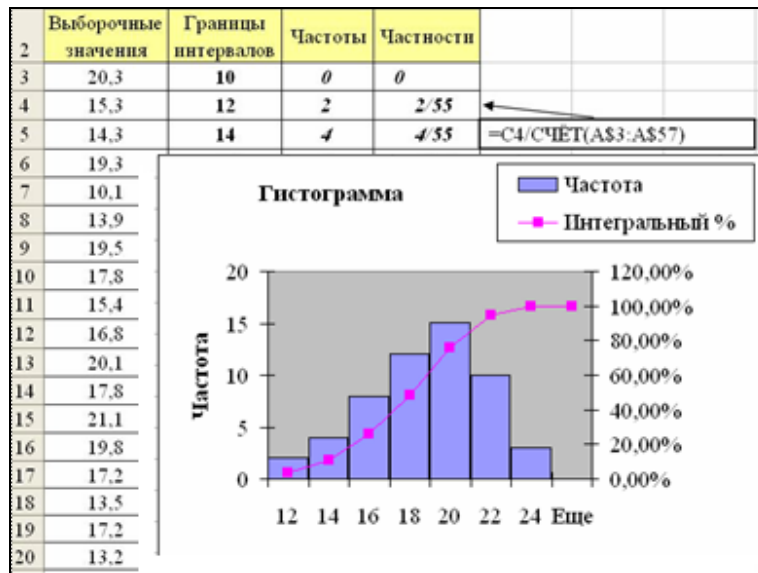




Ю.Е. Воскобойников

Е.И. Тимошенко

## МАТЕМАТИЧЕСКАЯ СТАТИСТИКА (с примерами в Excel)



НОВОСИБИРСК 2006

Ю.Е. Воскобойников, Е.И. Тимошенко

## МАТЕМАТИЧЕСКАЯ СТАТИСТИКА (с примерами в Excel)

УЧЕБНОЕ ПОСОБИЕ

2 издание, переработанное и дополненное

НОВОСИБИРСК 2006

УДК 519.2  
ББК 22.172  
В650

**Воскобойников Ю. Е.**

МАТЕМАТИЧЕСКАЯ СТАТИСТИКА (С ПРИМЕРАМИ В EXCEL) : учеб. пособие / Ю. Е. Воскобойников, Е. И. Тимошенко ; Новосиб. гос. архитектур.-строит. ун-т (Сибстрин). – 2-е изд., перераб. и доп. – Новосибирск : НГАСУ (Сибстрин), 2006. – 152 с.

**ISBN 5-7795-0292-7**

Данное учебное пособие содержит наиболее важные разделы математической статистики: точечное и интервальное оценивание параметров распределений, проверку различных статистических гипотез. Приведено большое количество примеров, которые позволят студентам лучше усвоить не только общетеоретические положения, но и возможные области приложения математической статистики.

Учебное пособие написано в соответствии с программой курса "Математическая статистика" для студентов специальности 080502 "Экономика и управление на предприятии (в строительстве)". Также оно будет полезно студентам других специальностей строительных вузов.

Печатается по решению издательско-библиотечного совета  
НГАСУ (Сибстрин)

Рецензенты:

- С.М. Зеркаль, д-р техн. наук, профессор, вед. науч. сотр. (Институт математики СО РАН);
- А.В. Федоров, д-р физ.-мат. наук, профессор, завлабораторией (Институт теоретической и прикладной механики СО РАН)

**ISBN 5-7795-0292-7**

- © Воскобойников Ю.Е., Тимошенко Е.И., 2006
- © Новосибирский государственный архитектурно-строительный университет (Сибстрин), 2006

## ОГЛАВЛЕНИЕ

ВВЕДЕНИЕ .....	5
1. МАТЕМАТИЧЕСКАЯ СТАТИСТИКА И ЕЕ ЗАДАЧИ.....	9
1.1. Задачи математической статистики .....	9
1.2. Решение задач математической статистики в табличном процессоре Excel .....	11
2. ГЕНЕРАЛЬНАЯ И ВЫБОРОЧНАЯ СОВОКУПНОСТИ. ВЫБОРОЧНЫЕ ХАРАКТЕРИСТИКИ .....	13
2.1. Генеральная и выборочная совокупности .....	13
2.2. Свойства выборочной совокупности .....	14
2.3. Вариационные ряды.....	16
2.4. Выборочная функция распределения. Гистограмма .....	19
2.5. Выборочное среднее и выборочная дисперсия .....	24
2.6. Вычисление выборочных характеристик в Excel .....	29
3. ТОЧЕЧНЫЕ ОЦЕНКИ НЕИЗВЕСТНЫХ ПАРАМЕТРОВ .....	40
3.1. Определение и свойства точечной оценки .....	40
3.2. Точечная оценка математического ожидания .....	45
3.3. Точечные оценки дисперсии.....	47
3.4. Точечная оценка вероятности события.....	51
3.5. Метод максимального правдоподобия.....	52
3.6. Вычисление точечных оценок в Excel .....	61
4. ИНТЕРВАЛЬНЫЕ ОЦЕНКИ НЕИЗВЕСТНЫХ ПАРАМЕТРОВ .....	72
4.1. Некоторые распределения выборочных характеристик.....	72
4.2. Понятие интервальной оценки параметра случайной величины .....	78
4.3. Интервальные оценки математического ожидания нормального распределения.....	79
4.4. Интервальные оценки дисперсии нормального распределения.....	84

4.5. Интервальная оценка вероятности события .....	86
4.6. Вычисление границ доверительных интервалов в Excel.....	89
5. ПРОВЕРКА СТАТИСТИЧЕСКИХ ГИПОТЕЗ .....	92
5.1. Понятие статистической гипотезы. Основные этапы проверки гипотезы .....	92
5.2. Проверка гипотезы о числовом значении математического ожидания нормального распределения .....	100
5.3. Проверка гипотезы о числовом значении дисперсии нормального распределения .....	106
5.4. Проверка гипотезы о числовом значении вероятности события.....	109
5.5. Проверка гипотезы о равенстве математических ожиданий двух нормальных распределений .....	113
5.6. Проверка гипотезы о равенстве математических ожиданий двух произвольных распределений по выборкам большого объема.....	116
5.7. Проверка гипотезы о равенстве математических ожиданий двух нормальных распределений с неизвестными, но равными дисперсиями .....	117
5.8. Проверка гипотезы о равенстве дисперсий двух нормальных распределений .....	121
5.9. Проверка гипотезы о законе распределения с применением критерия согласия Пирсона.....	125
5.10. Проверка гипотезы о независимости двух генеральных совокупностей с применением критерия $\chi^2$ .....	133
5.11. Проверка статистических гипотез в Excel.....	136
6. ВОПРОСЫ ДЛЯ САМОКОНТРОЛЯ .....	146
ЗАКЛЮЧЕНИЕ.....	145
БИБЛИОГРАФИЧЕСКИЙ СПИСОК .....	149
ПРИЛОЖЕНИЕ .....	150

## ВВЕДЕНИЕ

Цель науки – описание, объяснение и предсказание явлений действительности на основе установленных законов, что позволяет находить решения в типичных ситуациях. Многие явления окружающего мира взаимно связаны и влияют одно на другое. Проследить все связи и определить влияние каждой из них на явление не всегда представляется возможным. Поэтому ограничиваются изучением влияния лишь основных факторов, определяющих изучаемое явление. В основе выявления этих связей лежит наблюдение. При этом для обнаружения общих закономерностей, которым подчиняется явление, необходимо многократно его наблюдать в одинаковых условиях, т.е. соблюдать во всех наблюдениях практически одинаковые значения основных факторов. После накопления полученных таким образом данных возникает главный вопрос: как обработать результаты наблюдений и сделать обоснованные выводы об изучаемых закономерностях? Ответы на этот вопрос и другие вопросы, связанные с обработкой данных, дает *математическая статистика*.

Математическая статистика – наука, изучающая методы обработки результатов наблюдений массовых случайных явлений, обладающих статистической устойчивостью, закономерностью, с целью выявления этой закономерности. Выводы о закономерностях, которым подчиняются явления, изучаемые методами математической статистики, всегда основываются на ограниченном числе наблюдений. Для вынесения обоснованного заключения о закономерностях изучаемого явления математическая статистика опирается на *теорию вероятностей*, которая имеет дело с математическими моделями случайных явлений. Обработав результаты наблюдений, исследователь выдвигает ряд гипотез (предположений) о том, что рассматриваемое явление можно описать той или иной вероятностной теоретической моделью. Далее, используя математико-статистические методы, можно дать ответ на вопрос, какую из гипотез или моделей следует принять, которая и будет считаться искомой закономерностью изучаемого явления. Правомерен такой вывод или нет, покажет практика использования выбранной модели. Таково типичное содержание *математико-статистического исследования*.

# 1. МАТЕМАТИЧЕСКАЯ СТАТИСТИКА И ЕЕ ЗАДАЧИ

## 1.1. Задачи математической статистики

*Математическая статистика – наука, изучающая методы исследования закономерностей в массовых случайных явлениях и процессах по данным, полученным из конечного числа наблюдений за ними.*

Построенные на основании этих методов закономерности относятся не к отдельным испытаниям, из повторения которых складывается данное массовое явление, а представляют собой утверждения об общих вероятностных характеристиках данного процесса. Такими характеристиками могут быть вероятности, плотности распределения вероятностей, математические ожидания, дисперсии и т.п.

Найденные характеристики позволяют построить *вероятностную модель* изучаемого явления. Применяя к этой модели методы теории вероятностей, исследователь может решать технико-экономические задачи, например, определять вероятность безотказной работы агрегата в течение заданного отрезка времени. Таким образом, теория вероятностей по вероятностной модели процесса предсказывает его поведение, а математическая статистика по результатам наблюдений за процессом строит его вероятностную модель. В этом состоит тесная взаимосвязь между данными науками.

Очевидно, что для обнаружения закономерностей случайного массового явления необходимо провести сбор статистических сведений, т.е. сведений, характеризующих отдельные единицы каких-либо массовых явлений. Пусть, например, мы располагаем материалом о числе дефектных изделий в изготовленной в определенных условиях партии продукции. Проблемы возникают тогда, когда на основании этой информации мы захотим сделать выводы относительно качества производства продукции, выпускаемой предприятием. Нас может интересовать вероятность производства дефектного изделия, средняя долговечность всех выпускаемых изделий и т.д. Собранный материал рассматривается лишь как некоторая пробная группа, одна из многих возможных пробных групп. Конечно, выводы, сделанные на основании этого ограниченного

числа наблюдений, отражают данное массовое явление лишь приближенно. Математическая статистика указывает, как наилучшим способом использовать имеющуюся информацию для получения по возможности более точных характеристик массового явления.

