathbf{M}Y_0Y_0^T = \mathbf{M}AX_0X_0^T A^T = \mathbf{M}(X_0X_0^T)A^T = AR_X A^T = I_n.$$

Так как ковариационная матрица является положительно определённой, то матрица линейного преобразования  $A$  всегда существует и невырождена, т.е. преобразование  $A$  является обратимым. Поэтому можно классифицировать гауссовский случайный вектор  $Y$  с условными математическими ожиданиями  $\mathbf{M}_1^Y = \mathbf{A}\mathbf{M}_1$  и  $\mathbf{M}_2^Y = \mathbf{A}\mathbf{M}_2$  и ковариационной матрицей

$$R_Y = I_n, \text{ сравнивая расстояния } \|Y - \mathbf{M}_1^Y\|^2 \text{ и } \|Y - \mathbf{M}_2^Y\|^2.$$

## 6.2 Линейная разделяющая функция, минимизирующая вероятность ошибки

Во многих прикладных задачах предположение о равенстве ковариационных функций при наблюдении объектов различных классов не выполняется. Рассмотрим задачу построения линейного классификатора, который определяется *линейной разделяющей функцией*

$$h(x) = V^T x + v_0 : \begin{cases} < 0 \rightarrow x \in \omega_1, \\ \geq 0 \rightarrow x \in \omega_2. \end{cases}$$

Задача синтеза классификатора заключается в определении коэффициентов  $V^T = [v_1 \dots v_n]$  и  $v_0$ , оптимальных по некоторому критерию. При нормальных распределениях с одинаковыми ковариационными матрицами линейная разделяющая функция  $h(x)$  совпадает с отношением правдоподобия.

Построим линейную разделяющую функцию, минимизирующую вероятность ошибки классификации (как мы видели выше, в общем случае разделяющая функция не является линейной). Даже если наблюдения не подчиняются нормальному закону, то в силу центральной предельной теоремы величина  $h(X)$  имеет распределение, близкое к нормальному при  $n \rightarrow \infty$ . Условные математическое ожидание и дисперсия скалярной величины  $h(X)$  для первого и второго классов объектов равны:

$$\begin{aligned} m_i &= \mathbf{M}(h(X) | \omega_i) = \mathbf{M}(V^T X | \omega_i) + v_0 = V^T \mathbf{M}(X | \omega_i) + v_0 = \\ &= V^T \mathbf{M}_i + v_0, \quad i = \overline{1, 2}, \end{aligned}$$

$$\begin{aligned}\sigma_i^2 &= \mathbf{D}\{h(X) | \omega_i\} = \mathbf{D}\{V^T X | \omega_i\} = \\ &= \mathbf{M}\{V^T (X - m_i)(X - m_i)^T V | \omega_i\} = V^T R_i V, \quad i = \overline{1, 2},\end{aligned}$$

где  $\mathbf{M}_i$  и  $R_i$  – условное математическое ожидание и условная корреляционная матрица (при условии принадлежности к классу  $i$ ) вектора признаков  $X$ .

Вероятность ошибки:

$$\begin{aligned}\varepsilon &= \mathbf{p}_1 \int_0^{+\infty} \frac{1}{\sigma_1 \sqrt{2\pi}} \exp\left(-\frac{(x - m_1)^2}{2\sigma_1^2}\right) dx + \\ &+ \mathbf{p}_2 \int_{-\infty}^0 \frac{1}{\sigma_2 \sqrt{2\pi}} \exp\left(-\frac{(x - m_2)^2}{2\sigma_2^2}\right) dx \rightarrow \min_{v_0, v_1, \dots, v_n}.\end{aligned}$$

Здесь  $m_1 = m_1(v_0, v_1, \dots, v_n)$ ,  $\sigma_1 = \sigma_1(v_0, v_1, \dots, v_n)$ ,

$m_2 = m_2(v_0, v_1, \dots, v_n)$ ,  $\sigma_2 = \sigma_2(v_0, v_1, \dots, v_n)$ .

В общем случае эта оптимизационная задача может быть решена численно с использованием специально разработанных итерационных процедур (см., например, [2, с. 105-106]).

**Пример 6.1.** Построить линейный классификатор для двух классов независимых гауссовских наблюдений двумерного случайного

вектора  $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$  с математическими ожиданиями  $\mathbf{a}_1 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$  и

$\mathbf{a}_2 = \begin{pmatrix} 0 \\ 4 \end{pmatrix}$ , корреляционными матрицами  $\mathbf{R}_1 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ ,  $\mathbf{R}_2 = \begin{pmatrix} 4 & 0 \\ 0 & 1 \end{pmatrix}$  и

равными априорными вероятностями классов:  $\mathbf{p}_1 = \mathbf{p}_2 = 1/2$ . Сравнить вероятность ошибки классификации построенного линейного классификатора с вероятностью ошибки байесовской классификации.

### Линейный классификатор.

Линейная разделяющая функция:  $h(x) = v_0 + v_1x_1 + v_2x_2$ . Из соображений симметрии следует, что  $v_1 = 0$ , то есть  $h(x) = v_0 + v_2x_2$ , т.е. уравнение разделяющей поверхности:  $x_2 = -\frac{v_0}{v_2}$ .

Так как параметры  $v_0, v_2$  являются детерминированными (неслучайными) величинами, то случайная величина  $X = v_0 + v_2x_2$  имеет нормальное распределение для объектов  $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$  как первого, так и второго классов. Условные математические ожидания и дисперсии случайной величина  $X$  :

$$m_1 = \mathbf{M}(X | \omega_1) = \mathbf{M}(v_0 + v_2x_2 | \omega_1) = v_0, \quad ,$$

$$m_2 = \mathbf{M}(X | \omega_2) = \mathbf{M}(v_0 + v_2x_2 | \omega_2) = v_0 + 4v_2,$$

$$\sigma_1^2 = \mathbf{D}(X | \omega_1) = \mathbf{D}(v_0 + v_2x_2 | \omega_1) = v_2^2,$$

$$\sigma_2^2 = \mathbf{D}(X | \omega_2) = \mathbf{D}(v_0 + v_2x_2 | \omega_2) = v_2^2.$$

Уравнению разделяющей поверхности  $h(x) = v_0 + v_2x_2 = 0$  соответствует значение случайной величины  $X = 0$ .

Вероятность ошибочной классификации:

$$\begin{aligned} \varepsilon &= \mathbf{p}_1\varepsilon_1 + \mathbf{p}_2\varepsilon_2 = \frac{1}{2} \int_0^{+\infty} f_X(x | \omega_1) dx + \frac{1}{2} \int_{-\infty}^0 f_X(x | \omega_2) dx = \\ &= \frac{1}{2} \int_0^{+\infty} \frac{1}{\sigma_1 \sqrt{2\pi}} \exp\left(-\frac{(x-m_1)^2}{2\sigma_1^2}\right) dx + \frac{1}{2} \int_{-\infty}^0 \frac{1}{\sigma_2 \sqrt{2\pi}} \exp\left(-\frac{(x-m_2)^2}{2\sigma_2^2}\right) dx = \\ &= \frac{1}{2} \int_0^{+\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(x-v_0)^2}{2v_2^2}\right) \frac{dx}{v_2} + \frac{1}{2} \int_{-\infty}^0 \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(x-v_0-4v_2)^2}{2v_2^2}\right) \frac{dx}{v_2} = \\ &= \left\{ u = \frac{x}{v_2}, \quad du = \frac{dx}{v_2}, \quad z = \frac{v_0}{v_2} \right\} = \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{2} \int_0^{+\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(u-z)^2\right) du + \frac{1}{2} \int_{-\infty}^0 \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(u-z-4)^2\right) du = \\
&= \frac{1}{2\sqrt{2\pi}} \int_{-z}^{+\infty} \exp\left(-\frac{u^2}{2}\right) du + \frac{1}{2\sqrt{2\pi}} \int_{-\infty}^{-z-4} \exp\left(-\frac{u^2}{2}\right) du.
\end{aligned}$$

Приравниваем нулю производную  $\varepsilon$  по  $z$  и ищем вероятность ошибки в точке минимума:

$$\frac{\partial \varepsilon}{\partial z} = -\left(-\frac{1}{2\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right)\right) + \left(-\frac{1}{2\sqrt{2\pi}} \exp\left(-\frac{(z+4)^2}{2}\right)\right) = 0 \rightarrow$$

$$\rightarrow \frac{1}{2\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) = \frac{1}{2\sqrt{2\pi}} \exp\left(-\frac{(z+4)^2}{2}\right) \rightarrow$$

$$\rightarrow (z+4)^2 = z^2 \rightarrow z = -2 \rightarrow \frac{v_0}{v_2} = -2;$$

$$\begin{aligned}
\varepsilon &= \frac{1}{2\sqrt{2\pi}} \int_2^{+\infty} \exp\left(-\frac{x^2}{2}\right) dx + \frac{1}{2\sqrt{2\pi}} \int_{-\infty}^{-2} \exp\left(-\frac{x^2}{2}\right) dx = \\
&= \frac{1}{2}(1 - \Phi(2)) + \frac{1}{2}\Phi(-2) = \Phi(-2) = 0.02275013,
\end{aligned}$$

где  $\Phi(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right) du$  – интеграл вероятности (интегральная функция стандартного нормального распределения).

Использованы также соотношения:

$$\frac{1}{\sqrt{2\pi}} \int_x^{+\infty} \exp\left(-\frac{u^2}{2}\right) du = 1 - \Phi(x), \quad \Phi(-x) = 1 - \Phi(x).$$

Вероятность ошибочной классификации:  $\varepsilon = \mathbf{0.02275013}$ .

Уравнение разделяющей поверхности:  $x_2 = 2$ .

Области принятия решений:  $G_1 = \{(x_1, x_2): x_2 < 2\}$ ,

$G_2 = \{(x_1, x_2): x_2 > 2\}$ .

## Оптимальный байесовский классификатор

Теперь для сравнения построим уравнение оптимальной разделяющей поверхности при байесовской классификации и рассчитаем вероятность ошибки.

$$\mathbf{p}_1 f(\mathbf{x}|\omega_1) = \mathbf{p}_2 f(\mathbf{x}|\omega_2), \quad \mathbf{p}_1 = \mathbf{p}_2 = 1/2,$$

$$\frac{1}{2} \frac{1}{\sqrt{|\mathbf{R}_1|}} \exp\left[-\frac{1}{2}(\mathbf{x}-\mathbf{a}_1)^T \mathbf{R}_1^{-1}(\mathbf{x}-\mathbf{a}_1)\right] = \frac{1}{2} \frac{1}{\sqrt{|\mathbf{R}_2|}} \exp\left[-\frac{1}{2}(\mathbf{x}-\mathbf{a}_2)^T \mathbf{R}_2^{-1}(\mathbf{x}-\mathbf{a}_2)\right],$$

$$-\frac{1}{2}(\mathbf{x}-\mathbf{a}_1)^T \mathbf{R}_1^{-1}(\mathbf{x}-\mathbf{a}_1) + \frac{1}{2}(\mathbf{x}-\mathbf{a}_2)^T \mathbf{R}_2^{-1}(\mathbf{x}-\mathbf{a}_2) = \ln \frac{1}{2},$$

$$-\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}^T \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} x_1 \\ x_2 - 4 \end{pmatrix}^T \begin{pmatrix} \frac{1}{4} & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 - 4 \end{pmatrix} = -2 \ln 2,$$

$$-x_1^2 - x_2^2 + \frac{1}{4}x_1^2 + (x_2 - 4)^2 = -2 \ln 2, \quad -\frac{3}{4}x_1^2 - 8x_2 + 16 = -2 \ln 2.$$

Уравнение разделяющей поверхности в виде  $x_2 = x_2(x_1)$ :

$$x_2 = -\frac{3x_1^2}{32} + 2 + \frac{\ln 2}{4} \quad (\text{парабола}).$$

Области принятия решений:

$$G_1 = \left\{ (x_1, x_2) : x_2 < -\frac{3x_1^2}{32} + 2 + \frac{\ln 2}{4} \right\},$$

$$G_2 = \left\{ (x_1, x_2) : x_2 > -\frac{3x_1^2}{32} + 2 + \frac{\ln 2}{4} \right\}.$$

Вероятность ошибочной классификации:

$$\begin{aligned} \varepsilon &= \mathbf{p}_1 \Pr\{\text{ошибка}|\omega_1\} + \mathbf{p}_2 \Pr\{\text{ошибка}|\omega_2\} = \\ &= \frac{1}{2}(\Pr\{X \in G_2|\omega_1\} + \Pr\{X \in G_1|\omega_2\}) = \frac{1}{2}(\varepsilon_1 + \varepsilon_2), \end{aligned}$$

Рассчитаем численно вероятности ошибок первого и второго рода:

$$\begin{aligned} \varepsilon_1 &= \iint_{G_2} f_1(\mathbf{x}) d\mathbf{x} = \iint_{G_2} \frac{1}{2\pi\sqrt{|\mathbf{R}_1|}} \exp\left[-\frac{1}{2}(\mathbf{x}-\mathbf{a}_1)^T \mathbf{R}_1^{-1}(\mathbf{x}-\mathbf{a}_1)\right] d\mathbf{x} = \\ &= \iint_{x_2 > -\frac{3x_1^2}{32} + 2 + \frac{\ln 2}{4}} \frac{1}{2\pi} \exp\left[-\frac{1}{2}(x_1^2 + x_2^2)\right] dx_1 dx_2 = 0.01980355. \end{aligned}$$



$$\begin{aligned} \varepsilon_2 &= \iint_{G_1} f_2(\mathbf{x}) d\mathbf{x} = \iint_{G_1} \frac{1}{2\pi\sqrt{|\mathbf{R}_2|}} \exp\left[-\frac{1}{2}(\mathbf{x}-\mathbf{a}_2)^T \mathbf{R}_2^{-1}(\mathbf{x}-\mathbf{a}_2)\right] d\mathbf{x} = \\ &= \iint_{G_1} \frac{1}{2\pi\sqrt{4}} \exp\left[-\frac{1}{2}\begin{pmatrix} x_1 \\ x_2-4 \end{pmatrix}^T \begin{pmatrix} \frac{1}{4} & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2-4 \end{pmatrix}\right] dx_1 dx_2 = \\ &= \iint_{x_2 < -\frac{3x_1^2}{32} + 2 + \frac{\ln 2}{4}} \frac{1}{2 \cdot 2\pi} \exp\left[-\frac{1}{2}\left(\frac{x_1^2}{4} + (x_2-4)^2\right)\right] dx_1 dx_2 = 0.02026693. \end{aligned}$$

Общая вероятность ошибки  $\varepsilon = \mathbf{0.02003524}$  (приблизительно на 12% меньше, чем у линейного классификатора).

```
# R version 3.4.2 (2017-09-28) + prasma 2017 (c) MacroSoft
# Линейное и байесовское разделение двух гауссовых классов с
равными
# априорными вероятностями и некоррелированными компонентами:
# N1: m=(0;0), D=(1,0; 0,1), Pr1=1/2
# N2: m=(0;4), D=(4,0; 0,1), Pr2=1/2

require(prasma)

e0 <- pnorm(-2) # Ошибка линейного классификатора e0

f1 <- function(x1, x2) # Ошибка первого рода байесовского клас-
сификатора e1
  ifelse(x2 > -3*x1*x1/32 + 2 + log(2)/4, exp(-x1^2/2 -
x2^2/2)/(2*pi), 0)
e1 <- dblquad(f1, -Inf, Inf, -Inf, Inf, dim=1)

f2 <- function(x1, x2) # Ошибка второго рода байесовского клас-
сификатора e2
  ifelse(x2 > -3*x1*x1/32 + 2 + log(2)/4, 0, exp(-x1^2/8 - (x2-
4)^2/2)/(2*pi))/2
e2 <- dblquad(f2, -Inf, Inf, -Inf, Inf, dim=1)

e <- .5*(e1+e2) # Общая ошибка байесовского классификатора e

fprintf ("\nОшибка линейного классификатора e0 = %10.8f\n", e0)
fprintf ("\nОшибка первого рода e1 = %10.8f", e1)
fprintf ("\nОшибка второго рода e2 = %10.8f", e2)
fprintf ("\nОбщая ошибка e = %10.8f", e)
```