Конкретизируем задачи, решение которых будет рассмотрено в данном пособии.

1. *Оценка неизвестной функции распределения и функции плотности.* По результатам  $n$  независимых испытаний над случайной величиной  $X$  получены ее значения

$$x_1, x_2, \dots, x_n.$$

Требуется оценить, хотя бы приближенно, неизвестные функции распределения  $F(x)$  и плотности  $p(x)$ .

2. *Оценка неизвестных параметров распределения.* Поясним задачу на примере нормального распределения генеральной совокупности, зависящей от двух параметров  $\alpha$  и  $\sigma$ . Требуется на основании имеющихся данных приближенно найти значение этих параметров. Для этого изучаются некоторые случайные величины и на основе их свойств определяется точность полученных оценок. Мы будем различать два случая: когда имеется достаточно большое количество статистических данных и когда их набор ограничен. Во втором случае будем строить интервалы со случайными границами, на которые попадают неизвестные параметры распределения.
3. *Проверка статистических гипотез.* Предположим, например, что игральная кость подбрасывается  $n$  раз, причем  $n_i$  ( $i = 1, \dots, 6$ ) означает количество появлений  $i$  очков. Если кость симметрична, то любое количество очков должно появиться практически одинаковое число раз (при условии, что  $n$  достаточно велико). Это следует из известной теоремы Бернулли, утверждающей, что относительная частота  $\frac{n_i}{n}$  близка к вероятности  $p = \frac{1}{6}$ . Однако между числами  $\frac{n_i}{n}$  могут быть различия. Возникает вопрос: насколько эти различия согласованы с гипотезой о симметричности игральной кости? Разра-

ботаны методы, позволяющие дать ответы на подобные вопросы с заданной надежностью.

При обращении к понятиям теории вероятностей мы будем опираться на учебное пособие [1].

### **1.2. Решение задач математической статистики в табличном процессоре Excel**

Решение задач математической статистики обуславливает существенный объем вычислений, связанный с численной реализацией необходимого вычислительного алгоритма и графической интерпретацией результатов решения. Этому моменту в учебной литературе уделяется крайне мало внимания, что затрудняет использование методов математической статистики на практике. Поэтому одной из основных целей данного пособия является *изложение численных методик решения задач математической статистики в вычислительной среде табличного процессора Excel 2003*. Для каждой из рассматриваемых задач математической статистики кроме теоретических положений даются фрагменты документов Excel 2003, реализующих алгоритмы решения задачи. При этом алгоритм решения может быть реализован путем программирования необходимых выражений в ячейках электронной таблицы или путем обращения к стандартным функциям или модулям Excel 2003. В учебном пособии будут использоваться обе рассмотренные возможности реализации требуемого вычислительного алгоритма. Поэтому предполагается, что читатель имеет достаточные навыки для реализации вычислений в Excel с использованием:

- программирования арифметических выражений в ячейках электронной таблицы;
- функций Excel (в основном математических и статистических).

**Замечание 1.1.** При описании той или иной функции в качестве *формальных параметров* используются *имена переменных*, определенные в тексте пособия. При обращении к функции в качестве *фактических параметров* могут использоваться *константы, адреса ячеек, диапазоны адресов и арифметические выражения*. Например, описание функции для вычисления среднего арифметического значения (выборочного среднего) имеет вид:

$$\text{СРЗНАЧ}(x_1; x_2; \dots; x_m),$$

где  $x_1, x_2, \dots, x_m$  – формальные параметры, число которых не превышает 30 ( $m \leq 30$ ). Для вычисления среднего значения величин, находящихся в ячейках В3, В4, В5, В6, С3, С4, С5, С6, обращение к функции в соответствующей ячейке имеет вид:

$$=\text{СРЗНАЧ}(\text{В3:В6}; \text{С3:С6}),$$

т.е. в качестве фактических параметров используются два диапазона ячеек. ♦

**Замечание 1.2.** Так как в запрограммированной ячейке выводится результат вычислений и не видно самого запрограммированного выражения, то в некоторых случаях рядом с результатом приводится (в другой ячейке) запрограммированное выражение (своеобразный комментарий к выполняемым вычислениям). В случаях, когда не очевидно, к какой ячейке относится приводимое выражение, используется стрелка, указывающая на нужную ячейку. ♦

## 2. ГЕНЕРАЛЬНАЯ И ВЫБОРОЧНАЯ СОВОКУПНОСТИ. ВЫБОРОЧНЫЕ ХАРАКТЕРИСТИКИ

### 2.1. Генеральная и выборочная совокупности

Для обнаружения закономерностей, описывающих исследуемое массовое явление, необходимо иметь опытные данные, полученные в результате обследования соответствующих объектов, отображающих изучаемое явление. Например, для определения плотности распределения диаметра шлифованного валика необходимо располагать набором возможных значений его диаметра.

Зачастую реально существующую совокупность объектов (например, валики, изготовленные в течение января) можно мысленно дополнить любым количеством таких же однородных объектов (например, валики, изготовленные в тех же условиях в феврале, марте и т.д.). Такие совокупности объектов будем называть *генеральными совокупностями*.

Каждой генеральной совокупности соответствует случайная величина, определяемая изучаемым признаком объекта. В нашем примере – это диаметр валика. Так как понятия генеральной совокупности и соответствующей случайной величины связаны с наблюдениями (измерениями) в неизменных условиях, то для ее обозначения (по аналогии с курсом теории вероятностей) будем использовать прописные буквы латинского алфавита (например,  $X, Y$ ).

Часть отобранных объектов из генеральной совокупности называется *выборочной совокупностью*, или *выборкой*.

Результаты измерений изучаемого признака  $n$  объектов выборочной совокупности порождают  $n$  значений  $x_1, x_2, \dots, x_n$  случайной величины  $X$ . Число  $n$  называется *объемом выборки*.

Наряду с генеральной совокупностью  $X$  будем рассматривать  $n$  независимых случайных величин, обозначаемых той же буквой, что и генеральная совокупность, и имеющих точно такое же распределение, как генеральная совокупность. Итак,  $X_1, X_2, \dots, X_n$  –  $n$  независимых экземпляров  $X$ . Если  $F(x)$  – функция распределения генеральной совокупности  $X$ , то у каждой случайной величины  $X_i$  функция распределения также равна  $F(x)$ . Понятно, что

получить  $n$  значений случайной величины  $X$  – все равно что получить одно значение  $n$ -мерной случайной величины  $(X_1, X_2, \dots, X_n)$ . Поэтому каждую выборку  $x_1, x_2, \dots, x_n$  объема  $n$  мы можем рассматривать как одно значение  $n$ -мерной случайной величины  $(X_1, \dots, X_n)$ .

Поясним сказанное на примере. Пусть  $X$  – дискретная случайная величина, принимающая значения 1, 2, 3, 4, 5, 6, каждое с вероятностью  $p = \frac{1}{6}$ . Данную случайную величину, или в новой терминологии – генеральную совокупность, мы можем вообразить как урну, содержащую одинаковое количество шаров с номерами от 1 до 6. Производя выбор с возвращением трех шаров и записывая их номера, мы получим выборку объема 3 из генеральной совокупности  $X$ . Вообразим себе три урны того же содержания, т.е. три копии  $X_1, X_2, X_3$  урны  $X$ . Выберем из каждой урны по одному шару. Получим выборку  $x_1, x_2, x_3$  из генеральной совокупности  $X$ .

### 2.2. Свойства выборочной совокупности

Для того чтобы по отобранным значениям некоторого количественного показателя можно было достаточно уверенно судить обо всей совокупности, полученная выборка должна быть *репрезентативной (представительной)*, т.е. правильно отражать пропорции генеральной совокупности. Предположим, например, что вся совокупность состоит из равного большого количества белых и черных шаров, помещенных в ящик, на дне которого имеется отверстие. Если черные шары сосредоточены в нижней части ящика, а белые – в верхней, то, открывая некоторое небольшое количество раз заслонку в отверстии ящика, мы получим выборку только из черных шаров. На основании такого способа отбора шаров мы не сможем сделать правильных выводов о содержании всей совокупности шаров, т.е. такая выборка не будет репрезентативной. Выборка будет представительной лишь тогда, когда все объекты генеральной совокупности будут иметь *одинаковую вероятность попасть в выборку*. Для этого шары должны быть перемешаны. Другими словами, *репрезентативность выборки обеспечивается случайностью отбора объектов в выборку*.

Существует несколько способов отбора, обеспечивающих репрезентативность выборки.

Пусть небольшие по размеру объекты генеральной совокупности находятся, например, в ящике. Каждый раз после тщательного перемешивания (если оно не вызывает разрушения объектов) из ящиков наудачу берут один объект. Эту операцию повторяют до тех пор, пока не образуется выборка нужного объема. Очевидно, что такая техника отбора невозможна, если генеральная совокупность состоит из больших (по размерам) или хрупких объектов, например из мощных электромоторов. В этих случаях поступают следующим образом. Все объекты генеральной совокупности нумеруют и каждый номер записывают на отдельную карточку. После этого карточки с номерами тщательно перемешивают и из пачки карточек выбирают одну. Объект, номер которого совпал с номером выбранной карточки, включают в выборку. Номера объектов можно "отбирать" с помощью таблиц случайных чисел – это целесообразно при большом объеме генеральной совокупности.

Принципиально, что при отборе объектов в выборочную совокупность возможны два варианта:

1. Объект возвращается в генеральную совокупность. Выборочная совокупность, полученная таким образом, называется *случайной выборкой с возвратом* (или *повторной выборкой*).
2. Объект, включенный в выборку, не возвращается в генеральную совокупность. Образованная выборка называется *случайной выборкой без возврата* (или *бесповторной выборкой*).

Очевидно, что в повторной выборке возможна ситуация, когда один и тот же объект будет обследован несколько раз. Если объем генеральной совокупности велик, то различие между повторной и бесповторной выборками (которые составляют небольшую часть генеральной совокупности) незначительно и это практически не сказывается на окончательных результатах. В таких случаях, как правило, используют выборку без возврата. Если генеральная совокупность имеет не очень большой объем, то различие между указанными выборками будет существенным.

### 2.3. Вариационные ряды

После получения (тем или иным способом) выборочной совокупности все ее объекты обследуются по отношению к определенной случайной величине, т.е. обследуемому признаку объекта. В результате этого получают наблюдаемые данные, которые представляют собой множество чисел, расположенных в беспорядке. Анализ таких данных весьма затруднителен, и для изучения закономерностей полученные данные подвергаются определенной обработке.

♦ **Пример 2.1.** На телефонной станции проводились наблюдения над числом  $X$  неправильных соединений в минуту. Наблюдения в течение часа дали следующие 60 значений:

3; 1; 3; 1; 4; | 1; 2; 4; 0; 3; | 0; 2; 2; 0; 1; | 1; 4; 3; 1; 1;

4; 2; 2; 1; 1; | 2; 1; 0; 3; 4; | 1; 3; 2; 7; 2; | 0; 0; 1; 3; 3;

1; 2; 1; 2; 0; | 2; 3; 1; 2; 5; | 1; 2; 4; 2; 0; | 2; 3; 1; 2; 5. ●

Очевидно, что число  $X$  является дискретной случайной величиной, а полученные данные есть значения этой случайной величины. Анализ исходных данных в таком виде весьма затруднителен.

Простейшая операция – *ранжирование* опытных данных, результатом которого являются значения, расположенные в порядке *неубывания*. Если среди элементов встречаются одинаковые, то они объединяются в одну группу. Значение случайной величины, соответствующее отдельной группе сгруппированного ряда наблюдаемых данных, называется *вариантом*, а изменение этого значения – *варьированием*. Варианты будем обозначать строчными буквами с соответствующими порядковому номеру группы индексами  $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ , где  $m$  – число групп. При этом  $x^{(1)} < x^{(2)} < \dots < x^{(m)}$ .

Численность отдельной группы сгруппированного ряда данных называется *частотой*  $n_i$ , где  $i$  – индекс варианта, а отношение частоты данного варианта к общей сумме частот называется *частотностью* (или *относительной частотой*) и обозначается  $\omega_i$ ,  $i = 1, \dots, m$ , т.е.

$$\omega_i = \frac{n_i}{\sum_{i=1}^m n_i}, \quad (2.1)$$

при этом  $\sum_{i=1}^m n_i = n$ .

Дискретным вариационным рядом называется ранжированная совокупность вариантов  $x^{(i)}$  с соответствующими им частотами  $n_i$  или частностями  $\omega_i$ .

♦ **Пример 2.2.** Для данных примера 2.1 были выполнены операции ранжирования и группировки. В результате были получены семь значений случайной величины (варианты): 0; 1; 2; 3; 4; 5; 7. При этом значение 0 в этой группе встречается 8 раз, значение 1 – 17 раз, значение 2 – 16 раз, значение 3 – 10 раз, значение 4 – 6 раз, значение 5 – 2 раза, значение 7 – 1 раз. Вычисленные значения частот и частностей приведены в табл. 2.1.

Таблица 2.1

Индекс	$i$	1, 2, 3, 4, 5, 6, 7
Вариант	$x^{(i)}$	0, 1, 2, 3, 4, 5, 7
Частота	$n_i$	8, 17, 16, 10, 6, 2, 1
Частность	$\omega_i$	$\frac{8}{60}, \frac{17}{60}, \frac{16}{60}, \frac{10}{60}, \frac{6}{60}, \frac{2}{60}, \frac{1}{60}$

Таким образом, получен дискретный ряд:

$$0(8); 1(17); 2(16); 3(10); 4(6); 5(2); 7(1),$$

где в скобках указаны соответствующие частоты. В отличие от исходных данных (см. пример 2.1), этот ряд позволяет делать некоторые выводы о статистических закономерностях. ☉

Если среди  $n$  наблюдаемых значений  $x_i$  отсутствуют одинаковые значения, то  $m = n, n_i = 1$ , а дискретный вариационный ряд имеет вид

$$x^{(1)} < x^{(2)} < \dots < x^{(n-1)} < x^{(n)}.$$

Если число возможных значений дискретной случайной величины достаточно велико или наблюдаемая случайная величина является непрерывной, то строят *интервальный вариационный ряд*, под которым понимают упорядоченную совокупность интервалов варьирования значений случайной величины с соответствующими частотами или частностями попаданий в каждый из них значений случайной величины.

Как правило, частичные интервалы, на которые разбивается весь интервал варьирования, имеют одинаковую длину и представляются в виде

$$[z_i, z_i + h), \quad i = 1, 2, \dots, m, \quad (2.2)$$

где  $m$  – число интервалов.

Длину  $h$  следует выбирать так, чтобы построенный ряд не был громоздким, но в то же время позволял выявлять характерные изменения случайной величины.

Для вычисления  $h$  рекомендуется использовать следующую формулу:

$$h = \frac{x_{\max} - x_{\min}}{1 + 3.222 \lg n},$$

где  $x_{\max}, x_{\min}$  – наибольшее и наименьшее значения случайной величины. Если окажется, что  $h$  – дробное число, то за длину интервала следует принять либо ближайшую простую дробь, либо ближайшую целую величину. При этом необходимо выполнение условий:

$$z_1 \leq x_{\min}; \quad z_m + h \geq x_{\max}. \quad (2.3)$$

После нахождения частных интервалов определяется, сколько значений случайной величины попало в каждый конкретный интервал. При этом в интервал включают значения, большие или равные нижней границе и меньшие верхней границы.

♦ **Пример 2.3.** При изменении диаметра валика после шлифовки была получена следующая выборка (объемом  $n = 55$ ):



20.3	15.4	17.2	19.2	23.3	18.1	21.9
15.3	16.8	13.2	20.4	16.5	19.7	20.5
14.3	20.1	16.8	14.7	20.8	19.5	15.3
19.3	17.8	16.2	15.7	22.8	21.9	12.5
10.1	21.1	18.3	14.7	14.5	18.1	18.4
13.9	19.8	18.5	20.2	23.8	16.7	20.4
19.5	17.2	19.6	17.8	21.3	17.5	19.4
17.8	13.5	17.8	11.8	18.6	19.1	

Необходимо построить интервальный вариационный ряд, состоящий из семи интервалов.

*Решение.* Так как наибольшая варианта равна 23.8, а наименьшая 10.1, то вся выборка попадает в интервал (10,24). Мы расширили интервал (10.1,23.8) для удобства вычислений. Длина каждого частичного интервала равна  $\frac{24-10}{7} = 2$ . Получаем следующие семь интервалов:

[10,12);[12,14);[14,16);[16,18);[18,20);[20,22);[22,24),

а соответствующий интервальный вариационный ряд представлен в табл. 2.2.

Таблица 2.2

$X$	10–12	12–14	14–16	16–18	18–20	20–22	22–24
$\omega_i$	$\frac{2}{55}$	$\frac{4}{55}$	$\frac{8}{55}$	$\frac{12}{55}$	$\frac{15}{55}$	$\frac{11}{55}$	$\frac{3}{55}$

## 2.4. Выборочная функция распределения. Гистограмма

В теории вероятностей для характеристики распределения случайной величины  $X$  служит функция распределения

$$F(x) = P(X < x),$$

равная вероятности события  $\{X < x\}$ , где  $x$  – любое действительное число.

Одной из основных характеристик выборки является *выборочная (эмпирическая) функция распределения*

$$F_n^*(x) = \frac{n_x}{n}, \quad (2.4)$$

где  $n_x$  – количество элементов выборки, меньших чем  $x$ . Другими словами,  $F_n^*(x)$  есть относительная частота появления события  $A = \{X < x\}$  в  $n$  независимых испытаниях. Главное различие между  $F(x)$  и  $F_n^*(x)$  состоит в том, что  $F(x)$  определяет вероятность события  $A$ , а выборочная функция распределения  $F_n^*(x)$  – относительную частоту этого события.

Из определения (2.4) имеем следующие свойства функции  $F_n^*(x)$ :

1.  $0 \leq F_n^*(x) \leq 1$ .
2.  $F_n^*(x)$  – неубывающая функция.
3.  $F_n^*(-\infty) = 0$ ;  $F_n^*(\infty) = 1$ .

Напоминаем, что такими же свойствами обладает и функция распределения  $F(x)$  (вспомните эти свойства и сравните).

Функция  $F_n^*(x)$  является "ступенчатой", имеются разрывы в точках, которым соответствуют наблюдаемые значения вариантов. Величина скачка равна относительной частоте варианта.

Аналитически  $F_n^*(x)$  задается следующим соотношением:

$$F_n^*(x) = \begin{cases} 0 & \text{при } x \leq x^{(1)}; \\ \sum_{j=1}^{i-1} \omega_j & \text{при } x^{(i-1)} < x \leq x^{(i)}, \quad i = 1, 2, \dots, m; \\ 1 & \text{при } x > x^{(m)}, \end{cases} \quad (2.6)$$

где  $\omega_i$  – соответствующие относительные частоты, определяемые выражением (2.1);  $x^{(i)}$  – элементы вариационного ряда (варианты).

*Замечание.* В случае интервального вариационного ряда под  $x^{(i)}$  понимается середина  $i$ -го частичного интервала.

Перед вычислением  $F_n^*(x)$  полезно построить дискретный или интервальный вариационный ряд.

♦ **Пример 2.4.** Построить выборочную функцию распределения по наблюдаемым данным, приведенным в примере 2.1.

*Решение.* Используя соответствующий этим данным дискретный вариационный ряд (см. табл. 2.1), вычислим значения  $F_n^*(x)$  по формуле (2.6) и занесем их в табл. 2.3.

Таблица 2.3

$x$	$F_{60}^*(x)$
$x \leq 1$	0
$0 < x \leq 1$	$\omega_1 = \frac{8}{60}$
$1 < x \leq 2$	$\omega_1 + \omega_2 = \frac{25}{60}$
$2 < x \leq 3$	$\omega_1 + \omega_2 + \omega_3 = \frac{41}{60}$
$3 < x \leq 4$	$\omega_1 + \omega_2 + \omega_3 + \omega_4 = \frac{51}{60}$
$4 < x \leq 5$	$\omega_1 + \omega_2 + \omega_3 + \omega_4 + \omega_5 = \frac{57}{60}$
$5 < x \leq 7$	$\omega_1 + \omega_2 + \omega_3 + \omega_4 + \omega_5 + \omega_6 = \frac{59}{60}$
$x > 7$	$\omega_1 + \omega_2 + \omega_3 + \omega_4 + \omega_5 + \omega_6 + \omega_7 = \frac{60}{60} = 1$

Из графика  $F_{60}^*(x)$  (рис. 2.1) видно, что  $F_{60}^*(x)$  удовлетворяет свойствам (2.5). ☺

**Задача 2.1.** Построить выборочную функцию распределения по наблюдаемым данным, приведенным в примере 2.3.

Напомним, что  $F_n^*(x)$  равна относительной частоте появления события  $A = \{X < x\}$  и, следовательно, при любом значении  $x$  величина  $F_n^*(x)$  является случайной. Тогда конкретной выборке  $(x_1, x_2, \dots, x_n)$  объема  $n$  соответствует функция распределения  $F_n^*(x)$ , которая в силу своей случайности будет отличаться от

$F_n^*(x)$ , построенной по другой выборке из той же генеральной совокупности.