Рис. 6.1. Фрагмент R-программы для вычисления вероятностей ошибок

## Линейная и оптимальная разделяющие поверхности

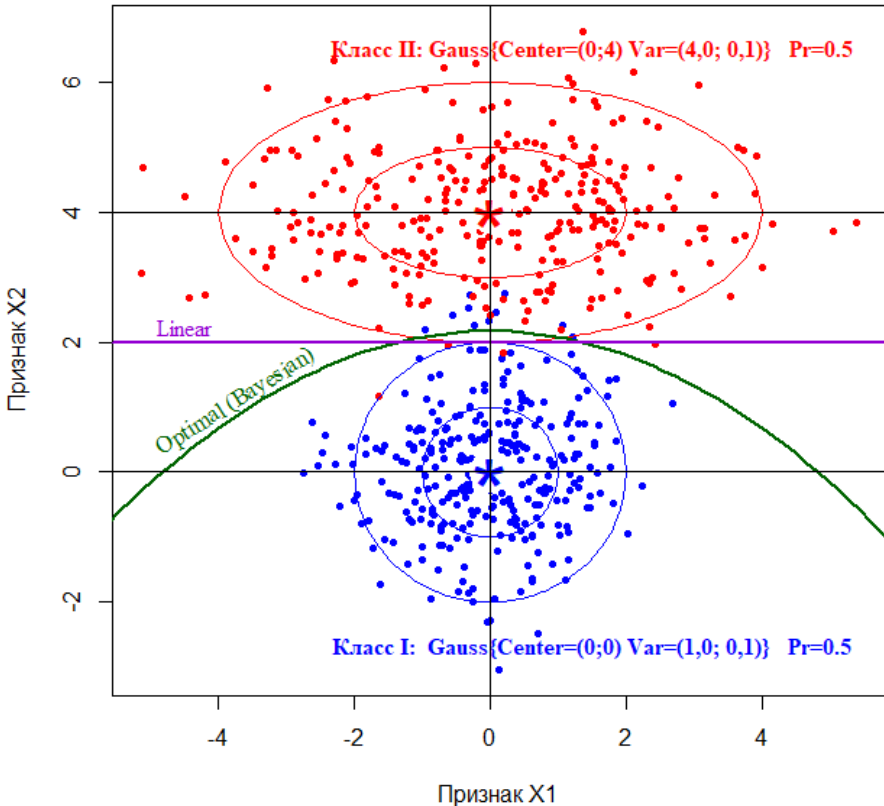


Рис. 6.2. Разделяющие поверхности для оптимальной байесовской и для оптимальной линейной классификации

**Пример 6.2.** Построить линейный классификатор для двух классов независимых гауссовских наблюдений двумерного случайного вектора  $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$  с математическими ожиданиями  $\mathbf{a}_1 = \begin{pmatrix} 2 \\ 0 \end{pmatrix}$  и  $\mathbf{a}_2 = \begin{pmatrix} 0 \\ 2 \end{pmatrix}$ , корреляционными матрицами  $\mathbf{R}_1 = \begin{pmatrix} 0.2 & 0.1 \\ 0.1 & 0.2 \end{pmatrix}$ ,

$\mathbf{R}_2 = \begin{pmatrix} 0.2 & -0.1 \\ -0.1 & 0.2 \end{pmatrix}$  и одинаковыми априорными вероятностями

классов:  $\mathbf{p}_1 = \mathbf{p}_2 = 0.5$ .

Линейная разделяющая функция:  $h(x) = v_0 + v_1x_1 + v_2x_2$ . Так как параметры  $v_0, v_1, v_2$  являются детерминированными (неслучайными) величинами, то случайная величина  $X = v_0 + v_1x_1 + v_2x_2$  имеет

нормальное распределение для объектов  $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$  как первого,

так и второго классов. Условные математические ожидания и дисперсии случайной величина  $X$  :

$$m_1 = \mathbf{M}(X | \omega_1) = \mathbf{M}(v_0 + v_1x_1 + v_2x_2 | \omega_1) =$$

$$= v_0 + v_1\mathbf{M}(x_1 | \omega_1) + v_2\mathbf{M}(x_2 | \omega_1) = v_0 + 2v_1,$$

$$m_2 = \mathbf{M}(X | \omega_2) = \mathbf{M}(v_0 + v_1x_1 + v_2x_2 | \omega_2) =$$

$$= v_0 + v_1\mathbf{M}(x_1 | \omega_2) + v_2\mathbf{M}(x_2 | \omega_2) = v_0 + 2v_2,$$

$$\sigma_1^2 = \mathbf{D}(X | \omega_1) = \mathbf{D}(v_0 + v_1x_1 + v_2x_2 | \omega_1) =$$

$$= v_1^2\mathbf{D}(x_1 | \omega_1) + v_2^2\mathbf{D}(x_2 | \omega_1) + 2v_1v_2\mathbf{R}(x_1, x_2 | \omega_1) = 0.2v_1^2 + 0.2v_2^2 + 0.2v_1v_2 = 0.2(v_1^2 + v_2^2 + v_1v_2),$$

$$\sigma_2^2 = \mathbf{D}(X | \omega_2) = \mathbf{D}(v_0 + v_1x_1 + v_2x_2 | \omega_2) =$$

$$= v_1^2\mathbf{D}(x_1 | \omega_2) + v_2^2\mathbf{D}(x_2 | \omega_2) + 2v_1v_2\mathbf{R}(x_1, x_2 | \omega_2) =$$

$$= 0.2v_1^2 + 0.2v_2^2 - 0.2v_1v_2 = 0.2(v_1^2 + v_2^2 - v_1v_2).$$

Область принятия решений в пользу первого класса

$$G_1 = \{(x_1, x_2): v_0 + v_1x_1 + v_2x_2 < 0\}.$$

Область принятия решений в пользу второго класса

$$G_2 = \{(x_1, x_2): v_0 + v_1x_1 + v_2x_2 > 0\}.$$

Вероятность ошибочной классификации:

$$\begin{aligned}
\boldsymbol{\varepsilon} &= \mathbf{p}_1 \boldsymbol{\varepsilon}_1 + \mathbf{p}_2 \boldsymbol{\varepsilon}_2 = \frac{1}{2} \int_0^{+\infty} f_X(x | \boldsymbol{\omega}_1) dx + \frac{1}{2} \int_{-\infty}^0 f_X(x | \boldsymbol{\omega}_2) dx = \\
&= \frac{1}{2} \int_0^{+\infty} \frac{1}{\sigma_1 \sqrt{2\pi}} \exp\left(-\frac{(x-m_1)^2}{2\sigma_1^2}\right) dx + \frac{1}{2} \int_{-\infty}^0 \frac{1}{\sigma_2 \sqrt{2\pi}} \exp\left(-\frac{(x-m_2)^2}{2\sigma_2^2}\right) dx = \\
&= \frac{1}{2} \int_0^{+\infty} \frac{1}{\sqrt{2\pi \cdot 0.2(v_1^2 + v_2^2 + v_1 v_2)}} \exp\left(-\frac{(x-v_0-2v_1)^2}{2 \cdot 0.2(v_1^2 + v_2^2 + v_1 v_2)}\right) dx + \\
&\left\{ z = \frac{x-v_0-2v_1}{\sqrt{0.2(v_1^2 + v_2^2 + v_1 v_2)}}; dz = \frac{dx}{\sqrt{0.2(v_1^2 + v_2^2 + v_1 v_2)}} \left[ \begin{array}{l} x = +\infty \Rightarrow z = +\infty \\ x = 0 \Rightarrow z = \frac{-v_0-2v_1}{\sqrt{0.2(v_1^2 + v_2^2 + v_1 v_2)}} \end{array} \right] \right\} \\
&+ \frac{1}{2} \int_{-\infty}^0 \frac{1}{\sqrt{2\pi \cdot 0.2(v_1^2 + v_2^2 - v_1 v_2)}} \exp\left(-\frac{(x-v_0-2v_2)^2}{2 \cdot 0.2(v_1^2 + v_2^2 - v_1 v_2)}\right) dx = \\
&\left\{ z = \frac{x-v_0-2v_2}{\sqrt{0.2(v_1^2 + v_2^2 - v_1 v_2)}}; dz = \frac{dx}{\sqrt{0.2(v_1^2 + v_2^2 - v_1 v_2)}} \left[ \begin{array}{l} x = -\infty \Rightarrow z = -\infty \\ x = 0 \Rightarrow z = \frac{-v_0-2v_2}{\sqrt{0.2(v_1^2 + v_2^2 - v_1 v_2)}} \end{array} \right] \right\} \\
&= \frac{1}{2} \int_{\frac{-v_0-2v_1}{\sqrt{v_1^2 + v_2^2 + v_1 v_2}}}^{+\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) dx + \frac{1}{2} \int_{-\infty}^{\frac{-v_0-2v_2}{\sqrt{v_1^2 + v_2^2 - v_1 v_2}}} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) dx = \\
&= \frac{1}{2} - \frac{1}{2} \Phi\left(\frac{-v_0-2v_1}{\sqrt{0.2(v_1^2 + v_2^2 + v_1 v_2)}}\right) + \frac{1}{2} \Phi\left(\frac{-v_0-2v_2}{\sqrt{0.2(v_1^2 + v_2^2 - v_1 v_2)}}\right) \rightarrow \min_{v_0, v_1, v_2},
\end{aligned}$$

где  $\Phi(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right) du$  – интеграл вероятности (инте-

гральная функция стандартного нормального распределения).

Численная минимизация даёт результат:  $v_0 = 0.453$ ,

$v_1 = -0.927$ ,  $v_2 = 0.927$ . Разделяющая поверхность

$0.453 - 0.927x_1 + 0.927x_2 = 0$ , или  $x_2 = x_1 - 0,489$ . Мини-

мальная вероятность ошибки равна  $0.0005101382$ .

```

# Filename: LinClass2.r R version 3.4.0 (2017-04-21) MacroSoft
(c) 2017
# Вероятность ошибки линейного классификатора
# p1 = 1/2; p2 = 1/2; M1 = (2 0); M2 = (0 2);
# R1 = [.2, .1; .1, .2]; R2 = [.2, -.1; -.1, .2]

require(mvtnorm)

fr <- function(x) {
  v1 <- x[1]; v2 <- x[2]; v0 <- x[3]
  (1 - pnorm((-v0-2*v1)/sqrt(.2*(v1*v1+v2*v2+v1*v2)))
  + pnorm((-v0-2*v2)/sqrt(.2*(v1*v1+v2*v2-v1*v2))))/2
}

opt_res <- optim(c(-1, 1, 0), fr) # Filename: LinClass2.r R ver-
sion 3.4.0 (2017-04-21) MacroSoft (c) 2017
# Вероятность ошибки линейного классификатора
# p1 = 1/2; p2 = 1/2; M1 = (2 0); M2 = (0 2);
# R1 = [.2, .1; .1, .2]; R2 = [.2, -.1; -.1, .2]

require(mvtnorm)

fr <- function(x) {
  v1 <- x[1]; v2 <- x[2]; v0 <- x[3]
  (1 - pnorm((-v0-2*v1)/sqrt(.2*(v1*v1+v2*v2+v1*v2)))
  + pnorm((-v0-2*v2)/sqrt(.2*(v1*v1+v2*v2-v1*v2))))/2
}

opt_res <- optim(c(-1, 1, 0), fr)
print(opt_res$par); print(opt_res$value)
windows(600, 600); n <- 200
M1 <- c(2, 0); R1 <- .2*array(c(1, .5, .5, 1), dim=c(2,2))
M2 <- c(0, 2); R2 <- .2*array(c(1, -.5, -.5, 1), dim=c(2,2))
xy <- rbind(rmvnorm(n,M1,R1), rmvnorm(n,M2,R2))
plot(xy, xlab='X1', ylab='X2')
abline(h=0); abline(v=0); abline(a=-.489, b=1)
text(2, 2.4, "x2 = x1 - 0.489")
print(opt_res$par); print(opt_res$value)
windows(600, 600); n <- 200
M1 <- c(2, 0); R1 <- .2*array(c(1, .5, .5, 1), dim=c(2,2))
M2 <- c(0, 2); R2 <- .2*array(c(1, -.5, -.5, 1), dim=c(2,2))
xy <- rbind(rmvnorm(n,M1,R1), rmvnorm(n,M2,R2))
plot(xy, xlab='X1', ylab='X2')
abline(h=0); abline(v=0); abline(a=-.489, b=1)
text(2, 2.4, "x2 = x1 - 0.489")

```

Рис. 6.3. Фрагмент R-программы для расчёта  
и визуализации минимальной вероятности ошибки

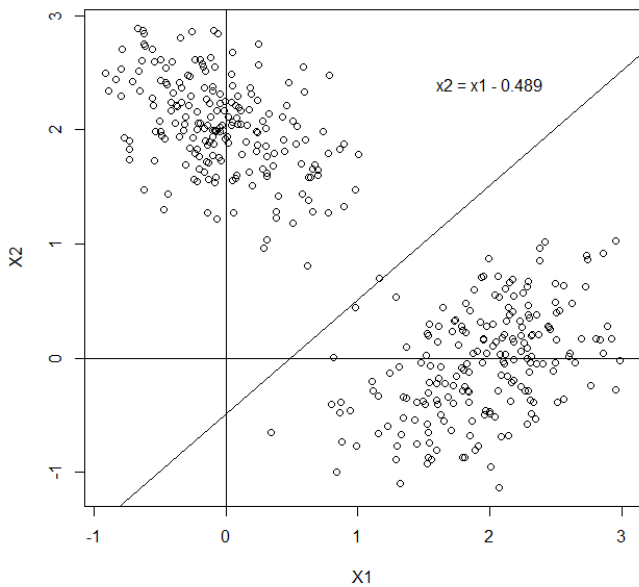


Рис. 6.4. Разделяющая поверхность для оптимальной линейной классификации

## 7 ДИСКРИМИНАНТНЫЙ АНАЛИЗ

### 7.1 Матрицы рассеяния и критерий разделимости

*Дискриминант* (лат. *discriminans* (*dis-criminantis*)) переводится с латыни как различающий, разделяющий). Целью дискриминантного анализа является расчёт некоторого коэффициента, влияющего на эффективность классификации (вероятности ошибок первого и второго рода) при использовании того или иного набора признаков для классификации. Этот коэффициент не должен зависеть от вероятностного распределения признаков, от используемых методов классификации и от соответствующей вероятности ошибки классификации. Рассмотрим два класса объектов с  $n$ -мерными признаками  $X_1, X_2$  и априорными вероятностями классов  $p_1 = \Pr\{\omega_1\}$ ,  $p_2 = \Pr\{\omega_2\}$ .

В дискриминантном анализе критерии разделимости классов формулируются с использованием матриц рассеяния внутри классов и матриц рассеяния между классами [2].

*Матрица рассеяния между классами*  $S_1$  показывает разброс объектов между классами. Эта матрица определяется несколькими способами. Мы используем самый простой:

$S_1 = (M_1 - M_2)(M_1 - M_2)^T$  – расстояние между центрами классов.

Другие методы определения межклассового разброса признаков:

$$S'_1 = \sum_{i=1}^2 p_i \left[ (M_i - M_0)(M_i - M_0)^T \right] - \text{среднее расстояние до}$$

общего центра, где  $M_0$  – математическое ожидание смеси:

$$M_0 = p_1 M_1 + p_2 M_2;$$

$$S''_1 = \mathbf{M} \left( X^{(1)} - X^{(2)} \right) \left( X^{(1)} - X^{(2)} \right)^T - \text{математическое ожида-}$$

ние расстояния между объектами первого и второго классов.