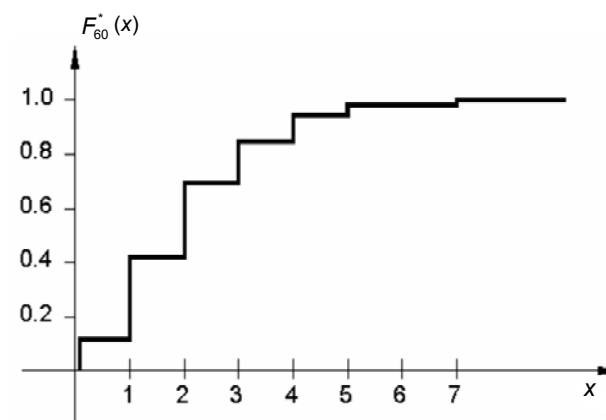


Рис. 2.1. График выборочной функции распределения (пример 2.4)

Возникает вопрос: зачем нужна такая характеристика, меняющаяся от выборки к выборке? Ответ получаем на основе следующих рассуждений.

По теореме Бернулли относительная частота появления события  $A$  в  $n$  независимых опытах сходится по вероятности к вероятности  $P(X < x)$  этого события при увеличении  $n$ . Следовательно, при больших объемах выборки выборочная функция распределения  $F_n^*(x)$  близка к теоретической функции  $F(x)$ . Точнее, имеет место следующая теорема.

**Теорема В.И. Гливенко.** Для любого действительного числа  $x$  и любого  $\varepsilon > 0$

$$\lim_{n \rightarrow \infty} P(|F_n^*(x) - F(x)| > \varepsilon) = 0.$$

Таким образом, по функции  $F_n^*(x)$  мы можем получить приближенно функцию  $F(x)$ , т.е. функция  $F_n^*(x)$  является оценкой  $F(x)$ .

В качестве оценки плотности распределения вероятности непрерывной случайной величины используют *гистограмму относительных частот*.

Гистограммой относительных частот называется система прямоугольников, каждый из которых основанием имеет  $i$ -й интервал интервального вариационного ряда; площадь, равную относительной частоте  $\omega_i$ , а высота  $y_i$  определяется по формуле

$$y_i = \frac{\omega_i}{h_i}, \quad i = 1, 2, \dots, m,$$

где  $h_i = z_{i+1} - z_i$  — длина  $i$ -го частичного интервала. Если длина частичных интервалов одинакова, то  $h_i = h$  (см. (2.2), (2.3)).

Очевидно, что *сумма площадей всех прямоугольников равна 1* (докажите это свойство).

Площадь прямоугольника  $\omega_i$  равна относительной частоте попадания элементов выборочной совокупности объема  $n$  на  $i$ -й интервал, т.е.

$$\omega_i = \omega_n^*(z_i \leq X < z_{i+1}).$$

С другой стороны, если  $y = p(x)$  — плотность вероятности случайной величины  $X$ , то вероятность

$$p_i = P(z_i \leq X < z_{i+1})$$

по теореме Бернулли близка при большом значении  $n$  к относительной частоте.

Поэтому значение  $\omega_i$  близко к

$$p_i = P(z_i \leq X < z_{i+1}) = \int_{z_i}^{z_{i+1}} p(x) dx. \quad (2.7)$$

Пусть  $y_i$  — высота  $i$ -го прямоугольника. По теореме о среднем интеграл, выражающий вероятность в формуле (2.7), можно записать в виде

$$p_i = \int_{z_i}^{z_{i+1}} p(x) dx = (z_{i+1} - z_i) \cdot p(u_i), \quad (2.8)$$

где  $u_i$  — некоторое число из промежутка  $[z_i, z_{i+1})$ . Так как  $\omega_i = (z_{i+1} - z_i)y_i$ , то значения  $y_i$  и  $p(u_i)$  близки друг к другу. Практически это означает, что график плотности распределения генеральной совокупности  $X$  проходит вблизи верхних границ прямоугольников, образующих гистограмму. Поэтому при больших объемах выборок и удачном выборе длины частичных интервалов гистограмма напоминает график плотности распределения  $p(x)$ .

♦ **Пример 2.5.** Построим гистограмму относительных частот выборочной совокупности из примера 2.3.

*Решение.* Используя интервальный вариационный ряд (см. табл. 2.2), находим высоты  $y_i$  по формуле  $y_i = \omega_i / 2$ . График построенной гистограммы приведен на рис. 2.2. Здесь же штриховой линией отмечен предполагаемый график неизвестной плотности  $p(x)$ . ☺

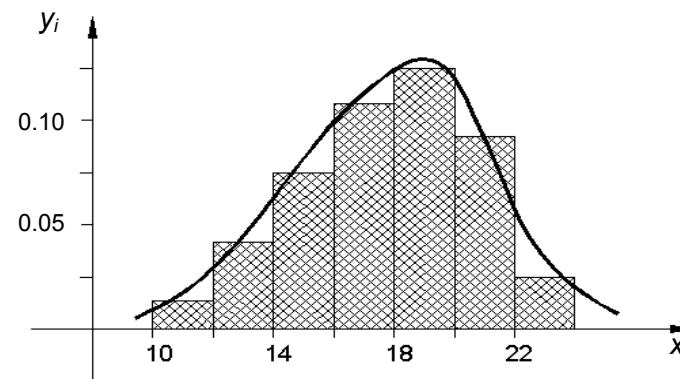


Рис. 2.2. График гистограммы частностей (пример 2.5)

## 2.5. Выборочное среднее и выборочная дисперсия

Рассмотренная выборочная функция распределения и гистограмма позволяют делать выводы о закономерностях исследуемого массового явления. Однако они неудобны для описания группиро-

вания и рассеивания наблюдаемых данных. Для этого используются так называемые *числовые характеристики выборочной совокупности*, из которых рассмотрим *выборочное среднее* и *выборочную дисперсию*.

Выборочным средним  $\bar{X}_g$  называется случайная величина, определенная формулой

$$\bar{X}_g = \frac{X_1 + X_2 + \dots + X_n}{n}. \quad (2.9)$$

Так как конкретная выборка  $x_1, \dots, x_n$  является реализацией значений случайных величин  $X_1, \dots, X_n$ , то среднее значение выборки

$$\bar{x}_g = \frac{x_1 + x_2 + \dots + x_n}{n} \quad (2.10)$$

является одной из реализаций случайной величины  $\bar{X}_g$ . Другими словами,  $\bar{x}_g$  есть одно из значений случайной величины  $\bar{X}_g$ .

Если данные представлены в виде вариационного ряда, то для вычисления выборочного среднего целесообразно применить одно из следующих соотношений:

- для дискретного вариационного ряда

$$\bar{x}_g = \frac{\sum_{i=1}^m x^{(i)} n_i}{\sum_{i=1}^m n_i} = \sum_{i=1}^m x^{(i)} \omega_i; \quad (2.11)$$

- для интервального вариационного ряда

$$\bar{x}_g = \frac{\sum_{i=1}^m z_i^* n_i}{\sum_{i=1}^m n_i} = \sum_{i=1}^m \omega_i z_i^*, \quad (2.12)$$

где  $\omega_i$  – частность (относительная частота), соответствующая  $i$ -й variante или  $i$ -му частичному интервалу;  $z_i^*$  – середина  $i$ -го частичного интервала, т.е.

$$z_i^* = \frac{(z_i + z_{i+1})}{2}, \quad i = 1, 2, \dots, m.$$

Сравним математическое ожидание дискретной случайной величины  $X$ , вычисляемое по формуле

$$M(X) = \sum_{i=1}^m x_i p_i, \quad (2.13)$$

и значение выборочного среднего, определяемое (2.11). Прежде всего, очевидна их внешняя схожесть. Однако в формуле (2.13)  $x_i$  – возможные значения случайной величины, а  $p_i$  – вероятности. В формуле (2.11)  $x^{(i)}$  – варианты случайной величины, полученные в результате наблюдений,  $\omega_i$  – их относительная частота. Далее, математическое ожидание не является случайной величиной, а выборочное среднее – случайная величина, значение которой меняется от выборки к выборке. Несмотря на это, как будет показано ниже, выборочное среднее при определенных условиях выступает как "хорошая" оценка математического ожидания.

♦ **Пример 2.6.** Вычислим значение выборочного среднего по выборке примера 2.1.

*Решение.* Используя дискретный вариационный ряд (см. табл. 2.1) и соотношение (2.1), имеем

$$\bar{x}_g = 0 \cdot \frac{8}{60} + 1 \cdot \frac{17}{60} + 2 \cdot \frac{16}{60} + 3 \cdot \frac{10}{60} + 4 \cdot \frac{6}{60} + 5 \cdot \frac{2}{60} + 7 \cdot \frac{1}{60} = 2.0. \quad \bullet$$

Так как значение выборочного среднего есть выборочный аналог математического ожидания, то имеет смысл ввести характеристику, которая бы оценивала величину рассеивания значений  $x_1, x_2, \dots, x_n$  относительно  $\bar{x}_g$ , а именно

$$d_g = \sum_{i=1}^n \frac{(x_i - \bar{x}_g)^2}{n}. \quad (2.14)$$

Число  $d_g$  является значением случайной величины

$$D_{\epsilon} = \sum_{i=1}^n \frac{(X_i - \bar{X}_{\epsilon})^2}{n}, \quad (2.15)$$

которую мы будем называть выборочной дисперсией.

Если данные представлены в виде вариационного ряда, то целесообразно для вычислений  $d_{\epsilon}$  вместо (2.14) использовать следующие соотношения:

- для дискретного вариационного ряда

$$d_{\epsilon} = \frac{\sum_{i=1}^m (x^{(i)} - \bar{x}_{\epsilon})^2 n_i}{n} = \sum_{i=1}^m (x^{(i)} - \bar{x}_{\epsilon})^2 \omega_i; \quad (2.16)$$

- для интервального вариационного ряда

$$d_{\epsilon} = \frac{\sum_{i=1}^m (z_i^* - \bar{x}_{\epsilon})^2 n_i}{n} = \sum_{i=1}^m (z_i^* - \bar{x}_{\epsilon})^2 \omega_i, \quad (2.17)$$

где  $\omega_i, z_i^*$  – те же, что и в формулах (2.11), (2.12).

Можно показать справедливость следующих выражений, являющихся аналогами (2.14), (2.16), (2.17) соответственно:

$$d_{\epsilon} = \frac{1}{n} \sum_{i=1}^n (x^{(i)})^2 - (\bar{x}_{\epsilon})^2; \quad (2.18)$$

$$d_{\epsilon} = \sum_{i=1}^m (x^{(i)})^2 \omega_i - (\bar{x}_{\epsilon})^2; \quad (2.19)$$

$$d_{\epsilon} = \sum_{i=1}^m (z_i^*)^2 \omega_i - (\bar{x}_{\epsilon})^2. \quad (2.20)$$

Приведенные соотношения (2.18)–(2.20) оказываются более удобными для программной реализации вычислений значения  $d_{\epsilon}$ .

Однако если генеральная дисперсия  $\sigma^2$  существенно меньше

квадрата математического ожидания, т.е.  $\sigma^2 \ll (M(x))^2$ , то из-за ошибок округления при машинном счете по этим формулам возможна ситуация  $d_{\epsilon} < 0$ . Тогда следует положить  $d_{\epsilon} = 0$ .

Сравним формулу (2.16) с формулой дисперсии дискретной случайной величины

$$D(X) = \sum_{i=1}^m (x_i - M(X))^2 p_i. \quad (2.21)$$

Различие между этими формулами состоит в том, что: а) величина  $D(X)$  не случайна,  $d_{\epsilon}$  – значение случайной величины, которое может меняться от выборки к выборке; б) в формуле (2.21)  $x_i$  – возможные значения случайной величины  $X$ ,  $p_i$  – их вероятности,  $M(X)$  – математическое ожидание. В формуле (2.16)  $x^{(i)}$  – варианты случайной величины,  $\omega_i$  – их относительные частоты, а  $\bar{x}_{\epsilon}$  – значения выборочного среднего. Несмотря на различия, между этими двумя формулами много общего. Во-первых, обе они являются мерой рассеивания. Во-вторых, кроме внешнего сходства формул, соответствующие дисперсии обладают схожими свойствами. В-третьих, как будет показано ниже, выборочная дисперсия при определенных условиях является хорошей оценкой для генеральной дисперсии  $D(X)$ .

♦ **Пример 2.7.** Необходимо вычислить значение выборочной дисперсии по выборке примера 2.1.

*Решение.* Воспользуемся формулой (2.19). Первоначально, используя дискретный вариационный ряд (см. табл. 2.1), вычислим

$$\sum_{i=1}^7 (x^{(i)})^2 \omega_i = 0 \cdot \frac{8}{60} + 1 \cdot \frac{17}{60} + 4 \cdot \frac{16}{60} + 9 \cdot \frac{10}{60} + 16 \cdot \frac{6}{60} + 25 \cdot \frac{2}{60} + 49 \cdot \frac{1}{60} = 6.09. \quad (2.22)$$

Так как значение  $\bar{x}_{\epsilon}$  было вычислено в примере 2.6 ( $\bar{x}_{\epsilon} = 2.0$ ), то

$$d_{\epsilon} = \sum_{i=1}^7 (x^{(i)})^2 \omega_i - (\bar{x}_{\epsilon})^2 = 6.09 - 4.0 = 2.09. \quad \bullet$$

## 2.6. Вычисление выборочных характеристик в Excel

**Вычисление частот.** Для вычисления частот  $n_i$  можно использовать **функцию ЧАСТОТА**, обращение к которой имеет вид:

$$=\text{ЧАСТОТА}(\text{массив\_данных}; \text{массив\_границ}),$$

где *массив\_данных* – адреса ячеек, для которых вычисляется частота  $n_i$ ; *массив\_границ* – адреса ячеек, в которых размещаются упорядоченные по возрастанию значения  $z_j$ ,  $j=1,2,\dots,m+1$ , где  $m$  – число интервалов.

При использовании этой функции необходимо помнить:

1. Функция ЧАСТОТА вводится как формула массива, т.е. предварительно выделяется интервал ячеек, в который будут помещены вычисленные частоты (число ячеек должно быть на 1 больше числа границ), затем вводится функция ЧАСТОТА с соответствующими аргументами, потом одновременно нажимаются клавиши [Ctrl] + [Shift] + [Enter].

2. Функция ЧАСТОТА игнорирует пустые ячейки и текстовые данные.

3. Если *массив\_границ* не содержит возрастающих значений границ и интервалов, то осуществляется автоматическое вычисление границ интервалов равной ширины, причем число интервалов равно корню квадратному из числа элементов *массива\_данных*.

**Результатом работы** является массив значений, определяемый по следующему правилу: первый элемент равен числу  $n_0$  элементов *массива\_данных* меньше  $z_1$ ; последний элемент равен числу  $n_{m+1}$  элементов *массива\_данных* больше  $z_{m+1}$ ; остальные элементы определяются как числа  $n_j$  элементов  $x_i$  *массива\_данных*, удовлетворяющих условию

$$z_j < x_i \leq z_{j+1}, \quad j=1,2,\dots,m.$$

Другими словами, кроме  $m$  значений частот  $n_j$ ,  $j=1,2,\dots,m$ , соответствующих  $m$  интервалам, вычисляются частоты  $n_0$  (число значений  $x_i$ , лежащих левее  $z_1$ ) и  $n_{m+1}$  (число значений  $x_i$ , лежащих правее  $z_{m+1}$ ).

♦ **Пример 2.8.** По выборке примера 2.3 вычислить частоты и частности для семи заданных интервалов [10,12); [12,14); [14,16); [16,18); [18,20); [20,22); [22,24), используя функцию ЧАСТОТА.

**Решение.** Первоначально, начиная с ячейки A3 (рис. 2.2), введем в столбец A 55 элементов выборки примера 2.3 (диапазон A3:A57). Затем, начиная с ячейки B3, введем границы заданных интервалов (см. рис. 2.2).

После подготовки этих данных выделяем ячейки C3:C11, вводим выражение

$$=\text{ЧАСТОТА}(A3:A57;B3:B10)$$

и нажимаем одновременно клавиши [Ctrl] + [Shift] + [Enter]. В ячейках C3:C11 появляется результат выполнения функции (см. рис. 2.2).

Для вычисления относительных частот  $\omega_j$  (частностей) необходимо частоты поделить на число элементов выборки. Эти вычисления реализованы в ячейках D3:D11 (см. рис. 2.2). Для контроля правильности вычисления частот и частностей в ячейках C12, D12 определены суммы (см. рис. 2.2):

$$\sum_{j=0}^{m+1=9} n_j = 55, \quad \sum_{j=0}^{m+1=9} \omega_j = 1. \quad \bullet$$

Для подсчета количества элементов выборки (т.е. объема выборки) использовалась **функция СЧЁТ**, обращение к которой имеет вид:

$$\text{СЧЁТ}(\text{массив\_данных}),$$

где *массив\_данных* – адреса ячеек или числовые константы.

**Результатом работы** является количество числовых величин в *массиве\_данных*. При этом игнорируются пустые ячейки, логические значения, тексты и значения ошибок.

	A	B	C	D	E	F	G
1							
2	Выборочные значения	Границы интервалов	Частоты	Частности			
3	20,3	10	0	0			
4	15,3	12	2	2/55			
5	14,3	14	4	4/55	=C4/CЧЕТ(A\$3:A\$57)		
6	19,3	16	8	8/55			
7	10,1	18	12	12/55			
8	13,9	20	15	3/11			
9	19,5	22	11	1/5			
10	17,8	24	3	3/55			
11	15,4		0	0			
12	16,8		55	1			
13	20,1				=СУММ(D3:D11)		
14	17,8				=СУММ(C3:C11)		
15	21,1						
16	19,8	={ЧАСТОТА(A3:A57:B3:B10)}					
17	17,2						
18	13,5						
19	17,2						
20	13,2						

Рис. 2.2. Фрагмент вычисления частот и частностей

данных выбрать режим *Гистограмма* и щелкнуть на кнопке ОК. Появится окно гистограммы, показанное на рис. 2.3. В окне задаются следующие параметры:

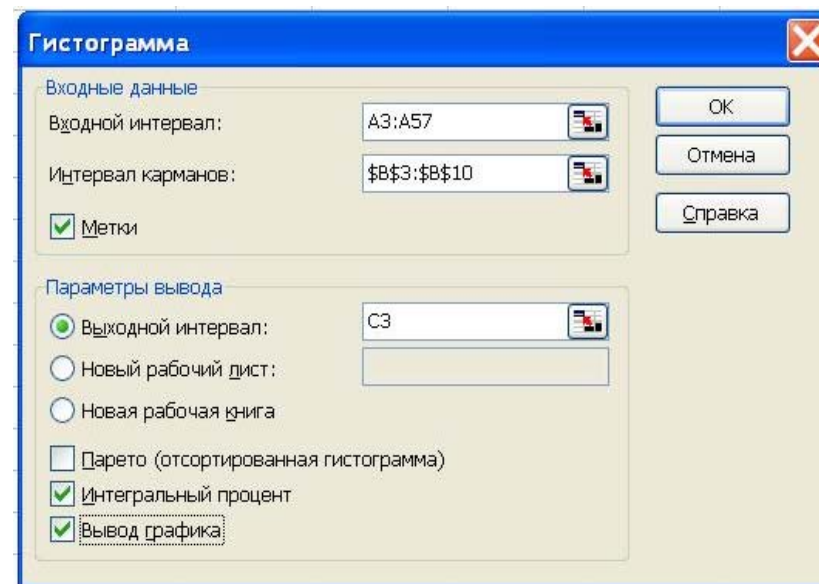


Рис. 2.3. Диалоговое окно режима *Гистограмма*

*Входной интервал:* – адреса ячеек, содержащие выборочные данные.

*Интервал карманов:* (необязательный параметр) – адреса ячеек, содержащие границы интервалов (кармана). Эти значения должны быть введены в возрастающем порядке.

*Метки* – флажок, включаемый, если первая строка во входных данных содержит заголовки. Если заголовки отсутствуют, то флажок следует выключить.

*Выходной интервал:* / *Новый рабочий лист:* / *Новая рабочая книга.* Включенный переключатель *Выходной интервал* требует ввода адреса верхней ячейки, начиная с которой будут размещаться вычисленные относительные частоты  $\omega_j$ . В положении переключателя *Новый рабочий лист:* открывается новый лист, в кото-

**Вычисление ненормированной гистограммы относительных частот.** Иногда в статистической (особенно зарубежной) литературе под гистограммой понимают систему прямоугольников, каждый из которых основанием имеет  $j$ -й интервал, а высота равна  $\omega_j$ . Очевидно, что сумма высот всех прямоугольников равна 1.

Заметим, что у ранее определенной гистограммы относительных частот сумма площадей прямоугольников равна 1, а высота прямоугольников равна  $u_j = \omega_j / h_j$ , где  $h_j$  – длина  $j$ -го интервала (т.е. выполнено нормирование). Поэтому первую гистограмму будем называть ненормированной гистограммой относительных частот.

Для построения ненормированной гистограммы необходимо обратиться к пункту **Сервис** строки меню Excel, а затем щелкнуть на команде *Анализ данных*, в появившемся окне диалога Анализ

ром начиная с ячейки A1 размещаются частности  $\omega_j$ . В положении переключателя *Новая рабочая книга* открывается новая книга, на первом листе которой начиная с ячейки A1 размещаются частности  $\omega_j$ .