*Матрица рассеяния внутри классов*  $S_2$  показывает разброс признаков объектов относительно векторов математических ожиданий классов:

$$S_2 = p_1 R_1 + p_2 R_2 = p_1 \mathbf{M} \left[ (X_1 - M_1)(X_1 - M_1)^T \right] + \\ + p_2 \mathbf{M} \left[ (X_2 - M_2)(X_2 - M_2)^T \right], \text{ где } M_1 = \mathbf{M}X_1, M_2 = \mathbf{M}X_2 .$$

*Критерием разделимости* классов должна быть функция от матриц рассеяния, определённых выше. Эта функция должна увеличиваться при увеличении рассеяния между классами  $S_1$  и при уменьшении рассеяния внутри классов  $S_2$ . Таким условиям удовлетворяет множество критериев, из которых наиболее простым является

$$J_1 = \text{tr} \left( S_2^{-1} S_1 \right).$$

Другими наиболее употребительными критериями разделимости являются следующие:

$$J_2 = \ln \left( |S_2^{-1} S_1| \right) = \ln \left( |S_1| / |S_2| \right),$$

$$J_3 = \text{tr} S_1 - \text{tr} S_2,$$

$$J_4 = \text{tr} S_1 / \text{tr} S_2 .$$



## 7.2 Разделимость выборки

Допустим, что имеется выборка признаков (обучающая выборка) двух классов объектов (априорные вероятности классов  $p_1$  и  $p_2$ , число объектов каждого класса  $N_1$  и  $N_2$ )

$$X_1 = \begin{pmatrix} x_1^1 \\ x_1^2 \\ x_1^3 \\ \dots \\ x_1^{N_1} \end{pmatrix}, \quad X_2 = \begin{pmatrix} x_2^1 \\ x_2^2 \\ x_2^3 \\ \dots \\ x_2^{N_2} \end{pmatrix},$$

где  $x_i^j = (x_{i1}^j \ x_{i2}^j \ \dots \ x_{in}^j)^T$  –  $n$ -мерный вектор признаков объектов  $i$ -го класса в  $j$ -м наблюдении. Требуется оценить критерий разделимости  $\mathbf{J}_1 = \text{tr}(S_2^{-1}S_1)$ .

Векторы центров классов:

$$M_1 = (M_{11} \ M_{12} \ \dots \ M_{1n})^T, \quad M_{1k} = \frac{1}{N_1} \sum_{j=1}^{N_1} x_{1k}^j,$$

$$M_2 = (M_{21} \ M_{22} \ \dots \ M_{2n})^T, \quad M_{2k} = \frac{1}{N_2} \sum_{j=1}^{N_2} x_{2k}^j.$$

*Матрица рассеяния между классами:*

$$S_1 = (M_1 - M_2)(M_1 - M_2)^T.$$

Ковариационные матрицы признаков первого и второго классов:

$$R_1 = \left\{ \frac{1}{N_1} \sum_{j=1}^{N_1} \tilde{x}_{1k}^j \tilde{x}_{1l}^j \right\}_{n \times n}, \quad R_2 = \left\{ \frac{1}{N_2} \sum_{j=1}^{N_2} \tilde{x}_{2k}^j \tilde{x}_{2l}^j \right\}_{n \times n},$$

где  $\tilde{x}_{ik}^j$  – центрированные значения признаков:

$$\tilde{x}_{ik}^j = x_{ik}^j - M(x_{ik}).$$

*Матрица рассеяния внутри классов:*

$$S_2 = p_1 R_1 + p_2 R_2 = \frac{p_1}{N_1} \sum_{j=1}^{N_1} \tilde{x}_{1k}^j \tilde{x}_{1l}^j + \frac{p_2}{N_2} \sum_{j=1}^{N_2} \tilde{x}_{2k}^j \tilde{x}_{2l}^j.$$

### 7.3 Выбор признаков, максимизирующих критерий $J_1$

#### 7.3.1 Линейное преобразование случайных векторов

Пусть случайный вектор-столбец  $X = (x_1 \ x_2 \ \dots \ x_n)^T$  имеет математическое ожидание  $\mathbf{M}X$  и ковариационную матрицу  $R_X$ . Тогда случайный вектор  $Y = (y_1 \ y_2 \ \dots \ y_m)^T$ , являющийся линейным преобразованием вектора  $X$ ,

$$Y = AX$$

имеет математическое ожидание

$$\mathbf{M}Y = \mathbf{A}MX$$

и ковариационную матрицу

$$R_Y = AR_X A^T,$$

так как

$$\begin{aligned} R_Y &= \mathbf{M}[(Y - \mathbf{M}Y)(Y - \mathbf{M}Y)^T] = \mathbf{M}[(AX - \mathbf{A}MX)(AX - \mathbf{A}MX)^T] = \\ &= \mathbf{M}[A(X - \mathbf{M}X)(X - \mathbf{M}X)^T A^T] = \mathbf{A} \mathbf{M}[(X - \mathbf{M}X)(X - \mathbf{M}X)^T] A^T = AR_X A^T. \end{aligned}$$

### 7.3.2 Диагоналирующее и декоррелирующее преобразования

Обозначим через  $\lambda_1, \dots, \lambda_n$  собственные значения, а через  $U_1, U_2, \dots, U_n$  – ортонормированные собственные векторы некоторой ковариационной матрицы  $R$  размера  $n \times n$ .

$$\lambda_i U_i = R U_i, \quad i = \overline{1, n}.$$

Так как ковариационная матрица  $R$  является симметричной и положительно определённой, то все её  $n$  собственных значений  $\lambda_1, \dots, \lambda_n$  являются вещественными и положительными. Справедливо матричное равенство:

$$U \Lambda = R U,$$

где  $U = [U_1 \ U_2 \ \dots \ U_n]$  – ортонормированная матрица собственных векторов матрицы  $R$ , расположенных по столбцам;  $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$  – диагональная матрица собственных значений матрицы  $R$ . Из матричной алгебры известно, что  $U^T U = I$ ,  $U U^T = I$ ,  $U^{-1} = U^T$ ,  $R = U \Lambda U^T$ ,  $U^T R U = \Lambda$ ,  $\Lambda^{-1/2} U^T R U \Lambda^{-1/2} = I$ .

Если в качестве матрицы  $A$  линейного преобразования случайного вектора  $X$  в случайный вектор  $Y$

$$Y = A X$$

выбрать транспонированную матрицу собственных векторов

$$A = U^T,$$

то ковариационная матрица результирующего вектора  $R_Y = A R_X A^T = U^T R_X U = \Lambda$ , то есть компоненты вектора  $Y$  бу-

дуг некоррелированными случайными величинами с дисперсиями  $\lambda_1, \dots, \lambda_n$  – собственными числами матрицы  $R_X$ . Такое преобразование  $A$  в этом разделе будем называть *диагонализующим преобразованием*.

Если также потребовать, чтобы компоненты вектора  $Y$  обладали единичной дисперсией, в качестве матрицы линейного преобразования следует выбрать

$$A = \Lambda^{-1/2} U^T.$$

Тогда  $R_Y = A R_X A^T = \Lambda^{-1/2} U^T R_X U \Lambda^{-1/2} = \Lambda^{-1/2} \Lambda \Lambda^{-1/2} = I$ . Преобразование  $A = \Lambda^{-1/2} U^T$  будем называть *декоррелирующим преобразованием*.

Декоррелированность декоррелированного случайного вектора сохраняется при *любом* ортонормированном преобразовании  $A = \Psi^T$ :  $Z = \Psi^T Y$ . То есть, если  $R_Y = I$ , то

$$R_Z = \Psi^T R_Y \Psi = \Psi^T I \Psi = \Psi^T \Psi = I.$$

Фактически, оба преобразования  $A = U^T$  и  $A = \Lambda^{-1/2} U^T$  являются одновременно и диагонализующими, и декоррелирующими.

### 7.3.3 Одновременная диагонализация двух ковариационных матриц

Пусть  $R_{1X}$  и  $R_{2X}$  – две в общем случае различные ковариационные матрицы случайных векторов  $X_1$  и  $X_2$ . Сначала применим декоррелирующее преобразование  $Y_1 = \Lambda_{1X}^{-1/2} U_{1X}^T X_1$  к вектору  $X_1$ . Тогда

$$R_{Y_1} = \Lambda_{1X}^{-1/2} U_{1X}^T R_{1X} U_{1X} \Lambda_{1X}^{-1/2} = I.$$

То же самое преобразование (согласованное с  $R_{X_1}$ )  $Y_2 = \Lambda_{X_1}^{-1/2} U_{X_1}^T X_2$  применим к вектору  $X_2$ . Получим:

$$R_{Y_2} = \Lambda_{X_1}^{-1/2} U_{X_1}^T R_{X_2} U_{X_1} \Lambda_{X_1}^{-1/2}.$$

В общем случае матрица  $R_{Y_2}$  не будет диагональной.

Теперь применим ортонормированное преобразование  $Z_2 = U_{Y_2}^T Y_2$  для диагонализации матрицы  $R_{Y_2}$  ( $U_{Y_2}^T$  – ортонормированная матрица собственных векторов ковариационной матрицы  $R_{2Y}$  случайного вектора  $Y_2$ ):

$$R_{Z_2} = U_{Y_2}^T R_{Y_2} U_{Y_2} = \Lambda_{Y_2}.$$

И, наконец, применим последнее ортонормированное преобразование  $Z_1 = U_{X_2}^T Y_1$  к декоррелированному вектору  $Y_1$  (при этом её декоррелированность сохраняется, что доказано выше):

$$R_{1Z} = U_{2X}^T R_{1Y} U_{2X} = U_{2X}^T I U_{2X} = U_{2X}^T U_{2X} = I.$$

Таким образом, искомая матрица преобразования, одновременно декоррелирующая  $X_1$  и  $X_2$ , равна:

$$A = U_{Y_2}^T \Lambda_{X_1}^{-1/2} U_{X_1}^T : Z_1 = (U_{Y_2}^T \Lambda_{X_1}^{-1/2} U_{X_1}^T) X_1, Z_2 = (U_{Y_2}^T \Lambda_{X_1}^{-1/2} U_{X_1}^T) X_2.$$

$$\begin{aligned} R_{Z_1} &= U_{Y_2}^T \Lambda_{X_1}^{-1/2} (U_{X_1}^T R_{X_1} U_{X_1}) \Lambda_{X_1}^{-1/2} U_{Y_2} = U_{Y_2}^T \Lambda_{X_1}^{-1/2} (\Lambda_{X_1}) \Lambda_{X_1}^{-1/2} U_{Y_2} = \\ &= U_{Y_2}^T U_{Y_2} = I. \end{aligned}$$

$$R_{Z_2} = U_{Y_2}^T (\Lambda_{X_1}^{-1/2} U_{X_1}^T R_{X_2} U_{X_1} \Lambda_{X_1}^{-1/2}) U_{Y_2} = U_{Y_2}^T (R_{Y_2}) U_{Y_2} = \Lambda_{Y_2}$$

Полученный алгоритм поясняет следующая схема (рис. 7.1).

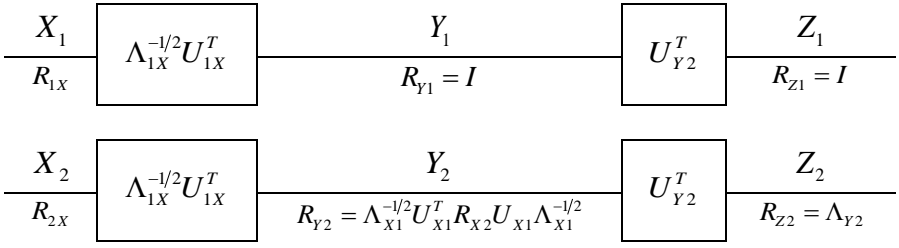


Рис. 7.1. Схема одновременной диагонализации двух ковариационных матриц

Докажем, что полученная матрица преобразования  $A$  является транспонированной матрицей собственных векторов матрицы  $R_{X1}^{-1} R_{X2}$ .

Собственные числа  $\lambda$  матрицы  $R_{Y2}$  определяются из уравнения  $|R_{Y2} - \lambda I| = 0$ . Подставляя формулы декоррелирующего преобразования, приведённые выше, и учитывая, что матрица декоррелирующего преобразования  $\Lambda_{X1}^{-1/2} U_{X1}^T$  и ковариационная матрица  $R_{X1}$  не вырождены, получаем:

$$\begin{aligned} & \left| \Lambda_{X1}^{-1/2} U_{X1}^T R_{X2} U_{X1} \Lambda_{X1}^{-1/2} - \lambda \Lambda_{X1}^{-1/2} U_{X1}^T R_{X1} U_{X1} \Lambda_{X1}^{-1/2} \right| = 0, \\ & \left| \left( \Lambda_{X1}^{-1/2} U_{X1}^T \right) \left( R_{X2} - \lambda R_{X1} \right) \left( U_{X1} \Lambda_{X1}^{-1/2} \right) \right| = 0, \quad |R_{X2} - \lambda R_{X1}| = 0, \quad |R_{X1}^{-1} R_{X2} - \lambda I| = 0. \end{aligned}$$

Это означает, что собственные числа  $\lambda_i$  матрицы  $R_{X1}^{-1} R_{X2}$  совпадают с собственными числами матрицы  $R_{Y2}$ . Подставив в уравнение  $U_{Y2}^T R_{Y2} U_{Y2} = \Lambda$  значение  $R_{Y2} = \Lambda_{X1}^{-1/2} U_{X1}^T R_{X2} U_{X1} \Lambda_{X1}^{-1/2}$ , получаем:

$$\begin{aligned} & U_{Y2}^T \Lambda_{X1}^{-1/2} U_{X1}^T R_{X2} U_{X1} \Lambda_{X1}^{-1/2} U_{Y2} = \Lambda_{Y2}, \\ & R_{2X} U_{1X} \Lambda_{1X}^{-1/2} U_{2Y} = \left( U_{2Y}^T \Lambda_{1X}^{-1/2} U_{1X}^T \right)^{-1} \Lambda_{Y2}, \\ & R_{1X}^{-1} R_{2X} U_{1X} \Lambda_{1X}^{-1/2} U_{2Y} = U_{1X} \Lambda_{1X}^{-1} U_{1X}^T \left( U_{2Y}^T \Lambda_{1X}^{-1/2} U_{1X}^T \right)^{-1} \Lambda_{Y2}, \\ & R_{1X}^{-1} R_{2X} \left( U_{1X} \Lambda_{1X}^{-1/2} U_{2Y} \right) = \left( U_{1X} \Lambda_{1X}^{-1/2} U_{2Y} \right) \Lambda_{Y2}. \end{aligned}$$

Таким образом, матрица диагонализующего преобразования  $A$  равна транспонированной матрице собственных векторов  $(U_{1X}\Lambda_{1X}^{-1/2}U_{2Y})^T$  матрицы  $R_{1X}^{-1}R_{2X}$ .

### 7.3.4 Оптимизация сокращения размерности пространства признаков

Критерий  $J_1 = \text{tr}(S_{X_2}^{-1}S_{X_1})$  может быть максимизирован с использованием линейного преобразования пространства признаков:

$$Y = AX,$$

где матрица преобразования  $A$  имеет размер  $m \times n$  ( $m \leq n$ ),  $n$  – размерность исходного пространства признаков  $X$ ,  $m$  – новая размерность пространства признаков.

Идея заключается в нахождении такой матрицы преобразования  $A$ , которая позволяет одновременно диагонализировать две матрицы рассеяния  $S_{X_1}$  и  $S_{X_2}$ . При любом выборе матриц рассеяния  $S_{X_1}$  и  $S_{X_2}$  в пространстве  $X$  соответствующие матрицы рассеяния в пространстве  $Y$  в силу линейности преобразования равны

$$S_{1Y} = AS_{1X}A^T, S_{2Y} = AS_{2X}A^T.$$

Пусть  $\lambda_i, U_i, i = 1, 2, 3, \dots, n$  – собственные значения и собственные векторы матрицы  $S_{2X}^{-1}S_{1X}$ ,  $\mu_j, V_j, j = 1, 2, 3, \dots, m$  – собственные значения и собственные векторы матрицы  $S_{2Y}^{-1}S_{1Y}$ .