*Парето (отсортированная гистограмма)* – устанавливается в активное состояние, чтобы представить  $\omega_j$  в порядке их убывания. Если параметр выключен, то  $\omega_j$  приводятся в порядке следования интервалов.

*Интегральный процент* – устанавливается в активное состояние для расчета выраженных в процентах накопленных относительных частот (процентный аналог значений выборочной функции распределения (2.6) при  $x_i = z_j$ ,  $j = 1, 2, \dots, m + 1$ ).

*Вывод графика* – устанавливается в активное состояние для автоматического создания встроенной диаграммы на листе, содержащем относительные частоты  $\omega_j$ .

При использовании режима *Гистограмма* модуля *Анализ данных* необходимо помнить:

1. Относительные частоты  $\omega_j$  вычисляются как количество элементов  $x_i$  выборки, удовлетворяющих условию

$$z_j < x_i \leq z_{j+1}.$$

2. Если границы интервалов не заданы, то автоматически будет создан набор интервалов с одинаковой длиной

$$h = \frac{x_{\max} - x_{\min}}{[k] - 1},$$

где  $[k]$  – целая часть величины  $k = 1 + 3.322 \cdot \lg n$ ,  $n$  – объем выборки.

♦ **Пример 2.9.** По выборке примера 2.3 построить ненормированную гистограмму относительных частот, используя режим *Гистограмма* модуля *Анализ данных*.

*Решение.* Первоначально, начиная с ячейки A3 (рис. 2.4), введем в столбец A 55 элементов выборки (диапазон A3:A57). Затем обратимся к пункту **Сервис**, команде *Анализ данных*, режиму *Гистограмма*. В появившемся диалоговом окне Гистограмма установим значения параметров, показанные на рис. 2.3, и после этого

щелкнем на кнопке ОК. В ячейках D4:D11 выводятся вычисленные значения  $\omega_j$ , а в ячейках E4:E11 – значения интегрального процента. В этом же листе строится диаграмма, на которой отображаются вычисленные характеристики. ☺

**Замечание 2.1.** Как правило, гистограммы изображаются в виде смежных прямоугольных областей. Поэтому столбики гистограммы на рис. 2.4 целесообразно расширить до соприкосновения друг с другом. Для этого необходимо щелкнуть мышью на диаграмме, далее на панель инструментов *Диаграмма*, раскрыть список инструментов и выбрать элемент *Ряд 'Частота'*, после чего щелкнуть на кнопке **Формат ряда**. В появившемся одноименном диалоговом окне необходимо активизировать закладку *Параметры* и в поле *Ширина зазора* установить значение 0.

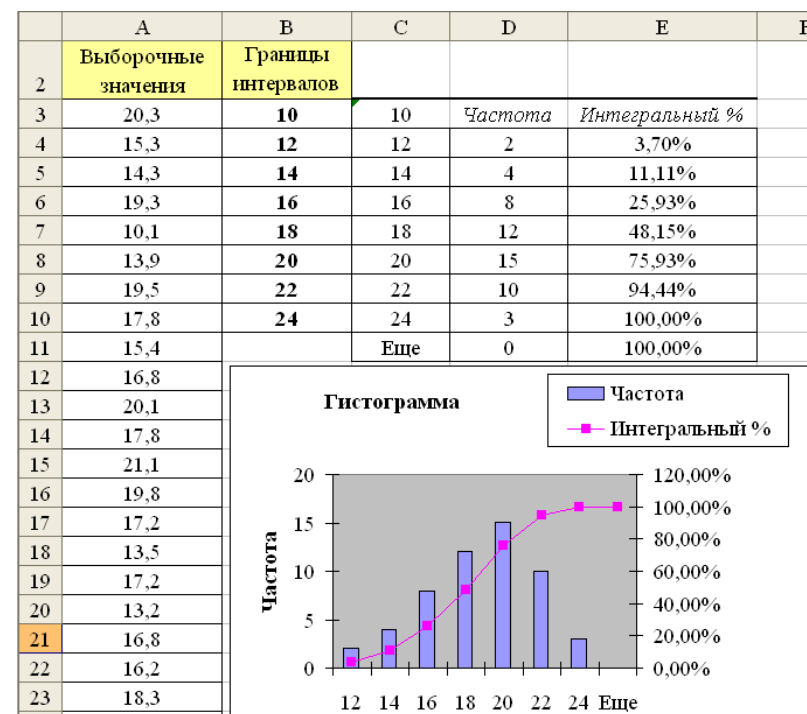


Рис. 2.4. Фрагмент построения гистограммы



На рис. 2.5 показана гистограмма, полученная из гистограммы (см. рис. 2.4) путем действий, описанных в замечании 2.1. ♦

**Замечание 2.2.** Ненормированная гистограмма относительных частот не может служить оценкой для плотности распределения случайной величины, из значений которой была сформирована выборка (особенно в случае неравных длин интервалов), из-за того, что сумма площадей прямоугольников  $\neq 1$ . В качестве такой оценки может рассматриваться гистограмма относительных частот. ♦

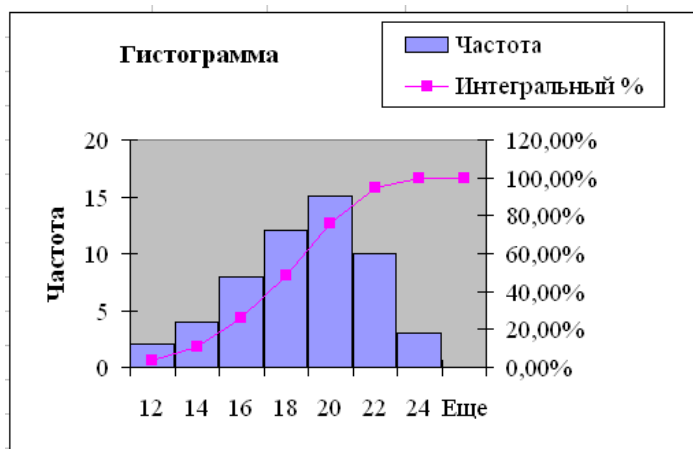


Рис. 2.5. График построенной гистограммы

**Вычисление гистограммы относительных частот.** Для вычисления такой гистограммы достаточно первоначально вычислить относительные частоты (частности), а затем полученные значения поделить на длину  $h_j$  соответствующего интервала, т.е. получить высоту соответствующего прямоугольника  $y_j = \omega_j / h_j$ . Для получения соприкасающихся прямоугольников выполнить операции, описанные в замечании 2.1 для соответствующего элемента.

♦ **Пример 2.10.** По выборке примера 2.3 построить гистограмму относительных частот.

**Решение.** Как и в примере 2.8, введем выборочные значения и, используя функцию ЧАСТОТА, вычислим частоты и частности. Затем, используя формулу  $y_j = \omega_j / h_j$ , где  $h_j = 2$ , вычислим высо-

ты прямоугольников (ячейки E3:E9) и середины интервалов (ячейки B3:B9). Для проверки правильности вычислений в ячейках D10, E10 определим суммы  $\sum \omega_j$ ,  $\sum y_j$ . Очевидно, что  $2 \cdot \sum y_j = 1$ .

В заключение по данным столбцов В, Е строим гистограмму (рис. 2.6). ☺

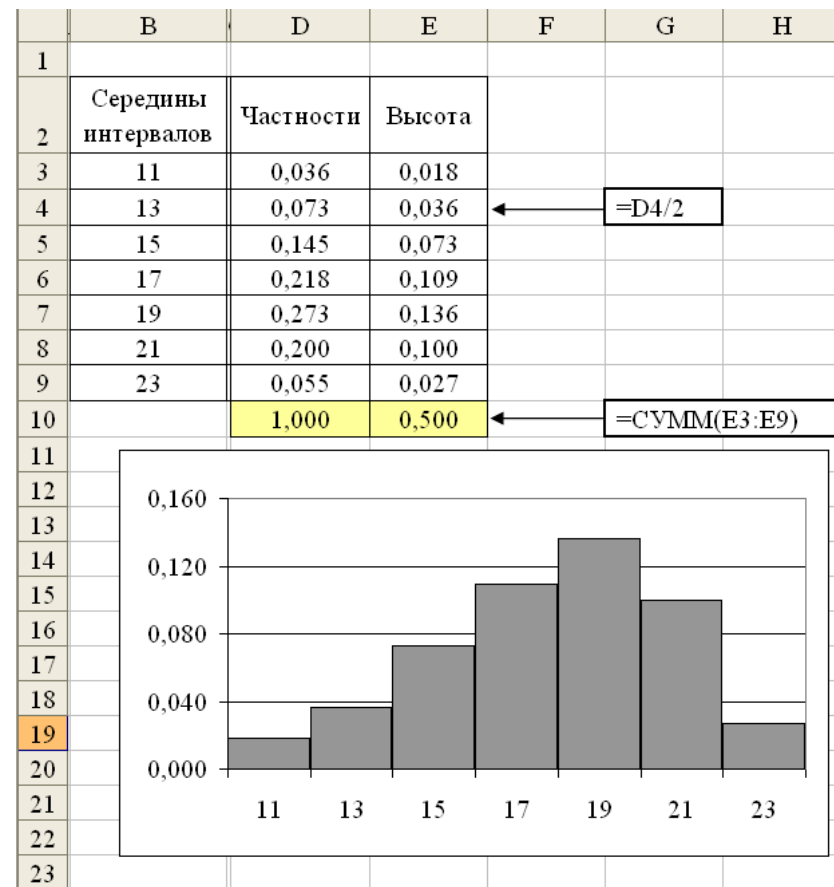


Рис. 2.6. Построение гистограммы относительных частот

**Вычисление выборочных среднего и дисперсии.** Для вычисления выборочного среднего (2.10) используется **функция СРЗНАЧ**, обращение к которой имеет вид:

$$=СРЗНАЧ(арг1; арг2; ...; арг30),$$

где  $арг1; арг2; ...; арг30$  – числа или адреса ячеек, содержащих числовые данные. Если ячейка содержит текстовые, логические значения или ячейка пуста, то такие ячейки игнорируются при подсчете среднего значения по формуле

$$\bar{x}_b = \frac{\sum_{i=1}^n x_i}{n}.$$

Здесь и в дальнейшем запись  $арг1; арг2; ...; арг30$  означает наличие от 1 до 30 аргументов функции Excel.

Для вычисления выборочной дисперсии (2.14) используется **функция ДИСПР**, обращение к которой имеет вид:

$$=ДИСПР(арг1; арг2; ...; арг30),$$

где  $арг1; арг2; ...; арг30$  – числа или адреса ячеек, содержащих числовые данные. Ячейки, содержащие текстовые, логические данные или пустые, при вычислении выборочной дисперсии игнорируются.