За счёт выбора матрицы преобразования  $A$  можно провести одновременную диагонализацию матриц  $S_{1X}$  и  $S_{2X}$ .

$$AS_{1X}A^T = \Lambda, AS_{2X}A^T = I_n.$$

При этом  $J_{1X} = \text{tr } S_{2X}^{-1} S_{1X} = \text{tr } U \Lambda U^T = \text{tr } U^T U \Lambda = \text{tr } \Lambda = \sum_{i=1}^n \lambda_i$ .

Аналогично:  $J_{1Y} = \text{tr } S_{2Y}^{-1} S_{1Y} = \sum_{i=1}^m \mu_i$ .

Здесь использовано свойство:  $\text{tr } AB = \text{tr } BA$  для любых матриц согласованного размера. Можно доказать [2, п.9.2.2], что критерий  $J_{1Y}$  при  $m < n$  достигает максимального значения при выборе линейного преобразования  $A = [U_1 \ U_2 \ \dots \ U_m]^T$ , где  $U_1, U_2, \dots, U_m$  – собственные векторы, соответствующие  $m$  максимальным собственным значениям матрицы  $S_{2X}^{-1} S_{1X}$ , при этом  $\mu_i = \lambda_i$ ,  $i = 1, 2, 3, \dots, m$ . Собственные значения считаются упорядоченными по убыванию:  $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_n$ .



## 8 МЕТОД ОПОРНЫХ ВЕКТОРОВ (SVM)

Под методом опорных векторов (в англоязычной литературе *SVM – Support Vectors Machine – машина опорных векторов*) понимается набор похожих алгоритмов обучения с учителем, разработанных различными авторами в конце прошлого века [4] и использующихся для задач классификации и регрессионного анализа. Мы сосредоточимся на задаче классификации. Обозначения и формулировки в этом разделе приводятся в соответствии с лекциями Воронцова [3].

*Формальная постановка задачи.* Пусть  $X$  – пространство объектов (признаков);  $Y$  – множество классов;  $X \rightarrow Y$  – функциональная зависимость, известная только на объектах обучающей выборки:  $X^L = (x_i, y_i)_{i=1}^L$ ,  $y_i = y(x_i)$ ;  $L$  – количество объектов в обучающей выборке. Требуется построить алгоритм  $a: X \rightarrow Y$ , аппроксимирующий целевую функциональную зависимость на всём пространстве  $X$ .

### 8.1 Линейно разделяемая выборка

Рассмотрим классификацию на два класса объектов, описываемых  $n$ -мерными векторами:  $X = \mathbb{R}^n \rightarrow Y = \{-1, +1\}$ . Уравнение разделяющей гиперплоскости:

$$\sum_{j=1}^n w_j x^j = b, \quad (8.1)$$

где  $x = (x^1, x^2, \dots, x^n)$  – признаки объекта  $x$ ;  
 $w = (w^1, w^2, \dots, w^n) \in \mathbb{R}^n$ ,  $b \in \mathbb{R}$  – параметры уравнения гиперплоскости.

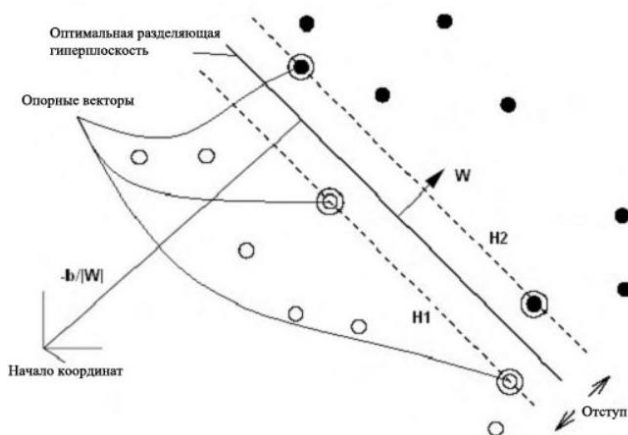


Рис. 8.1. Разделяющая гиперплоскость

В векторном виде уравнение разделяющей гиперплоскости:  $(w, x) = b$ , где через  $(\cdot, \cdot)$  здесь и далее обозначается скалярное произведение векторов. Алгоритм классификации:

$$a(x) = \text{sgn}((w, x) - b) = \begin{cases} +1, & x \text{ находится с одной стороны} \\ & \text{от гиперплоскости,} \\ -1, & x \text{ находится с другой стороны} \\ & \text{от гиперплоскости.} \end{cases}$$

От значений коэффициентов уравнения гиперплоскости  $w$ ,  $b$  зависит общее число ошибок классификации в обучающей выборке:

$$Q(w, b) = \sum_{i=1}^L \left[ y_i \cdot ((w, x_i) - b) < 0 \right],$$

где обозначено:  $[условие] = \begin{cases} 1, & \text{если условие истинно,} \\ 0, & \text{если условие ложно.} \end{cases}$

Если  $Q(w, b) = 0$ , то ошибки отсутствуют при некотором положении разделяющей гиперплоскости. Это положение не единственное (см. рис. 8.1). Поэтому для более уверенной классификации потребуем максимизации зазора (margin, отступа) между классами.

Заметим, что алгоритм классификации  $a(x) = \text{sgn}((w, x_i) - b)$  не изменится, если  $w$  и  $b$  умножить на одну и ту же положительную константу. Удобно выбрать эту константу так, чтобы на всех *граничных* объектах  $x_i$  обучающей выборки выполнялись условия

$$(w, x_i) - b = y_i.$$

Это можно сделать, так как при оптимальном положении разделяющей гиперплоскости все граничные объекты находятся от неё на одинаковом расстоянии. Остальные объекты находятся дальше. Таким образом, для всех  $x_i \in X^L$ :

$$\begin{cases} (w, x_i) - b \leq -1, & \text{если } y_i = -1; \\ (w, x_i) - b \geq +1, & \text{если } y_i = +1, \end{cases} \quad \text{или } y_i \cdot ((w, x_i) - b) = 1, \quad i = \overline{1, L}. \quad (8.2)$$

Условие  $-1 < (w, x) - b < +1$  задаёт полосу, разделяющую классы. Ни одна из точек обучающей выборки не может лежать внутри этой полосы. Границами полосы служат две параллельные гиперплоскости с направляющим вектором  $w$ . Точки, ближайшие к разделяющей гиперплоскости, лежат в точности на границах полосы. При этом сама разделяющая гиперплоскость проходит ровно по середине полосы. *Ширина разделяющей по-*

лосы должна быть максимальной. Пусть  $x_+$  и  $x_-$  – две произвольные точки классов  $+1$  и  $-1$ , лежащие на границе полосы. Тогда ширина полосы есть

$$\left( (x_+ - x_-), \frac{w}{\|w\|} \right) = \frac{(w, x_+) - (w, x_-)}{\|w\|} = \frac{(b+1) - (b-1)}{\|w\|} = \frac{2}{\|w\|} \rightarrow \min_w.$$

Ширина полосы максимальна, когда норма вектора  $w$  минимальна. Требуется найти такие значения параметров  $w$  и  $b$ , при которых норма вектора  $w$  минимальна при условии (8.1). Это задача квадратичного программирования, которая заключается в минимизации квадратичной формы при линейных ограничениях-неравенствах вида (8.1) относительно  $(n+1)$  переменных  $w^1, w^2, \dots, w^n, b$ :

$$\begin{cases} (w, w) \rightarrow \min_{w, b}; \\ y_i \cdot ((w, x_i) - b) \geq 1, \quad i = \overline{1, L}. \end{cases} \quad (8.3)$$

В скалярном виде:

$$\begin{cases} w_1^2 + w_2^2 + \dots + w_n^2 \rightarrow \min_{w_1, \dots, w_n, b}; \\ y_i \cdot (w_1 x_i^1 + w_2 x_i^2 + \dots + w_n x_i^n - b) \geq 1, \quad i = \overline{1, L}. \end{cases}$$

По условию Куна–Таккера в теории нелинейного программирования эта задача эквивалентна задаче поиска седловой точки функции Лагранжа:

$$\begin{cases} \mathbf{L}(w, b; \lambda) = \frac{1}{2}(w, w) - \sum_{i=1}^L \lambda_i (y_i \cdot ((w, x_i) - b) - 1) \rightarrow \min_{w, b} \max_{\lambda}; \\ \lambda_i \geq 0, \quad i = 1, 2, \dots, L; \\ \lambda_i = 0 \quad \text{или} \quad (w, x_i) - b = y_i, \quad i = 1, 2, \dots, L, \end{cases}$$

где  $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_n)$  – вектор двойственных переменных.

Условие  $\lambda_i = 0$  или  $(w, x_i) - b = y_i$ ,  $i = 1, 2, \dots, L$ , называется *условием дополняющей нежесткости*.

Необходимым условием седловой точки является равенство нулю производных Лагранжиана. Из уравнений  $\frac{\partial \mathbf{L}(w, b; \lambda)}{\partial w} = 0$  и  $\frac{\partial \mathbf{L}(w, b; \lambda)}{\partial b} = 0$  следует:

$$w = \sum_{i=1}^L \lambda_i y_i x_i, \quad (8.4)$$

$$\sum_{i=1}^L \lambda_i y_i = 0. \quad (8.5)$$

Из (8.4) следует, что искомый вектор весов  $w$  является линейной комбинацией векторов обучающей выборки, причём только тех, для которых  $\lambda_i \neq 0$ . Согласно условию дополняющей нежесткости на этих векторах  $x_i$  ограничения-неравенства обращаются в равенства:  $\langle w, x_i \rangle - b = y_i$ , следовательно эти векторы находятся на границе разделяющей полосы. Все остальные векторы отстоят дальше от границы, для них  $\lambda_i = 0$  и они не участвуют в суммировании (8.4).

**Определение.** Если  $\lambda_i > 0$  и  $(w, x_i) - b = y_i$ , то объект обучающей выборки  $x_i$  называется опорным вектором (*support vector*).

Подставляя (8.4) и (8.5) обратно в Лагранжиан, получим эквивалентную задачу квадратичного программирования, содержащую только двойственные переменные:

$$\left\{ \begin{array}{l} -\mathbf{L}(\lambda) = \frac{1}{2} \sum_{i=1}^L \sum_j^L \lambda_i \lambda_j y_i y_j (x_i, x_j) - \sum_{i=1}^L \lambda_i \rightarrow \min_{\lambda}; \\ \lambda_i \geq 0, \quad i = 1, 2, \dots, L; \\ \sum_{i=1}^L \lambda_i y_i = 0. \end{array} \right. \quad (8.6)$$

Полученная задача (8.6) имеет единственное решение из-за выпуклости функционала и выпуклости ограничений. Вектор  $w$  вычисляется из (8.4). Для любого опорного вектора  $x_i$  класса  $y_i$  порог  $b = (w, x_i) - y_i$ . На практике для повышения помехоустойчивости следует брать в качестве  $b$  среднее по всем опорным векторам, а ещё лучше медиану:

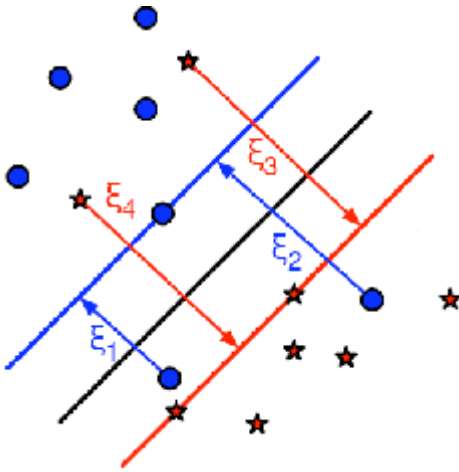
$$b = \text{med} \left\{ (w, x_i) - y_i : \lambda_i > 0, \quad i = \overline{1, L} \right\}. \quad (8.7)$$

В итоге искомый алгоритм классификации может быть записан в следующем виде:

$$a(x) = \text{sgn} \left( \sum_{i=1}^L \lambda_i y_i \langle x_i, x \rangle - b \right). \quad (8.8)$$

## 8.2 Линейно неразделимая выборка

Чтобы обобщить SVM на случай линейной неразделимости, позволим алгоритму допускать ошибки на обучающих объектах, но при этом постараемся, чтобы ошибок было поменьше. Введём набор дополнительных переменных  $\xi_i > 0$ , характеризующих величину ошибки на объектах  $\xi_i$ ,  $i = \overline{1, L}$ . Будем исходить из исходной задачи минимизации квадратичной формы (8.3)



$$\begin{cases} (w, w) \rightarrow \min_{w,b}; \\ y_i \cdot ((w, x_i) - b) \geq 1, \quad i = \overline{1, L}. \end{cases}$$

Смягчим в ней ограничения-неравенства, и одновременно введём в минимизируемый функционал штраф за суммарную ошибку:

$$\begin{cases} \frac{1}{2}(w, w) + C \sum_{i=1}^L \xi_i \rightarrow \min_{w,b,\xi}; \\ y_i \cdot ((w, x_i) - b) \geq 1 - \xi_i, \quad i = \overline{1, L}; \\ \xi_i \geq 0. \end{cases} \quad (8.9)$$

Положительная константа  $C$  является управляющим параметром метода и позволяет находить компромисс между максимизацией разделяющей полосы и минимизацией суммарной ошибки. Вернёмся к задаче (8.9) и запишем её функцию Лагранжа:

$$\mathbf{L}(w, b, \xi; \lambda, \eta) = \frac{1}{2}(w, w) - \sum_{i=1}^L \lambda_i (y_i \cdot ((w, x_i) - b) - 1) - \sum_{i=1}^L \xi_i \cdot (\lambda_i + \eta_i - C), \quad (8.10)$$

где  $\eta = (\eta_1, \eta_2, \dots, \eta_L)$  – вектор переменных, двойственных переменным  $\xi = (\xi_1, \xi_2, \dots, \xi_L)$ . Как и в прошлый раз, условия Куна–Таккера сводят задачу к поиску седловой точки функции Лагранжа:

$$\left\{ \begin{array}{l} \mathbf{L}(w, b, \xi; \lambda, \eta) \rightarrow \min_{w, b, \xi} \max_{\lambda, \eta}; \\ \xi_i \geq 0, \lambda_i \geq 0, \eta_i \geq 0, \quad i = \overline{1, L}; \\ \lambda_i = 0 \text{ или } y_i \cdot ((w, x_i) - b) = 1 - \xi_i, \quad i = \overline{1, L}, \\ \eta_i = 0 \text{ или } \xi_i = 0, \quad i = \overline{1, L}. \end{array} \right.$$

В последних двух строках записаны условия *дополняющей нежесткости*. Необходимым условием седловой точки является равенство нулю производных Лагранжиана. Из уравнений  $\frac{\partial \mathbf{L}}{\partial w} = \frac{\partial \mathbf{L}}{\partial b} = \frac{\partial \mathbf{L}}{\partial \xi_i} = 0$  получаются три полезных соотношения:

$$w = \sum_{i=1}^L \lambda_i y_i x_i, \quad (8.11)$$

$$\sum_{i=1}^L \lambda_i y_i = 0, \quad (8.12)$$

$$\eta_i + \lambda_i = C, \quad i = \overline{1, L}. \quad (8.13)$$

В силу соотношения (8.13) в Лагранжиане обнуляются все члены, содержащие переменные  $\xi_i$  и  $\eta_i$  и он принимает тот же вид, что и в случае линейной разделимости. Параметры разделяющей поверхности  $w$  и  $b$ , согласно формулам (8.11) и (8.12), также выражаются только через двойственные переменные  $\lambda_i$ . Таким образом, задача снова сводится к квадратичному программированию относительно двойственных переменных  $\lambda_i$ . Единственное отличие от линейно разделимого случая состоит в появлении ограничения сверху  $\lambda_i \leq C$ :



$$\left\{ \begin{array}{l} -\mathbf{L}(\lambda) = \frac{1}{2} \sum_{i=1}^L \sum_j^L \lambda_i \lambda_j y_i y_j (x_i, x_j) - \sum_{i=1}^L \lambda_i \rightarrow \min_{\lambda}; \\ 0 \leq \lambda_i \leq C, \quad i = \overline{1, L}; \\ \sum_{i=1}^L \lambda_i y_i = 0. \end{array} \right. \quad (8.14)$$

На практике для построения SVM решают именно эту задачу, а не (8.6), так как гарантировать линейную разделимость выборки в общем случае не представляется возможным. Этот вариант алгоритма называют *SVM с мягким зазором (soft-margin SVM)*, тогда как в линейно разделимом случае говорят об *SVM с жёстким зазором (hard-margin SVM)*.

Для алгоритма классификации сохраняется формула (8.8) с той лишь разницей, что теперь ненулевыми  $\lambda_i$  обладают не только опорные объекты, но и объекты-нарушители. В определённом смысле это недостаток SVM, поскольку нарушителями часто оказываются шумовые выбросы, а построенное на них решающее правило, по сути дела, опирается на шум.

Константу  $C$  обычно выбирают по критерию скользящего контроля. Это трудоёмкий способ, так как задачу приходится решать заново при каждом значении  $C$ . Если есть основания полагать, что выборка почти линейно разделима и лишь объекты-выбросы классифицируются неверно, то можно применить фильтрацию выбросов. Сначала задача решается при некотором  $C$  и из выборки удаляется небольшая доля объектов, имеющих наибольшую величину ошибки  $\xi_i$ . После этого задача решается заново по усечённой выборке. Возможно, придётся проделать несколько таких итераций, пока оставшиеся объекты не окажутся линейно разделимыми.

### 8.3 Ядра и спрямляющие пространства

Другой подход к решению проблемы линейной неразделимости основан на переходе от исходного пространства признаков описаний объектов  $X$  к новому пространству  $H$  с помощью некоторого преобразования  $\psi : X \rightarrow H$ . Если пространство  $H$  имеет достаточно высокую размерность, то можно надеяться, что в нём выборка окажется линейно разделимой. Пространство  $H$  называют *спрямляющим*. Если предположить, что признаковыми описаниями объектов являются векторы  $\psi(x_i)$ , а не векторы  $x_i$ , то построение *SVM* проводится точно так же, как и ранее. Единственное отличие состоит в том, что скалярное произведение  $(x, x')$  в пространстве  $X$  всюду заменяется скалярным произведением  $(\psi(x), \psi(x'))$  в пространстве  $H$ . Отсюда вытекает естественное требование: пространство  $H$  должно быть наделено скалярным произведением, в частности, подойдёт любое евклидово, а в общем случае и гильбертово, пространство.

**Определение.** Функция  $K : X \times X \rightarrow \mathbb{R}$  называется *ядром (kernel function)*, если она представима в виде  $K(x, x') = (\psi(x), \psi(x'))$  при некотором отображении  $\psi : X \rightarrow H$ , где  $H$  – пространство со скалярным произведением.

Постановка задачи (8.14), и сам алгоритм классификации (8.8) зависят только от скалярных произведений объектов, но не от самих признаков описаний. Это означает, что скалярное произведение  $\langle x, x' \rangle$  можно формально заменить ядром

$K(x, x')$ . Поскольку ядро в общем случае нелинейно, такая замена приводит к существенному расширению множества реализуемых алгоритмов  $a: X \rightarrow Y$ . Более того, можно вообще не строить спрямляющее пространство  $H$  в явном виде и вместо подбора отображения  $\psi: X \rightarrow H$  заниматься непосредственно подбором ядра. Можно пойти ещё дальше и вовсе отказаться от признаков описаний объектов. Во многих практических задачах объекты изначально задаются информацией об их попарном взаимоотношении, например, отношении сходства. Если эта информация допускает представление в виде двуместной функции  $K(x, x')$ , удовлетворяющей аксиомам скалярного произведения, то задача может решаться методом *SVM*. Для такого подхода недавно был придуман термин *беспризнаковое распознавание (featureless recognition)*, хотя многие давно известные метрические алгоритмы классификации (*kNN*, *RBF* и др.) также не требуют задания признаков описаний.

**Теорема Мерсера.** Функция  $K(x, x')$  является ядром тогда и только тогда, когда она симметрична,  $K(x, x') = K(x', x)$  и неотрицательно определена:

$$\int_X \int_X K(x, x') g(x) g(x') dx dx' \geq 0 \text{ для любой функции } g: X \rightarrow \mathbb{R}.$$

**Конструктивные способы построения ядер.** Следующие правила порождения позволяют строить ядра в практических задачах.

1. Произвольное скалярное произведение  $K(x, x') = (x, x')$  является ядром.
2. Константа  $K(x, x') = 1$  является ядром.

3. Произведение ядер  $K(x, x') = K_1(x, x')K_2(x, x')$  является ядром.
4. Для любой функции  $\psi: X \rightarrow \mathbb{R}$  произведение  $K(x, x') = \psi(x)\psi(x')$  является ядром.
5. Линейная комбинация ядер с неотрицательными коэффициентами  $K(x, x') = \alpha_1 K_1(x, x') + \alpha_2 K_2(x, x')$  является ядром.
6. Композиция произвольной функции  $\phi: X \rightarrow X$  и произвольного ядра  $K_0$  является ядром:  

$$K(x, x') = K_0(\phi(x), \phi(x')).$$
7. Если  $s: X \times X \rightarrow \mathbb{R}$  – произвольная симметричная интегрируемая функция, то  $K(x, x') = \int_x s(x, z)s(x', z)dz$  является ядром.
8. Функция вида  $K(x, x') = k(x - x')$  является ядром тогда и только тогда, когда Фурье-образ  $F[k](\omega) = (2\pi)^{n/2} \int_x e^{-i(\omega, x)} k(x) dx$  неотрицателен.
9. Композиция произвольного ядра  $K_0$  и произвольной функции  $f: \mathbb{R} \rightarrow \mathbb{R}$ , представимой в виде сходящегося степенного ряда с неотрицательными коэффициентами  $K(x, x') = f(K_0(x, x'))$ , является ядром. В частности, функции  $f(z) = e^z$  и  $f(z) = \frac{1}{1-z}$  от ядра являются ядрами.

## 9 ЗАДАЧИ ДЛЯ САМОСТОЯТЕЛЬНОГО РЕШЕНИЯ

1. Доказать, что при одинаковом или противоположном ранжировании факторов  $X$  и  $Y$  коэффициент ранговой корреляции Кендалла  $\tau_{xy} = \pm 1$ .

2. Построить МНК-оценки параметров  $a, b, c, d$  несимметричной билинейной функции двух переменных  $z(x, y) = axy + bx + cy + d$  по наблюдениям в девяти точках:

$$\begin{aligned}z_1 &= z(1; 1), & z_2 &= z(1; 0), & z_3 &= z(1; -1), \\z_4 &= z(0; 1), & z_5 &= z(0; 0), & z_6 &= z(0; -1), \\z_7 &= z(-1; 1), & z_8 &= z(-1; 0), & z_9 &= z(-1; -1).\end{aligned}$$

3. Доказать, что при удвоении числа наблюдений в каждой точке увеличивается точность оценивания параметров линейной регрессии. Рассмотреть пример МНК-оценивания параметров  $a, b$  линейной зависимости  $y = ax + b$  по независимым наблюдениям в трёх точках  $\{-1, 0, 1\}$  и по двойным независимым наблюдениям в трёх точках  $\{-1, -1, 0, 0, 1, 1\}$ . Как сильно уменьшается при этом погрешность оценивания параметров?

4. Построить МНК-оценки параметров  $a, b$  нелинейной функции  $y = ax^b$  по трём наблюдениям  $y_1, y_2, y_3$  в точках  $x_1 = 1, x_2 = 2, x_3 = 3$ .

*Указание.* Применить логарифмирование исходного уравнения.

5. Найти несмещённую оценку дисперсии  $\sigma^2$  шума наблюдения  $\xi$  в стандартной классической модели линейной

регрессии  $y = ax + b + \xi$  по наблюдениям  $y_1, y_2, y_3, y_4, y_5$  в точках  $x_1 = -2, x_2 = -1, x_3 = 0, x_4 = 1, x_5 = 2$ .

**6.** Производится измерение веса предмета на цифровых весах, имеющих случайную погрешностью  $\xi$  (среднее значение погрешности равно нулю, дисперсия равна  $\sigma^2$ ). Определить, как изменяется погрешность измерения веса предмета при увеличении числа независимых измерений  $n$ .

**7.** Измерения веса двух предметов производятся на чашечных весах со стрелкой. Стрелка показывает разность весов предметов, располагающихся на двух чашках, с некоторой случайной погрешностью  $\xi$ , имеющей нулевое среднее и дисперсию  $\sigma^2$ . Требуется определить вес каждого предмета с максимальной точностью. Допускается провести четыре независимых измерения.

*Подсказка.* Рассмотреть случаи, когда предметы взвешиваются отдельно по одному и когда два предмета помещаются на различные чашки весов.

**8.** Сформулировать задачу классической линейной регрессии для определения весов  $\theta_1, \theta_2, \theta_3$  трёх предметов при взвешивании на чашечных весах со стрелкой. Построить матрицу эксперимента  $F$  и рассчитать погрешности определения весов предметов, если известно, что независимые измерения производятся с некоторой случайной погрешностью  $\xi$ , имеющей нулевое среднее и дисперсию  $\sigma^2$ . Рассчитать дисперсию ошибки оценивания веса каждого предмета: 1) при независимом взвешивании каждого из трёх предметов; 2) при размещении одновременно всех трёх предметов на весах в разных комбинациях. Общее число взвешиваний в каждом случае равно 6.

*Подсказка.* На чашечных весах определяется разность весов предметов, находящихся в двух чашках. На каждую чашку

весов может помещаться несколько предметов (или ни одного). Эксперимент определяется значениями базовой переменной

$$f_i = \begin{cases} +1, & \text{если } i^{\text{й}} \text{ предмет находится на левой чашке весов,} \\ -1, & \text{если } i^{\text{й}} \text{ предмет находится на правой чашке весов,} \\ 0, & \text{если } i^{\text{й}} \text{ предмет не участвует во взвешивании.} \end{cases}$$

(1)

$$F = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{pmatrix}, \quad F^T F = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{pmatrix} = 2 \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix},$$

$$D\hat{\theta} = \sigma^2 (F^T F)^{-1} = \frac{\sigma^2}{2} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

(2)

$$F = \begin{pmatrix} 1 & 1 & -1 \\ 1 & 1 & -1 \\ 1 & -1 & 1 \\ 1 & -1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}, \quad F^T F = \begin{pmatrix} 6 & 2 & 2 \\ 2 & 6 & -2 \\ 2 & -2 & 6 \end{pmatrix} = 2 \begin{pmatrix} 3 & 1 & 1 \\ 1 & 3 & -1 \\ 1 & -1 & 3 \end{pmatrix},$$

$$D\hat{\theta} = \sigma^2 (F^T F)^{-1} = \frac{\sigma^2}{2} \begin{pmatrix} 3 & 1 & 1 \\ 1 & 3 & -1 \\ 1 & -1 & 3 \end{pmatrix}^{-1} = \frac{\sigma^2}{32} \begin{pmatrix} 8 & ? & ? \\ ? & 8 & ? \\ ? & ? & 8 \end{pmatrix} = \frac{\sigma^2}{4} \begin{pmatrix} 1 & ? & ? \\ ? & 1 & ? \\ ? & ? & 1 \end{pmatrix}.$$

9. Доказать, что  $\overline{(x - \bar{x})^2} = \overline{x^2} - \bar{x}^2$ , где  $x^1, x^2, x^3, \dots, x^N$  – произвольный набор чисел,  $\overline{(\cdot)}$  – операция вычисления среднего арифметического.

10. Доказать общее соотношение для  $SS$ -технологии Фишера

$$SS_{\text{общ}} = SS_{\text{факт}} + SS_{\text{ош}}$$

на следующем примере:

Производится многократное измерение некоторого параметра с помощью  $m$  различных приборов. Каждым прибором выполняется  $n$  независимых измерений (всего  $N=mn$  измерений). Результаты измерений:  $y_1^1, y_1^2, \dots, y_1^n, y_2^1, y_2^2, \dots, y_2^n, y_m^1, y_m^2, \dots, y_m^n$ . Средние значения измерений для каждого прибора:  $\bar{y}_k = \sum_{i=1}^n y_k^i, k = \overline{1; m}$ .

*Подсказка:*

Общая вариабельность наблюдаемой переменной:

$$SS_{\text{общ}} = R_1^2 = \sum_{k=1}^m \sum_{i=1}^n (y_k^i - \bar{y})^2, \quad \bar{y} = \sum_{k=1}^m \sum_{i=1}^n y_k^i = \sum_{k=1}^m \sum_{i=1}^n \bar{y}_k - \text{общее}$$

среднее.

Изменчивость наблюдаемой переменной за счёт случайного шума наблюдения (остаточная сумма квадратов ошибок):

$$SS_{\text{ош}} = R_0^2 = \sum_{k=1}^m \sum_{i=1}^n (y_k^i - \bar{y}_k)^2.$$

Изменчивость наблюдаемой переменной за счёт изменения уровней входного фактора:

$$SS_{\text{факт}} = \sum_{k=1}^m \sum_{i=1}^n (\bar{y}_k - \bar{y})^2.$$



11. Дан набор из 8 точек на плоскости:

	A	B	C	D	E	F	G	H
X	1	3	4	5	1	4	1	2
Y	3	3	3	3	2	2	1	1

Выполнить кластеризацию точек по двум алгоритмам:  
 1) *k-means* (число кластеров равно 2, метрика – манхеттенская);  
 2) иерархическая *агломеративная* кластеризация (расстояние между кластерами измеряется по наиболее удалённому объекту, метрика – манхеттенская). Нарисовать дендрограмму для иерархической кластеризации.

12. Построить байесовский классификатор для двух классов наблюдений по одному признаку, распределённому по равномерному (первый класс) и гауссовскому (второй класс) законам с одинаковыми математическими ожиданиями и одинаковыми дисперсиями, если априорные вероятности классов одинаковы. Найти вероятность ошибочной классификации.

13. Найти разделяющие поверхности четырёх классов **I, II, III, IV** гауссовских наблюдений двумерного случайного вектора  $X$  с математическими ожиданиями

$$\mathbf{a}_1 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \mathbf{a}_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix},$$

$$\mathbf{a}_3 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \quad \mathbf{a}_4 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad \text{некоррелированными компонентами:}$$

$$\mathbf{R}_1 = \mathbf{R}_2 = \mathbf{R}_3 = \mathbf{R}_4 = \mathbf{R} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad \text{и априорными вероятностями}$$

$$\text{классов } \mathbf{p}_1 = 1/2, \quad \mathbf{p}_2 = 1/4, \quad \mathbf{p}_3 = 1/8, \quad \mathbf{p}_4 = 1/8.$$

14. Построить решающее правило для трёх признаков для трёх классов независимых гауссовских наблюдений трёхмерного случайного вектора  $X$  с математическими ожиданиями

$$\mathbf{a}_1 = (1 \ 0 \ 0)^T, \quad \mathbf{a}_2 = (0 \ 1 \ 0)^T, \quad \mathbf{a}_3 = (0 \ 0 \ 1)^T \quad \text{и одинаковыми}$$

корреляционными матрицами  $\mathbf{R}_1 = \mathbf{R}_2 = \mathbf{R}_3 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$ , если

априорные вероятности классов равны  $1/2$ ,  $1/3$ ,  $1/6$ .

**15.** Медицинский диагностический прибор анализирует заболевание на основе некоторого признака  $X$ , обеспечивая **чувствительность** диагностики по критерию Неймана–Пирсона на уровне 0.9. Найти уровень **специфичности**, если решение принимается исходя из значения признака  $X$ , распределённого по нормальному закону с параметрами (0;1) при отсутствии заболевания или с параметрами (1;1) при наличии заболевания.

**16.** Построить решающее правило Неймана–Пирсона для разделения двух классов объектов по признаку  $x$ , распределённому по закону Лапласа с плотностью вероятности  $f_1(x) = e^{-|x+1|}$ ,  $x \in (-\infty, \infty)$  (для первого класса) и  $f_2(x) = e^{-|x-1|}$ ,  $x \in (-\infty, \infty)$  (для второго класса). Установить связь между вероятностями ошибок первого и второго рода, считая, что пороговое значение  $h \in [-1; 1]$ .