Для вычисления суммы квадратов отклонений

$$\sum_{i=1}^n (x_i - \bar{x}_b)^2$$

используется **функция КВАДРОТКЛ**, обращение к которой имеет вид:

$$=КВАДРОТКЛ(арг1; арг2; ...; арг30),$$

где  $арг1, арг2, ..., арг30$  – числа или адреса ячеек, содержащих числовые данные.

♦ **Пример 2.11.** По выборке примера 2.3 вычислить выборочное среднее  $\bar{x}_b$  и выборочную дисперсию  $\bar{d}_b$  двумя способами:

**Способ 1.** Программируя в ячейках Excel необходимые вычисления.

**Способ 2.** Используя функции Excel СРЗНАЧ, ДИСПР.

**Решение.** Первоначально, начиная с ячейки А3, введем в столбец А 55 элементов выборки (диапазон А3:А57). Запрограммируем выражения (2.10), (2.14), используя функции СУММ,

КВАДРОТКЛ, аргументами, указанными на рис. 2.7. Затем вычислим характеристики (2.10), (2.14) с использованием статистических функций СРЗНАЧ, ДИСПР (см. рис. 2.7). Как и следовало ожидать, результаты вычислений двумя способами совпали. ☺

	А	В	С	Д	Е
1					
2	Выборочные значения				
3	20,3	Программирование			
4	15,3				
5	14,3	17,907			
6	19,3		=СУММ(А3:А57)/55		
7	10,1		=КВАДРОТКЛ(А3:А57)/55		
8	13,9	8,601			
9	19,5				
10	17,8				
11	15,4	Стандартные функции Excel			
12	16,8				
13	20,1		=СРЗНАЧ(А3:А57)		
14	17,8	17,907			
15	21,1	8,601			
16	19,8		=ДИСПР(А3:А57)		
17	17,2				
18	13,5				

Рис. 2.7. Вычисление выборочных среднего и дисперсии

**Задание 2.1.** По выборочным данным ( $n = 60$ ) примера 2.1 построить гистограмму относительных частот. Длину интервала определить по формуле

$$h = \frac{x_{\max} - x_{\min}}{(1 + 3.322 \cdot \lg n)}.$$

**Рекомендация.** При выполнении задания использовать пример 2.10. ♥

**Задание 2.2.** По выборочным данным ( $n = 60$ ) примера 2.1 построить ненормированную гистограмму относительных частот, используя режим *Гистограмма*.

*Рекомендация.* При выполнении задания использовать пример 2.9. ♥

**Задание 2.3.** По выборочным данным ( $n = 60$ ) примера 2.1 вычислить выборочные среднее и дисперсию, используя стандартные функции Excel.

*Рекомендация.* При выполнении задания использовать пример 2.11. ♥

Кроме приведенных функций при вычислении выборочных характеристик могут быть полезными следующие функции:

**Функция МАКС** вычисляет максимальное значение из заданных аргументов. Обращение к ней имеет вид:

$$=\text{МАКС}(\text{arg1}; \text{arg2}; \dots; \text{arg30}),$$

где  $\text{arg1}; \text{arg2}; \dots; \text{arg30}$  – числовые константы или адреса ячеек, содержащих числовые величины.

**Функция МИН** вычисляет минимальное значение из заданных аргументов. Обращение к ней имеет вид:

$$=\text{МИН}(\text{arg1}; \text{arg2}; \dots; \text{arg30}),$$

где  $\text{arg1}; \text{arg2}; \dots; \text{arg30}$  – числовые константы или адреса ячеек, содержащих числовые величины.

### 3. ТОЧЕЧНЫЕ ОЦЕНКИ НЕИЗВЕСТНЫХ ПАРАМЕТРОВ

#### 3.1. Определение и свойства точечной оценки

Большинство случайных величин, рассмотренных в курсе теории вероятностей, имели распределения, зависящие от одного или нескольких параметров. Так, биномиальное распределение зависит от параметров  $p$  и  $n$ , нормальное – от параметров  $a$  и  $\sigma$ , распределение Пуассона – от параметра  $\lambda$  и т.п. Одной из основных задач математической статистики (см. главу 1) является оценивание этих параметров по наблюдаемым данным, т.е. по выборочной совокупности. В главе 2 были рассмотрены выборочные среднее и дисперсия, которые интерпретировались как приближенные значения неизвестных значений математического ожидания и дисперсии изучаемой случайной величины  $X$ , т.е. являлись оценками этих неизвестных характеристик.

Выборочная характеристика, используемая в качестве приближенного значения неизвестного параметра генеральной совокупности, называется *точечной оценкой* этого параметра. В этом определении слово "точечная" означает, что значение оценки представляет собой число или точку на числовой оси.

Обозначим через  $\theta$  некоторый неизвестный параметр генеральной совокупности, а через  $\theta_n^*$  – точечную оценку этого параметра. Оценка  $\theta_n^*$  есть функция  $\varphi(X_1, X_2, \dots, X_n)$  от  $n$  независимых экземпляров  $X_1, X_2, \dots, X_n$  генеральной совокупности, где  $n$  – объем выборки (см. п. 2.1). Поэтому оценка  $\theta_n^*$ , как функция случайных величин, также является случайной, и свойства  $\theta_n^*$  можно исследовать с использованием понятий теории вероятностей.

В общем случае точечная оценка  $\theta_n^*$  не связана с оцениваемым параметром  $\theta$ . Поэтому естественно потребовать, чтобы  $\theta_n^*$  была близка к  $\theta$ . Это требование формулируется в терминах несмещенности, состоятельности и эффективности.

Оценка  $\theta_n^*$  параметра  $\theta$  называется *несмещенной*, если для любого фиксированного объема выборки  $n$  математическое ожидание оценки равно оцениваемому параметру, т.е.

$$M(\theta_n^*) = \theta. \quad (3.1)$$

Поясним смысл этого равенства следующим примером. Имеются два алгоритма вычисления оценок для параметра  $\theta$ . Значения оценок, построенных первым алгоритмом по различным выборкам объема  $n$  генеральной совокупности, приведены на рис. 3.1,а, а с использованием второго алгоритма – на рис. 3.1,б. Видим, что среднее значение оценок на рис. 3.1,а совпадает с  $\theta$ , и, естественно, такие оценки предпочтительнее по сравнению с оценками на рис. 3.1,б, которые концентрируются слева от значения  $\theta$  и для которых  $M(\theta_n^*) < \theta$ , т.е. эти оценки являются смещенными.

Оценка  $\theta_n^*$  называется *состоятельной*, если

$$\theta_n^* \xrightarrow{P} \theta,$$

т.е. для любого  $\varepsilon > 0$  при  $n \rightarrow \infty$

$$P\left(\left|\theta_n^* - \theta\right| < \varepsilon\right) \rightarrow 1. \quad (3.2)$$

Поясним смысл этого предельного соотношения. Пусть  $\varepsilon$  – очень малое положительное число. Тогда (3.2) означает, что чем больше число наблюдений  $n$ , тем больше уверенность (вероятность) в незначительном отклонении  $\theta_n^*$  от неизвестного параметра  $\theta$ . Очевидно, что "хорошая" оценка должна быть состоятельной, иначе она не имеет практического смысла, так как увеличение объема исходной информации не будет приближать нас к "истинному" значению  $\theta$ .

Предположим, что имеются две состоятельные и несмещенные оценки

$$\theta_n^{*(1)} = \varphi_1(x_1, \dots, x_n); \quad \theta_n^{*(2)} = \varphi_2(x_1, \dots, x_n) \quad (3.3)$$

одного и того же параметра  $\theta$ . Как из двух этих оценок выбрать лучшую? Каждая из них является случайной величиной, и мы не можем предсказать индивидуальное значение оценки в каждом частном случае. Однако, рассматривая в качестве меры концентрации распределения оценки  $\theta_n^*$  около значения параметра  $\theta$  величину  $M(\theta_n^* - \theta)^2$ , мы можем теперь точно охарактеризовать сравнительную эффективность оценок  $\theta_n^{*(1)}$  и  $\theta_n^{*(2)}$ . В качестве меры эффективности принимается отношение

$$e = \frac{M(\theta_n^{*(1)} - \theta)^2}{M(\theta_n^{*(2)} - \theta)^2}. \quad (3.4)$$

Если  $e > 1$ , то оценка  $\theta_n^{*(2)}$  более эффективна, чем  $\theta_n^{*(1)}$ . В случае несмещенных оценок  $M(\theta_n^{*(1)}) = \theta$ ,  $M(\theta_n^{*(2)}) = \theta$ , и поэтому

$$e = \frac{D(\theta_n^{*(1)})}{D(\theta_n^{*(2)})}, \quad (3.5)$$

где  $D(\theta_n^*)$  – дисперсия оценки  $\theta_n^*$ .

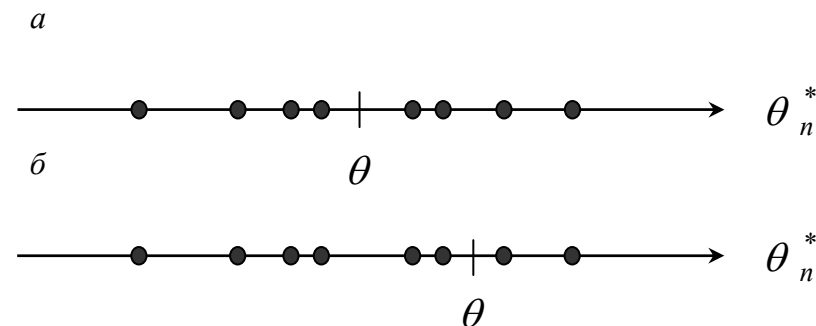


Рис. 3.1. К определению несмещенной оценки

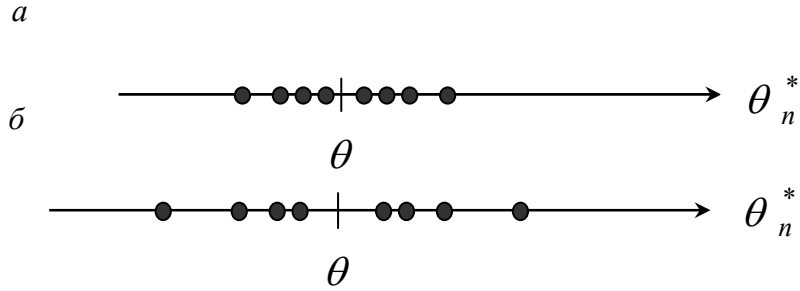


Рис. 3.2. К определению эффективной оценки

Таким образом, несмещенная оценка  $\theta_n^*$  параметра  $\theta$  называется *несмещенной эффективной*, если она среди всех других несмещенных оценок того же параметра обладает *наименьшей дисперсией*.