**17.** Производится обнаружение сигнала  $s$  на фоне шума  $\xi$ :

$$x = \begin{cases} s + \xi, & \text{если сигнал присутствует;} \\ \xi, & \text{если наблюдается только шум.} \end{cases}$$

Построить решающее правило Неймана–Пирсона, если шум является гауссовским с нулевым математическим ожиданием и дисперсией  $\sigma^2$  и задана вероятность ложной тревоги  $\varepsilon$ . Найти вероятность пропуска цели.

**18.** Приведите простейший пример, показывающий, что при оптимальной (байесовской) классификации центр (матема-

тическое ожидание) одного класса может оптимально классифицироваться в другой класс.

**19.** Построить линейный классификатор для двух классов независимых гауссовских наблюдений двумерного случайного вектора  $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$  с математическими ожиданиями  $\mathbf{a}_1 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$  и  $\mathbf{a}_2 = \begin{pmatrix} 0 \\ 4 \end{pmatrix}$ , корреляционными матрицами  $\mathbf{R}_1 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ ,  $\mathbf{R}_2 = \begin{pmatrix} 4 & 0 \\ 0 & 1 \end{pmatrix}$  и равными априорными вероятностями классов:  $\mathbf{p}_1 = \mathbf{p}_2 = 1/2$ . Сравнить вероятность ошибки классификации построенного линейного классификатора с вероятностью ошибки байесовской классификации.

**20.** Построить оптимальный линейный классификатор для двух классов объектов по двум признакам  $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$ , распределённых по нормальному закону с математическими ожиданиями  $\mathbf{m}_1 = \begin{pmatrix} 3 \\ 0 \end{pmatrix}$ ,  $\mathbf{m}_2 = \begin{pmatrix} 0 \\ 3 \end{pmatrix}$ , корреляционными матрицами  $\mathbf{R}_1 = \begin{pmatrix} 1 & 1/2 \\ 1/2 & 1 \end{pmatrix}$ ,  $\mathbf{R}_2 = \begin{pmatrix} 1 & -1/2 \\ -1/2 & 1 \end{pmatrix}$  и одинаковыми априорными вероятностями классов. Найти вероятность ошибочной классификации.

**21.** Построить оптимальный линейный классификатор для трёх классов объектов по двум признакам  $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$ , распределённых по нормальному закону с математическими ожиданиями  $\mathbf{m}_1 = \begin{pmatrix} 3 \\ 0 \end{pmatrix}$ ,  $\mathbf{m}_2 = \begin{pmatrix} 0 \\ 3 \end{pmatrix}$ ,  $\mathbf{m}_3 = \begin{pmatrix} 3 \\ 3 \end{pmatrix}$ , равными корреляционными

матрицами  $\mathbf{R}_1 = \mathbf{R}_2 = \mathbf{R}_3 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$  и априорными вероятностями классов  $1/2, 1/3, 1/6$ . Найти вероятность ошибочной классификации.

**22.** Какое расстояние должно быть между центрами двух классов гауссовых наблюдений с дисперсией  $\sigma^2$  для обеспечения вероятности ошибки классификации не более 1%?

**23.** Привести графический пример с минимальным количеством объектов с линейно разделимой выборкой, для которого линейное разделение даёт безошибочную классификацию на два класса, а метод классификации по ближайшему соседу даёт одну ошибку.

**24.** Сравнить с использованием критерия дискриминантного анализа  $\mathbf{J}_1$  эффективность классификации с использованием признаков задачи **19** и задачи **20**.

**25.** Случайный вектор  $X$  имеет ковариационную матрицу  $R_X = \begin{pmatrix} 1 & r \\ r & 1 \end{pmatrix}$ . Найти линейное преобразование  $A: Y = AX$ , диагонализующее ковариационную матрицу результирующего случайного вектора  $Y: R_Y = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ .

**26.** Какое максимальное значение может принимать вероятность ошибки при оптимальной байесовской классификации объектов на два класса?

**27.** Построить МНК-оценки параметров  $a, b$  уравнения плоскости  $z = ax + by$ , проходящей через точку  $(0,0,0)$ , по наблюдениям в  $N$  точках с известными координатами:

$$z_1 = z(x_1, y_1), \quad z_2 = z(x_2, y_2), \quad z_3 = z(x_3, y_3), \quad \dots, \\ z_N = z(x_N, y_N).$$

## Список литературы и интернет-ресурсы

1. Кендалл, М. Дж. Теория распределений / М.Дж. Кендалл, А. Стьюарт. – М.: Наука, 1966.
2. Ермаков, С.М. Математическая теория оптимального эксперимента: учеб. пособие / С.М. Ермаков, А.А. Жиглявский. – М.: Наука, 1987.
3. Шеффе, Г. Дисперсионный анализ / Г. Шеффе. – М.: Наука, 1980.
4. Ту, Дж. Принципы распознавания образов: пер. с англ. / Дж. Ту, Р. Гонсалес. – М.: Мир, 1978.
5. Фукунага, К. Введение в статистическую теорию распознавания образов: пер. с англ. / К. Фукунага. – М.: Наука, 1979.
6. Воронцов, К.В. Лекции по методу опорных векторов / К.В. Воронцов. – М.: Наука, 2007.
7. Vapnik, V.N. Statistical Learning Theory / V.N. Vapnik. – М.: Наука, 1998.

### Видеолекции

1. Support Vector Machines [Yaser Abu-Mostafa; 2012].  
<http://www.youtube.com/watch?v=eHsErIPJWUU&list=PLo2GY1dZ5AGmsdzWundK-ZeniCqfiqTxx&index=3>
2. Kernel Methods [Yaser Abu-Mostafa; 2012].  
<http://www.youtube.com/watch?v=XUj5JbQihIU&index=2&list=PLo2GY1dZ5AGmsdzWundK-ZeniCqfiqTxx>
3. How SVM (Support Vector Machine) algorithm works [Thales Sehn Körtling; 2014]  
<http://www.youtube.com/watch?v=1NxnPkZM9bc>
4. Support Vector Machine & In-depth Convex Analysis [Sanjeev Sharma]  
<http://www.youtube.com/watch?v=eh3sM4-3heo&index=1&list=PLo2GY1dZ5AGmsdzWundK-ZeniCqfiqTxx>

# Приложение А

## Статистические функции

### А.1 Квантиль распределения

Будем называть *квантилем распределения* непрерывной случайной величины функцию, обратную к функции распределения. Если  $f_{\xi}(x)$  и  $F_{\xi}(x)$  – соответственно плотность вероятности и функция распределения некоторой случайной величины  $\xi$ , то

$$F_{\xi}(x_p) = \int_{-\infty}^{x_p} f_{\xi}(u) du = p, \quad (\text{A.1})$$

где  $x_p = F_{\xi}^{-1}(p)$  – квантиль распределения этой случайной величины на уровне  $p$ ,  $0 < p < 1$ .

Ниже в таблицах используются следующие обозначения:

df – число степеней свободы;  $p$  – вероятность;  $q$  – квантиль.

### А.2 Нормальное распределение

Плотность вероятности нормальной случайной величины с математическим ожиданием  $a$  и дисперсией  $\sigma^2$ :

$$f_{\xi}(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(x-a)^2}{2\sigma^2}\right]. \quad (\text{A.2})$$

Функция нормального распределения (см. табл. А.1 и А.2):

$$F_{\xi}(x) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{(u-a)^2}{2\sigma^2}} du = \Phi\left(\frac{x-a}{\sigma}\right) = \frac{1}{2} + \Phi_0\left(\frac{x-a}{\sigma}\right) = \\ = \frac{1}{2} + \frac{1}{2} \operatorname{erf}\left(\frac{x-a}{\sigma\sqrt{2}}\right),$$

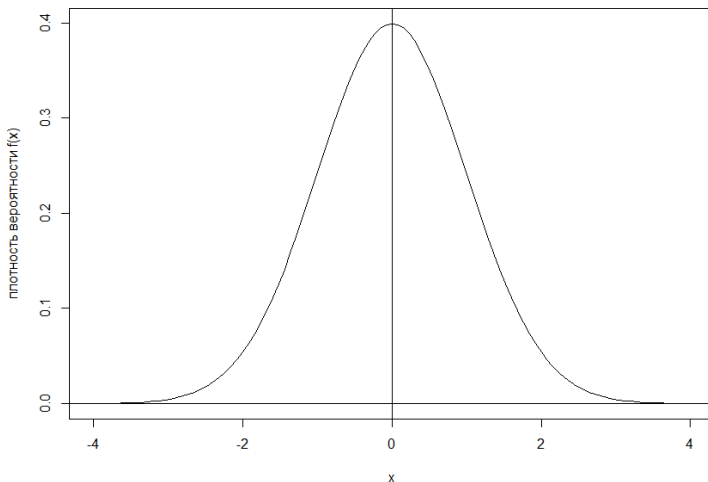
где  $\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-u^2/2} du$  – интеграл вероятности,

$$\Phi(x) = \frac{1}{2} + \Phi_0(x);$$

$\Phi_0(x) = \frac{1}{\sqrt{2\pi}} \int_0^x e^{-u^2/2} du$  – функция Лапласа,  $\Phi_0(x) = \frac{1}{2} \operatorname{erf}\left(\frac{x}{\sqrt{2}}\right)$ ;

$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-u^2} du$  – функция ошибок,  $\operatorname{erf}(x) = 2\Phi_0(x\sqrt{2})$ .

Стандартное нормальное распределение



Плотность вероятности нормального случайного вектора  $\xi = (\xi_1 \ \xi_2 \ \xi_3 \ \dots \ \xi_n)^T \sim N(\mathbf{a}; \mathbf{D}_{\xi})$ :

$$f_{\xi}(x_1, x_2, \dots, x_n) = (2\pi)^{-n/2} (\det \mathbf{D}_{\xi})^{-1/2} \exp\left[-\frac{1}{2}(\mathbf{x} - \mathbf{a})^T \mathbf{D}_{\xi}^{-1}(\mathbf{x} - \mathbf{a})\right], \quad (\text{A.3})$$

где  $\mathbf{a} = \mathbf{M}\xi$  – вектор математических ожиданий,  
 $\mathbf{D}_{\xi} = \mathbf{M}(\xi - \mathbf{a})(\xi - \mathbf{a})^T$  – дисперсионная (корреляционная) матрица,  $\mathbf{x} = (x_1 \ x_2 \ x_3 \ \dots \ x_n)^T$ .

Таблица А.1 – р-квантили стандартного нормального распределения

p =	0.999	0.998	0.995	0.990	0.980	0.950	0.900	0.800
q =	3.0902323	2.8781617	2.5758293	2.3263479	2.0537489	1.6448536	1.2815516	0.8416212

Таблица А.2 – Интеграл вероятности  
(функция стандартного нормального распределения)[z F(z)]

0.000	0.50000	0.050	0.51994	0.100	0.53983	0.150	0.55962	0.200	0.57926	0.250	0.59871	0.300	0.61791	0.350
0.63683	0.400	0.65542	0.400	0.67364	0.72575	0.650	0.74215	0.700	0.75804	0.750	0.77337	0.800	0.78814	0.850
0.500	0.69146	0.550	0.70884	0.600	0.80234	0.900	0.81594	0.950	0.82894	1.000	0.84134	1.050	0.85314	1.100
0.91149	1.400	0.91924	1.450	0.92647	0.86433	1.150	0.87493	1.200	0.88493	1.250	0.89435	1.300	0.90320	1.350
1.500	0.93319	1.550	0.93943	1.600	0.94520	1.650	0.95053	1.700	0.95543	1.750	0.95994	1.800	0.96407	1.850
0.96784	1.900	0.97128	1.950	0.97441	0.98214	2.150	0.98422	2.200	0.98610	2.250	0.98778	2.300	0.98928	2.350
2.000	0.97725	2.050	0.97982	2.100	0.99061	2.400	0.99180	2.450	0.99286	2.500	0.99379	2.520	0.99413	2.540
0.99585	2.660	0.99609	2.680	0.99632	0.99446	2.560	0.99477	2.580	0.99506	2.600	0.99534	2.620	0.99560	2.640
2.700	0.99653	2.720	0.99674	2.740	0.99693	2.760	0.99711	2.780	0.99728	2.800	0.99744	2.820	0.99760	2.840
0.99774	2.860	0.99788	2.880	0.99801	2.900	0.99813	2.910	0.99819	2.920	0.99825	2.930	0.99831	2.940	0.99836
2.950	0.99836	2.950	0.99841	2.960	0.99851	2.980	0.99856	2.990	0.99861	2.990	0.99861	2.991	0.99861	2.992
0.99864	2.998	0.99864	2.999	0.99865	0.99861	2.993	0.99862	2.994	0.99862	2.994	0.99862	2.995	0.99863	2.996
0.99863	2.997	0.99864	2.998	0.99865	3.000	0.99865	3.050	0.99886	3.100	0.99903	3.150	0.99918	3.200	0.99931
3.250	0.99931	3.250	0.99942	3.300	0.99960	3.400	0.99966	3.450	0.99972	3.500	0.99977	3.550	0.99981	3.600
0.99984	3.650	0.99987	3.700	0.99989	0.99994	3.900	0.99995	3.950	0.99996	4.000	0.99997	4.050	0.99997	4.100
0.99998	4.150	0.99998	4.200	0.99999	0.99999	4.400	0.99999	4.450	0.99999	4.500	0.99999	4.550	1.00000	4.600
1.00000	4.600	1.00000	4.700	1.00000	1.00000	4.800	1.00000	4.900	1.00000	5.000	1.00000	5.100	1.00000	5.200
1.00000	5.300	1.00000	5.400	1.00000										

### А.3 Распределение хи-квадрат

Если независимые случайные величины  $\xi_1, \xi_2, \dots, \xi_k$  имеют нормальное распределение с нулевым математическим



ожиданием и единичной дисперсией, то случайная величина  $\chi_k^2 = \sum_{i=1}^k \xi_i^2$  имеет  $\chi^2$ -распределение с  $k$  степенями свободы:

$$\sum_{i=1}^k \xi_i^2 \sim \chi^2(k). \quad (\text{A.4})$$

Если  $\xi = (\xi_1 \ \xi_2 \ \xi_3 \ \dots \ \xi_k)^T \sim N(\mathbf{0}; \mathbf{D}_\xi)$ , то

$$\xi^T \mathbf{D}_\xi^{-1} \xi \sim \chi^2(k). \quad (\text{A.5})$$

Если  $\xi = (\xi_1 \ \xi_2 \ \xi_3 \ \dots \ \xi_k)^T \sim N(\mathbf{0}; \mathbf{I}_k)$ , а  $k \times k$ -матрица  $\mathbf{A}$  является симметричной и идемпотентной:  $\mathbf{A}^T = \mathbf{A}$ ,  $\mathbf{A}^2 = \mathbf{A}$ , то  $\xi^T \mathbf{A} \xi \sim \chi^2(\text{tr } \mathbf{A})$ .



Плотность вероятности распределения хи-квадрат с  $k$  степенями свободы:

$$f_{\chi_k^2}(x) = \begin{cases} \frac{(1/2)^{k/2}}{\Gamma(k/2)} x^{k/2-1} e^{-x/2}, & x \geq 0, \\ 0, & x < 0, \end{cases} \quad (\text{A.6})$$

где  $\Gamma(s) = \int_0^{\infty} t^{s-1} e^{-t} dt$  – гамма-функция.

**Числовые характеристики распределения хи-квадрат:**

$$\mathbf{M}\chi_k^2 = k, \quad \mathbf{D}\chi_k^2 = 2k.$$

**Таблица А.3 – р-квантили хи-квадрат распределения с df степенями свободы**

p = 0.999	0.998	0.995	0.990	0.980	0.950	0.900	0.800												
df = 1	q = 10.8275662	9.5495357	7.8794386	6.6348966	5.4118944	3.8414588													
df = 2	q = 13.8155106	12.4292162	10.5966347	9.2103404	7.8240460	5.9914645													
df = 3	q = 16.2662362	14.7955171	12.8381565	11.3448667	9.8374093	7.8147279													
df = 4	q = 18.4668270	16.9237582	14.8602590	13.2767041	11.6678434	9.4877290													
df = 5	q = 20.5150057	18.9073774	16.7496023	15.0862725	13.3882226	11.0704977													
df = 6	q = 22.4577445	20.7911677	18.5475842	16.8118938	15.0332078	12.5915872													
df = 7	q = 24.3218863	22.6006709	20.2777399	18.4753069	16.6224219	14.0671404													
df = 8	q = 26.1244816	24.3520814	21.9549550	20.0902350	18.1682308	15.5073131													
df = 9	q = 27.8771649	26.0564333	23.5893508	21.6659943	19.6790161	16.9189776													
df = 10	q = 29.5882984	27.7216472	25.1881796	23.2092512	21.1607675	18.3070381													
df = 11	q = 31.2641336	29.3536376	26.7568489	24.7249703	22.6179408	19.6751376													
df = 12	q = 32.9094904	30.9569605	28.2995188	26.2169673	24.0539567	21.0260698													
df = 13	q = 34.5281790	32.5352145	29.8194712	27.6882496	25.4715091	22.3620325													
df = 14	q = 36.1232737	34.0913010	31.3193496	29.1412377	26.8727646	23.6847913													
df = 15	q = 37.6972982	35.6276001	32.8013206	30.5779142	28.2594963	24.9957901													
df = 16	q = 39.2523548	37.1460934	34.2671865	31.9999269	29.6331773	26.2962276													
	23.5418289	20.4650793																	

## Окончание табл. А.2

df = 17	q = 40.7902167	38.6484515	35.7184657	33.4086636	30.9950472	27.5871116
24.7690353	21.6145605					
df = 18	q = 42.3123963	40.1360984	37.1564515	34.8053057	32.3461609	28.8692994
25.9894231	22.7595458					
df = 19	q = 43.8201960	41.6102599	38.5822566	36.1908691	33.6874251	30.1435272
27.2035710	23.9004172					
df = 20	q = 45.3147466	43.0720001	39.9968463	37.5662348	35.0196255	31.4104328
28.4119806	25.0375056					
df = 22	q = 48.2679423	45.9618287	42.7956550	40.2893604	37.6594993	33.9244385
30.8132823	27.3014540					
df = 24	q = 51.1785978	48.8118021	45.5585119	42.9798201	40.2703610	36.4150285
33.1962443	29.5533152					
df = 26	q = 54.0519624	51.6268519	48.2898823	45.6416827	42.8558348	38.8851387
35.5631713	31.7946101					
df = 28	q = 56.8922854	54.4109680	50.9933763	48.2782358	45.4188474	41.3371382
37.9159225	34.0265651					
df = 30	q = 59.7030643	57.1674331	53.6719619	50.8921813	47.9618028	43.7729718
40.2560237	36.2501868					
df = 32	q = 62.4872191	59.8989864	56.3281150	53.4857718	50.4867045	46.1942595
42.5847451	38.4663128					
df = 34	q = 65.2472175	62.6079425	58.9639259	56.0609087	52.9952429	48.6023674
44.9031575	40.6756494					
df = 36	q = 67.9851676	65.2962777	61.5811791	58.6192145	55.4888599	50.9984602
47.2121739	42.8787986					
df = 38	q = 70.7028874	67.9656958	64.1814124	61.1620868	57.9687973	53.3835406
49.5125798	45.0762782					
df = 40	q = 73.4019575	70.6176778	66.7659618	63.6907398	60.4361336	55.7584793
51.8050572	47.2685377					
df = 45	q = 80.0767320	77.1794917	73.1660608	69.9568321	66.5552662	61.6562334
57.5053047	52.7288148					
df = 50	q = 86.6608152	83.6565574	79.4899785	76.1538912	72.6132524	67.5048065
63.1671210	58.1637966					
df = 55	q = 93.1675328	90.0613479	85.7489516	82.2921168	78.6191417	73.3114930
68.7962142	63.5772438					
df = 60	q = 99.6072331	96.4035454	91.9516982	88.3794189	84.5799493	79.0819445
74.3970057	68.9720687					

## А.4 Распределение Стьюдента

Если независимые случайные величины  $\xi_0, \xi_1, \xi_2, \dots, \xi_k$  имеют нормальное распределение с нулевым математическим ожиданием и единичной дисперсией, то случайная величина

$$t_k = \xi_0 \left( \frac{1}{k} \sum_{i=1}^k \xi_i^2 \right)^{-1/2} \sim t(k) \quad (\text{А.7})$$

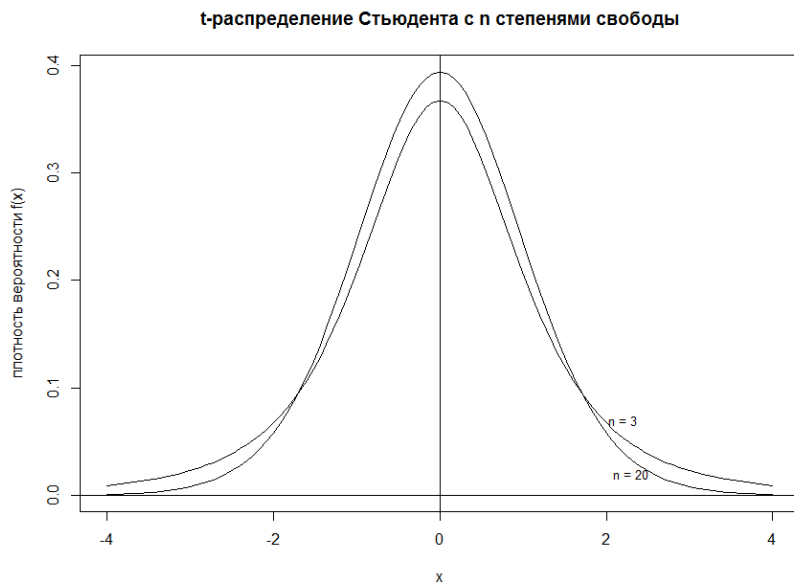
имеет  $t$ -распределение Стьюдента с  $k$  степенями свободы.

Другое определение. Если случайные величины  $\xi_0 \sim N(0;1)$ ,  $\chi_k^2 \sim \chi^2(k)$  независимы, то

$$\xi_0 / \sqrt{\chi_k^2 / k} \sim t(k). \quad (\text{A.8})$$

Плотность вероятности  $t$ -распределения Стьюдента:

$$f_{t_k}(x) = \frac{\Gamma((k+1)/2)}{\sqrt{\pi k} \Gamma(k/2)} \left(1 + \frac{x^2}{k}\right)^{-\frac{k+1}{2}}.$$



**Таблица А.4 – р-квантили t-распределения Стьюдента с df степенями свободы**

p =	0.999	0.998	0.995	0.990	0.980	0.950	0.900	0.800
df = 1	q = 318.3088390	159.1528487	63.6567412	31.8205160	15.8945448	6.3137515	3.0776835	1.3763819
df = 2	q = 22.3271248	15.7639146	9.9248432	6.9645567	4.8487322	2.9199856	1.8856181	1.0606602
df = 3	q = 10.2145319	8.0526131	5.8409093	4.5407029	3.4819088	2.3533634	1.6377444	0.9784723

## Продолжение табл. А.4

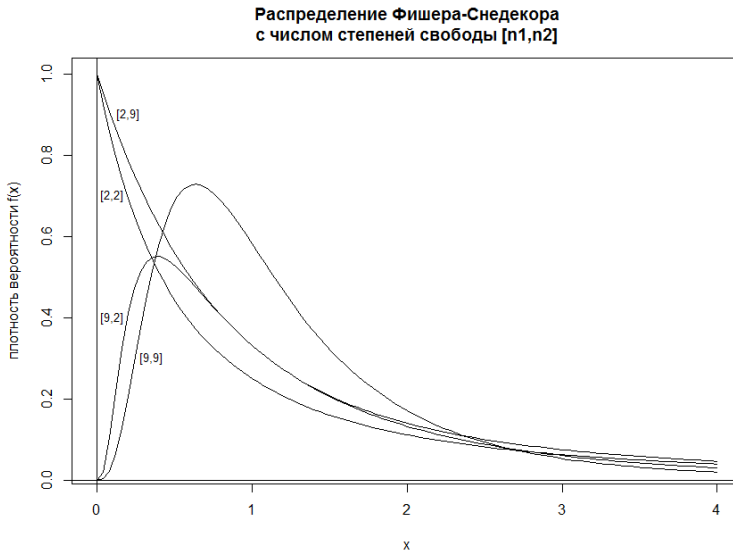
df = 4	q = 7.1731822	5.9513728	4.6040949	3.7469474	2.9985279	2.1318468
1.5332063	0.9409646					
df = 5	q = 5.8934295	5.0302444	4.0321430	3.3649300	2.7565085	2.0150484
1.4758840	0.9195438					
df = 6	q = 5.2076262	4.5241279	3.7074280	3.1426684	2.6122418	1.9431803
1.4397557	0.9057033					
df = 7	q = 4.7852896	4.2071246	3.4994833	2.9979516	2.5167524	1.8945786
1.4149239	0.8960296					
df = 8	q = 4.5007909	3.9909504	3.3553873	2.8964594	2.4489850	1.8595480
1.3968153	0.8888895					
df = 9	q = 4.2968057	3.8345103	3.2498355	2.8214379	2.3984410	1.8331129
1.3830287	0.8834039					
df = 10	q = 4.1437005	3.7162355	3.1692727	2.7637695	2.3593146	1.8124611
1.3721836	0.8790578					
df = 11	q = 4.0247010	3.6237693	3.1058065	2.7180792	2.3281398	1.7958848
1.3634303	0.8755300					
df = 12	q = 3.9296333	3.5495437	3.0545396	2.6809980	2.3027217	1.7822876
1.3562173	0.8726093					
df = 13	q = 3.8519824	3.4886730	3.0122758	2.6503088	2.2816036	1.7709334
1.3501713	0.8701515					
df = 14	q = 3.7873902	3.4378666	2.9768427	2.6244941	2.2637813	1.7613101
1.3450304	0.8680548					
df = 15	q = 3.7328344	3.3948290	2.9467129	2.6024803	2.2485403	1.7530504
1.3406056	0.8662450					
df = 16	q = 3.6861548	3.3579113	2.9207816	2.5834872	2.2353584	1.7458837
1.3367572	0.8646670					
df = 17	q = 3.6457674	3.3258988	2.8982305	2.5669340	2.2238453	1.7396067
1.3333794	0.8632790					
df = 18	q = 3.6104849	3.2978778	2.8784405	2.5523796	2.2137033	1.7340636
1.3303909	0.8620487					
df = 19	q = 3.5794001	3.2731475	2.8609346	2.5394832	2.2047014	1.7291328
1.3277282	0.8609506					
df = 20	q = 3.5518083	3.2511618	2.8453397	2.5279770	2.1966577	1.7247182
1.3253407	0.8599644					
df = 22	q = 3.5049920	3.2137824	2.8187561	2.5083246	2.1828926	1.7171444
1.3212367	0.8582661					
df = 24	q = 3.4667773	3.1831995	2.7969395	2.4921595	2.1715447	1.7108821
1.3178359	0.8568555					
df = 26	q = 3.4349972	3.1577165	2.7787145	2.4786298	2.1620289	1.7056179
1.3149719	0.8556652					
df = 28	q = 3.4081552	3.1361574	2.7632625	2.4671401	2.1539349	1.7011309
1.3125268	0.8546475					
df = 30	q = 3.3851849	3.1176818	2.7499957	2.4572615	2.1469663	1.6972609
1.3104150	0.8537673					
df = 32	q = 3.3653059	3.1016728	2.7384815	2.4486776	2.1409037	1.6938887
1.3085728	0.8529985					
df = 34	q = 3.3479343	3.0876678	2.7283944	2.4411496	2.1355813	1.6909243
1.3069516	0.8523212					
df = 36	q = 3.3326243	3.0753130	2.7194846	2.4344941	2.1308714	1.6882977
1.3055139	0.8517200					
df = 38	q = 3.3190297	3.0643332	2.7115576	2.4285676	2.1266740	1.6859545
1.3042302	0.8511828					
df = 40	q = 3.3068777	3.0545110	2.7044593	2.4232568	2.1229098	1.6838510
1.3030771	0.8506998					

df = 45	q = 3.2814798	3.0339594	2.6895850	2.4121159	2.1150048	1.6794274
1.3006493	0.8496819					
df = 50	q = 3.2614091	3.0176960	2.6777933	2.4032719	2.1087213	1.6759050
1.2987137	0.8488692					
df = 55	q = 3.2451491	3.0045060	2.6682160	2.3960811	2.1036068	1.6730340
1.2971343	0.8482054					
df = 60	q = 3.2317091	2.9935936	2.6602830	2.3901195	2.0993628	1.6706489
1.2958211	0.8476530					

## А.5 Распределение Фишера

Если независимые случайные величины  $\chi_k^2$  и  $\chi_l^2$  имеют  $\chi^2$ -распределение соответственно с  $k$  и  $l$  степенями свободы, то случайная величина  $F_{k,l} = \frac{\chi_k^2/k}{\chi_l^2/l}$  имеет распределение Фишера ( $F$ -распределение, распределение Фишера–Снедекора) с  $k$  и  $l$  степенями свободы:

$$\frac{\chi_k^2/k}{\chi_l^2/l} \sim F(k, l).$$



Если случайная величина  $t_k$  имеет распределение Стьюдента с  $k$  степенями свободы, то случайная величина  $t_k^2$  имеет распределение Фишера с 1 и  $k$  степенями свободы:

$$t_k^2 \sim F(1, k). \quad (\text{A.9})$$

Отсюда следует соотношение между квантилем  $t$ -распределения Стьюдента и квантилем распределения Фишера:

$$t_k^{1-\alpha/2} = \sqrt{F_{1,k}^{1-\alpha}}. \quad (\text{A.10})$$

Плотность вероятности распределения Фишера:

$$f_{F_{k,l}}(x) = \begin{cases} \frac{\Gamma((k+l)/2)}{\Gamma(k/2)\Gamma(l/2)} \left(\frac{k}{l}\right)^{k/2} x^{k/2-1} \left(1 + \frac{k}{l}x\right)^{-(k+l)/2}, & x \geq 0, \\ 0, & x < 0. \end{cases} \quad (\text{A.11})$$

**Таблица А.5 – р-квантили F-распределения Фишера с числом степеней свободы [n1,n2]**

$p = 0.900$

[, 1]	[, 2]	[, 3]	[, 4]	[, 5]	[, 6]	[, 7]	[, 8]	[, 9]
[1, ]	39.86346	8.526316	5.538319	4.544771	4.060420	3.775950	3.589428	3.457919
[2, ]	49.50000	9.000000	5.462383	4.324555	3.779716	3.463304	3.257442	3.113118
[3, ]	53.59324	9.161790	5.390773	4.190860	3.619477	3.288762	3.074072	2.923796
[4, ]	55.83296	9.243416	5.342644	4.107250	3.520196	3.180763	2.960534	2.806426
[5, ]	57.24008	9.292626	5.309157	4.050579	3.452982	3.107512	2.883344	2.726447
[6, ]	58.20442	9.325530	5.284732	4.009749	3.404507	3.054551	2.827392	2.668335
[7, ]	58.90595	9.349081	5.266195	3.978966	3.367899	3.014457	2.784930	2.624135
[8, ]	59.43898	9.366770	5.251671	3.954940	3.339276	2.983036	2.751580	2.589349
								2.469406

## Продолжение табл. А.5

[9,] 59.85759 9.380544 5.239996 3.935671 3.316281 2.957741 2.724678  
2.561238 2.440340

p = 0.950

[,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9]

[1,] 161.4476 18.51282 10.127964 7.708647 6.607891 5.987378 5.591448  
5.317655 5.117355

[2,] 199.5000 19.00000 9.552094 6.944272 5.786135 5.143253 4.737414  
4.458970 4.256495

[3,] 215.7073 19.16429 9.276628 6.591382 5.409451 4.757063 4.346831  
4.066181 3.862548

[4,] 224.5832 19.24679 9.117182 6.388233 5.192168 4.533677 4.120312  
3.837853 3.633089

[5,] 230.1619 19.29641 9.013455 6.256057 5.050329 4.387374 3.971523  
3.687499 3.481659

[6,] 233.9860 19.32953 8.940645 6.163132 4.950288 4.283866 3.865969  
3.580580 3.373754

[7,] 236.7684 19.35322 8.886743 6.094211 4.875872 4.206658 3.787044  
3.500464 3.292746

[8,] 238.8827 19.37099 8.845238 6.041044 4.818320 4.146804 3.725725  
3.438101 3.229583

[9,] 240.5433 19.38483 8.812300 5.998779 4.772466 4.099016 3.676675  
3.388130 3.178893

p = 0.980

[,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9]

[1,] 1012.545 48.50505 20.61798 14.03961 11.322754 9.876365 8.987714  
8.389477 7.960512

[2,] 1249.500 49.00000 18.85813 12.14214 9.454406 8.052094 7.202569  
6.636592 6.233995

[3,] 1350.505 49.16573 18.10970 11.34354 8.670194 7.286980 6.453946  
5.901384 5.509665

[4,] 1405.833 49.24874 17.69377 10.89942 8.233030 6.859440 6.034727  
5.488919 5.102656

[5,] 1440.612 49.29859 17.42876 10.61574 7.952932 6.584743 5.764727  
5.222715 4.839502

[6,] 1464.455 49.33184 17.24510 10.41859 7.757703 6.392778 5.575612  
5.035888 4.654497

[7,] 1481.803 49.35560 17.11028 10.27352 7.613672 6.250824 5.435476  
4.897197 4.516938

[8,] 1494.986 49.37342 17.00711 10.16225 7.502956 6.141480 5.327335  
4.789995 4.410456

[9,] 1505.341 49.38729 16.92561 10.07418 7.415156 6.054611 5.241280  
4.704564 4.325487

p = 0.990

[,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9]

[1,] 4052.181 98.50251 34.11622 21.19769 16.25818 13.745023  
12.246383 11.258624 10.561431



Окончание табл. А.5

[2,] 4999.500 99.00000 30.81652 18.00000 13.27393 10.924767 9.546578  
8.649111 8.021517  
[3,] 5403.352 99.16620 29.45670 16.69437 12.05995 9.779538 8.451285  
7.590992 6.991917  
[4,] 5624.583 99.24937 28.70990 15.97702 11.39193 9.148301 7.846645  
7.006077 6.422085  
[5,] 5763.650 99.29930 28.23708 15.52186 10.96702 8.745895 7.460435  
6.631825 6.056941  
[6,] 5858.986 99.33259 27.91066 15.20686 10.67225 8.466125 7.191405  
6.370681 5.801770  
[7,] 5928.356 99.35637 27.67170 14.97576 10.45551 8.259995 6.992833  
6.177624 5.612865  
[8,] 5981.070 99.37421 27.48918 14.79889 10.28931 8.101651 6.840049  
6.028870 5.467123  
[9,] 6022.473 99.38809 27.34521 14.65913 10.15776 7.976121 6.718752  
5.910619 5.351129

p = 0.995

[,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9]  
[1,] 16210.72 198.5013 55.55196 31.33277 22.78478 18.63500 16.235558  
14.688199 13.613609  
[2,] 19999.50 199.0000 49.79928 26.28427 18.31383 14.54411 12.403957  
11.042412 10.106714  
[3,] 21614.74 199.1664 47.46723 24.25912 16.52977 12.91660 10.882447  
9.596475 8.717055  
[4,] 22499.58 199.2497 46.19462 23.15450 15.55606 12.02753 10.050491  
8.805130 7.955885  
[5,] 23055.80 199.2996 45.39165 22.45643 14.93961 11.46370 9.522059  
8.301799 7.471158  
[6,] 23437.11 199.3330 44.83847 21.97458 14.51326 11.07304 9.155336  
7.951992 7.133850  
[7,] 23714.57 199.3568 44.43410 21.62169 14.20045 10.78592 8.885389  
7.694143 6.884908  
[8,] 23925.41 199.3746 44.12557 21.35198 13.96096 10.56576 8.678115  
7.495906 6.693300  
[9,] 24091.00 199.3885 43.88240 21.13908 13.77165 10.39149 8.513823  
7.338595 6.541090

## Приложение Б

### Сведения из теории матриц

#### Б.1 Спектральное разложение симметричной матрицы

Рассмотрим произвольную *симметричную*  $n \times n$ -матрицу  $\mathbf{A}$ ;  $\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_n$  – её собственные числа.

Матрица  $\mathbf{A}$  может быть представлена в виде:

$$\mathbf{A} = \mathbf{P}\mathbf{\Lambda}\mathbf{P}^T = \sum_{i=1}^n \lambda_i P_i P_i^T, \quad (\text{Б.1})$$

где  $P_1, P_2, P_3, \dots, P_n$  – ортонормированные собственные векторы-столбцы матрицы  $\mathbf{A}$ ,  $\mathbf{P} = [P_1 \ P_2 \ P_3 \ \dots \ P_n]$  – ортогональная  $n \times n$ -матрица собственных векторов матрицы  $\mathbf{A}$  ( $\mathbf{P}^{-1} = \mathbf{P}^T$ ,  $\mathbf{P}\mathbf{P}^T = \mathbf{P}^T\mathbf{P} = \mathbf{I}_n$ ,  $\mathbf{I}_n$  – единичная  $n \times n$ -матрица),  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_n)$  – диагональная матрица собственных чисел матрицы  $\mathbf{A}$ .

Определитель матрицы  $\mathbf{A}$ :

$$\det \mathbf{A} = \prod_{i=1}^n \lambda_i. \quad (\text{Б.2})$$

Если  $n \times n$ -матрица  $\mathbf{A}$  является положительно (неотрицательно) определённой, то

$$[\mathbf{A}]_{ii} > 0 \quad ([\mathbf{A}]_{ii} \geq 0), \quad \lambda_i > 0 \quad (\lambda_i \geq 0), \quad i = \overline{1, n}. \quad (\text{Б.3})$$

#### Б.2 Свойства операции *trace* (след)

Во всех рассмотренных ниже свойствах считается, что фигурирующие в формулах матрицы имеют согласованные размеры.

Если  $\mathbf{A}$  и  $\mathbf{B}$  – произвольные прямоугольные матрицы, то

$$\operatorname{tr} \mathbf{AB} = \operatorname{tr} \mathbf{BA}. \quad (\text{Б.4})$$

Если  $\mathbf{P}$  – ортогональная матрица ( $\mathbf{PP}^T = \mathbf{P}^T\mathbf{P} = \mathbf{I}_n$ ), то

$$\operatorname{tr} \mathbf{PAP}^T = \operatorname{tr} \mathbf{A}. \quad (\text{Б.5})$$

Если  $\mathbf{A}$  – симметричная матрица, то

$$\operatorname{tr} \mathbf{A}^k = \sum_{i=1}^n \lambda_i^k, \quad k = 0, \pm 1, \pm 2, \pm 3, \dots, \quad (\text{Б.6})$$

где  $\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_n$  – собственные числа матрицы  $\mathbf{A}$ .

### Б.3 Линейное преобразование случайного вектора

Если  $\xi$  – случайный вектор, а  $\mathbf{c}$  и  $\mathbf{B}$  – произвольные детерминированные вектор и матрица соответствующих размеров, то

$$\mathbf{M}(\mathbf{c} + \mathbf{B}\xi) = \mathbf{c} + \mathbf{B}\mathbf{M}\xi, \quad \mathbf{D}(\mathbf{c} + \mathbf{B}\xi) = \mathbf{B}\mathbf{D}\xi\mathbf{B}^T. \quad (\text{Б.7})$$

Если  $\xi \sim \mathbf{N}(\mathbf{a}, \mathbf{D}_\xi)$ ,  $\eta = \mathbf{c} + \mathbf{B}\xi$ , то  $\eta \sim \mathbf{N}(\mathbf{a} + \mathbf{c}, \mathbf{B}\mathbf{D}_\xi\mathbf{B}^T)$ .

## Приложение В

### Десять основных алгоритмов анализа данных

На IEEE International Conference on Data Mining (ICDM) in December 2006 выделены десять основных алгоритмов анализа данных (**The top 10 data mining algorithms** – табл. В.1).

Таблица В.1 – Десять основных алгоритмов анализа данных

<b>ID3</b>	<b>ID3 (Iterative Dichotomiser 3)</b> is an algorithm invented by Ross Quinlan used to generate a decision tree from a dataset
<b>C4.5</b>	Алгоритм для <b>построения</b> деревьев решений, разработанный в 1993 Джоном Квинланом (вариант машинного обучения с учителем). C4.5 является усовершенствованной версией алгоритма ID3 того же автора. В частности, в новую версию были добавлены отсечение ветвей, возможность работы с числовыми атрибутами, а также возможность построения дерева из неполной обучающей выборки, в которой отсутствуют значения некоторых атрибутов
<b>k-Means</b>	<b>Метод k-средних</b> (быстрый кластерный анализ) – наиболее популярный метод кластеризации. Был изобретён в 1950-х годах математиком Гуго Штейнгаузом <a href="https://ru.wikipedia.org/wiki/K-means_-_cite_note-1">https://ru.wikipedia.org/wiki/K-means - cite note-1</a> и Стюартом Ллойдом. Действие алгоритма таково, что он стремится минимизировать суммарное квадратичное отклонение точек кластеров от центров этих кластеров
<b>ИСОМАД (Isodata)</b>	<b>ИСОМАД</b> (Итеративный самоорганизующийся метод анализа данных, <b>Isodata</b> – Iterative Self-Organizing Data Analysis Techniques) аналогичен методу <b>k-Means</b> , однако обладает более широким набором параметров и вспомогательных эвристических процедур
<b>SVM</b>	<b>Метод опорных векторов (SVM – support vector machine)</b> – набор схожих алгоритмов <b>обучения с учителем</b> , использующихся для задач классификации и регрессионного анализа. Принадлежит к семейству линейных классификаторов, может также рассматриваться как специальный случай регуляризации по Тихонову. Особым свойством метода опорных

	<p>векторов является непрерывное уменьшение эмпирической ошибки классификации и увеличение зазора, поэтому метод также известен как метод классификатора с максимальным зазором. Основная идея метода – перевод исходных векторов в пространство более высокой размерности и поиск разделяющей гиперплоскости с максимальным зазором в этом пространстве. Две параллельных гиперплоскости строятся по обеим сторонам гиперплоскости, разделяющей наши классы. Разделяющей гиперплоскостью будет гиперплоскость, максимизирующая расстояние до двух параллельных гиперплоскостей. Алгоритм работает в предположении, что чем больше разница или расстояние между этими параллельными гиперплоскостями, тем меньше будет средняя ошибка классификатора</p>
<b>Apriori</b>	<p><b>Apriori</b> – алгоритм поиска ассоциативных правил. Apriori is an algorithm for frequent item set mining and association rule learning over transactional databases. It proceeds by identifying the frequent individual items in the database and extending them to larger and larger item sets as long as those item sets appear sufficiently often in the database. The frequent item sets determined by Apriori can be used to determine association rules which highlight general trends in the database: this has applications in domains such as market basket analysis</p>
<b>EM</b>	<p><b>EM-алгоритм</b> (<i>Expectation-maximization algorithm</i>) – алгоритм, используемый в математической статистике для нахождения оценок максимального правдоподобия параметров вероятностных моделей в случае, когда модель зависит от некоторых скрытых переменных. Каждая итерация алгоритма состоит из двух шагов. На <b>E-шаге (expectation)</b> вычисляется ожидаемое значение функции правдоподобия, при этом скрытые переменные рассматриваются как наблюдаемые. На <b>M-шаге (maximization)</b> вычисляется оценка максимального правдоподобия, таким образом увеличивается ожидаемое правдоподобие, вычисляемое на <b>E-шаге</b>. Затем это значение используется для <b>E-шага</b> на следующей итерации. Алгоритм выполняется до сходимости. Часто EM-алгоритм используют для разделения смеси гауссиан</p>

<b>PageRank</b>	<b>PageRank</b> – один из алгоритмов <b>ссылочного ранжирования</b> . Алгоритм применяется к коллекции документов, связанных гиперссылками (таких, как веб-страницы из всемирной паутины), и назначает каждому из них некоторое численное значение, измеряющее его «важность» или «авторитетность» среди остальных документов. Вообще говоря, алгоритм может применяться не только к веб-страницам, но и к любому набору объектов, связанных между собой взаимными ссылками, то есть к любому графу
<b>AdaBoost</b>	<b>AdaBoost</b> (сокращение от <i>Adaptive Boosting</i> ) – <b>алгоритм усиления классификаторов</b> путем объединения их в комитет, предложенный Йоавом Фройндом и Робертом Шапире. Этот алгоритм может использоваться в сочетании с несколькими алгоритмами классификации для улучшения их эффективности. AdaBoost является адаптивным в том смысле, что каждый следующий комитет классификаторов строится по объектам, неверно классифицированным предыдущими комитетами. AdaBoost чувствителен к шуму в данных и выбросам. Однако он менее подвержен переобучению по сравнению с другими алгоритмами машинного обучения
<b>kNN</b>	Метод <i>k</i> ближайших соседей (англ. <i>k-nearest neighbor algorithm</i> , <i>kNN</i> ) — метод <b>автоматической</b> классификации объектов. Основным принципом <b>метода ближайших соседей</b> является то, что объект присваивается тому классу, который является наиболее распространённым среди соседей данного элемента. Соседи берутся исходя из множества объектов, классы которых уже известны, и, исходя из ключевого для данного метода значения <i>k</i> высчитывается, какой класс наиболее многочислен среди них. Каждый объект имеет конечное количество атрибутов (размерностей). Предполагается, что существует определенный набор объектов с уже имеющейся классификацией
<b>Naive Bayes</b>	<b>Наивный байесовский классификатор</b> – простой вероятностный классификатор, основанный на применении теоремы Байеса со строгими (наивными) предположениями о независимости. В зависимости от точной природы вероятностной модели наивные байесовские классификаторы могут обучаться очень эффективно. Во многих практи-

	<p>ческих приложениях для оценки параметров для наивных байесовых моделей используют метод максимального правдоподобия; другими словами, можно работать с наивной байесовской моделью не веря в байесовскую вероятность и не используя байесовские методы. Несмотря на наивный вид и, несомненно, очень упрощенные условия, наивные байесовские классификаторы часто работают намного лучше во многих сложных жизненных ситуациях. Достоинством наивного байесовского классификатора является малое количество данных для обучения, необходимых для оценки параметров, требуемых для классификации</p>
<b>CART</b>	<p><b>Classification and Regression Trees</b> (Leo Breiman, Jerome Friedman, Richard Olshen, and Charles Stone, 1984) is a binary recursive partitioning procedure capable of processing continuous and nominal attributes both as targets and predictors. Data are handled in their raw form; no binning is required or recommended. Trees are grown to a maximal size without the use of a stopping rule and then pruned back (essentially split by split) to the root via cost-complexity pruning. The next split to be pruned is the one contributing least to the overall performance of the tree on training data (and more than one split may be removed at a time)</p>

Учебное издание

*Храмов Александр Григорьевич*

**МЕТОДЫ И АЛГОРИТМЫ  
ИНТЕЛЛЕКТУАЛЬНОГО АНАЛИЗА ДАННЫХ**

*Учебное пособие*

Редактор Н.С. Куприядова  
Компьютерная верстка Л.Р. Дмитриенко

Подписано в печать 22.07.2019. Формат 60x84 1/16.

Бумага офсетная. Печ. л. 11,0.

Тираж 25 экз. Заказ . Арт. – 27(Р1У)/2019.

ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ  
ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ  
«САМАРСКИЙ НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ  
УНИВЕРСИТЕТ ИМЕНИ АКАДЕМИКА С.П. КОРОЛЕВА»  
(САМАРСКИЙ УНИВЕРСИТЕТ)  
443086 Самара, Московское шоссе, 34.

---

Изд-во Самарского университета.  
443086 Самара, Московское шоссе, 34.