Приведенная на рис. 3.2,а оценка  $\theta_n^*$  является более эффективной по сравнению с оценкой, значения которой нанесены на рис. 3.2,б (почему?).

Как же выяснить, является ли несмещенная оценка эффективной? Очевидно, для этого необходимо сравнить дисперсию этой оценки с минимальной дисперсией.

Для широкого класса оценок *неравенство Рао–Крамера* указывает точную нижнюю границу для дисперсий различных оценок одного и того же параметра. Если существует оценка, дисперсия которой в точности равна этой нижней границе, то она называется *эффективной оценкой*. Оценка, имеющая наименьшую дисперсию среди оценок данного класса, называется *эффективной в данном классе оценок*. Поясним понятие эффективной оценки несколькими примерами.

Предположим, что генеральная совокупность распределена по нормальному закону с параметрами  $a$  и  $\sigma$ , причем  $a$  – математическое ожидание, подлежащее оценке, а  $\sigma^2$  – известная дисперсия. Оказывается, что для любой несмещенной регулярной оценки  $a^*$  имеет место неравенство

$$D(a^*) \geq \frac{\sigma^2}{n}, \quad (3.6)$$

где  $n$  – объем выборки, по которой производится оценивание. Если в качестве  $a^*$  принять  $\bar{X}_n$ , то дисперсия этой оценки, как будет показано ниже, равна  $\frac{\sigma^2}{n}$ , т.е.  $\bar{X}_n$  – эффективная оценка параметра  $a$ , так как для нее достигается нижняя грань в неравенстве (3.6).

Рассмотрим на примере понятие *эффективной в данном классе оценки*. Предположим, что один и тот же предмет, истинная величина которого равна  $l$ , измеряется  $n$  раз различными приборами, имеющими различную точность. Пусть  $X_i$  – результаты  $i$ -го измерения. Тогда

$$M(X_i) = l, \quad D(X_i) = \sigma_i^2,$$

если считать, что измерения проводятся без систематических ошибок. Дисперсия  $\sigma_i^2$  характеризует точность измерений. Для оценки истинного значения параметра  $l$  рассмотрим класс линейных оценок, т.е. оценок вида

$$l^* = c_1 X_1 + \dots + c_n X_n,$$

где  $c_1, \dots, c_n$  – некоторые неизвестные константы. Из всех несмещенных оценок данного класса нужно выбрать ту, которая имеет наименьшую дисперсию.

Из несмещенности оценок получим

$$M(l^*) = M\left(\sum_{i=1}^n c_i X_i\right) = \sum_{i=1}^n c_i M(X_i) = l \sum_{i=1}^n c_i.$$

Значит,

$$\sum_{i=1}^n c_i = 1. \quad (3.7)$$

Пользуясь свойствами дисперсии и независимостью проведенных измерений, получим

$$D(l^*) = \sum_{i=1}^n c_i^2 \sigma_i^2.$$

Числа  $c_1, \dots, c_n$  должны удовлетворять условию (3.7) и обеспечивать минимум функции

$$F(c_1, \dots, c_n) = \sum_{i=1}^n c_i^2 \sigma_i^2.$$

Мы получим задачу на условный экстремум, которую можно решить с помощью функции Лагранжа:

$$L(c_1, \dots, c_n) = F(c_1, \dots, c_n) - \lambda \left( \sum_{i=1}^n c_i - 1 \right).$$

Найдем критические точки функции Лагранжа:

$$\frac{\partial L}{\partial c_i} = 2c_i \sigma_i^2 - \lambda = 0, \quad i = 1, \dots, n;$$

$$\sum_{i=1}^n c_i - 1 = 0.$$

Отсюда находим значение

$$c_i = \frac{\frac{1}{\sigma_i^2}}{\sum_{i=1}^n \frac{1}{\sigma_i^2}}, \quad i = 1, \dots, n. \quad (3.8)$$

Полученный результат имеет простой физический смысл: чем меньше точность данного прибора, тем с меньшим значением коэффициента его результат должен входить в оценку.

Заметим, что если все приборы имеют одинаковую точность, т.е.  $\sigma_1^2 = \dots = \sigma_n^2$ , то  $c_i = 1/n$  и в качестве оценки получим  $l^* = \bar{X}_g$ .

### 3.2. Точечная оценка математического ожидания

Математическое ожидание  $M(X)$  генеральной совокупности  $X$  назовем *генеральной средней*  $\bar{x}_2$ , т.е.

$$\bar{x}_2 = M(X).$$

**Теорема 3.1.** Выборочное среднее  $\bar{X}_g$  есть *состоятельная и несмещенная* оценка генеральной средней  $\bar{x}_2$ .

*Доказательство.* Вначале покажем, что  $\bar{X}_g$  есть состоятельная оценка для  $\bar{x}_2$ , т.е.

$$\frac{X_1 + X_2 + \dots + X_n}{n} \xrightarrow{p} \bar{x}_2.$$

По следствию из теоремы Чебышева для одинаково распределенных случайных величин имеем

$$\frac{X_1 + X_2 + \dots + X_n}{n} \xrightarrow{p} M(X).$$

Так как  $M(X) = \bar{x}_2$ , то, используя свойства математического ожидания, получим

$$\begin{aligned} M(\bar{X}_g) &= M\left(\frac{X_1 + \dots + X_n}{n}\right) = \frac{M(X_1) + \dots + M(X_n)}{n} = \\ &= \frac{nM(X)}{n} = \bar{x}_2. \end{aligned}$$

Теорема доказана.

**Теорема 3.2.** Пусть случайная величина  $X$  имеет нормальное распределение  $N(a, \sigma)$ , где  $a$  – математическое ожидание,  $\sigma^2$  – дисперсия случайной величины  $X$ . Тогда выборочное среднее  $\bar{X}_g$  является *эффективной несмещенной* оценкой для  $\bar{x}_2$ .

*Доказательство.* Необходимо показать, что дисперсия  $D(\bar{X}_g)$  совпадает с минимальной дисперсией, равной в случае нормального распределения  $\sigma^2/n$ , а ее математическое ожидание  $M(\bar{X}_g)$  равно  $\bar{x}_2$ .

Найдем дисперсию  $D(\bar{X}_g)$ :

$$D(\bar{X}_g) = D\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n D(X_i) = \frac{nD(X)}{n^2} = \frac{\sigma^2}{n}. \quad (3.9)$$

Мы проверили при доказательстве теоремы 3.1, что  $M(\bar{X}_g) = \bar{x}_g$ . Так как дисперсия  $D(\bar{X}_g)$  равна минимальному значению, то выборочное среднее  $\bar{X}_g$  является эффективной несмещенной оценкой.

Теорема доказана.

Таким образом, показано, что выборочное среднее  $\bar{X}_g$  имеет все три свойства "хорошей" оценки. Этим и объясняется ее широкое использование в качестве оценки математического ожидания генеральной совокупности.

Напомним, что по конкретной выборке  $x_1, \dots, x_n$  вычисляется (см. (2.10)–(2.12)) "конкретное" значение  $\bar{x}_g$ , являющееся одним из множества возможных значений случайной величины  $\bar{X}_g$ .

### 3.3. Точечные оценки дисперсии

Дисперсию  $D(X)$  генеральной совокупности  $X$  будем называть генеральной дисперсией  $D_g$ , т.е.

$$D_g = D(X). \quad (3.10)$$

**Теорема 3.3.** Выборочная дисперсия  $D_g$  является состоятельной, но смещенной оценкой генеральной дисперсии  $D_g$ .

*Доказательство.* Получим сначала формулу для вычисления  $D_g$ . Согласно определению

$$D_g = \frac{\sum_{i=1}^n (X_i - \bar{X}_g)^2}{n}.$$

С другой стороны,

$$\begin{aligned} \sum_{i=1}^n (X_i - \bar{X}_g)^2 &= \sum_{i=1}^n (X_i^2 - 2\bar{X}_g X_i + \bar{X}_g^2) = \\ &= \sum_{i=1}^n X_i^2 - 2n\bar{X}_g^2 + n\bar{X}_g^2 = \sum_{i=1}^n X_i^2 - n\bar{X}_g^2. \end{aligned}$$

Тогда из определения дисперсии следует

$$D_g = \frac{\sum_{i=1}^n X_i^2 - n\bar{X}_g^2}{n} = \frac{\sum_{i=1}^n X_i^2}{n} - \bar{X}_g^2.$$

Воспользовавшись теперь следствием из теоремы Чебышева для одинаково распределенных случайных величин  $X_i^2$  и свойствами предела по вероятности, получаем

$$\frac{\sum_{i=1}^n X_i^2}{n} \xrightarrow{p} M(X_i^2) = M(X^2);$$

$$\bar{X}_g \xrightarrow{p} M(X)$$

и, значит,

$$D_g \xrightarrow{p} M(X^2) - M^2(X) = D(X) = D_g.$$

Следовательно, выборочная дисперсия  $D_g$  является состоятельной оценкой для генеральной дисперсии. Вычислим математическое ожидание  $D_g$  и убедимся, что  $M(D_g) \neq D_g$ . Имеем

$$\begin{aligned} M(D_g) &= M\left(\frac{\sum_{i=1}^n X_i^2}{n} - \bar{X}_g^2\right) = M\left(\frac{\sum_{i=1}^n X_i^2}{n}\right) - M(\bar{X}_g^2) = \\ &= M\left(\frac{\sum_{i=1}^n X_i^2}{n}\right) - M\left(\frac{X_1 + \dots + X_n}{n}\right)^2 = \\ &= M(X^2) - M\left(\frac{X_1^2 + X_2^2 + \dots + X_n^2 + \sum_{i \neq j} X_i X_j}{n^2}\right), \end{aligned}$$

где  $\sum_{i \neq j} X_i X_j$  означает сумму произведений величин  $X_i$  и  $X_j$  для всех значений  $i$  и  $j$  от 1 до  $n$ , но не равных между собой. Так как  $X_i$  и  $X_j$  независимы при  $i \neq j$ , то

$$M(X_i X_j) = M(X_i)M(X_j).$$

Поэтому, продолжая вычисления  $M(D_e)$ , получаем

$$\begin{aligned} M(D_e) &= M(X^2) - \frac{M(X_1^2) + \dots + M(X_n^2) + \sum_{i \neq j} M(X_i)M(X_j)}{n^2} = \\ &= M(X^2) - \frac{nM(X^2) + n(n-1)M^2(X)}{n^2} = \\ &= \frac{n-1}{n} [M(X^2) - M^2(X)] = \frac{n-1}{n} D_e. \end{aligned}$$

Множитель  $n(n-1)$  объясняется тем, что по правилу произведения количество различных пар  $(i, j)$  при  $1 \leq i \neq j \leq n$  равно  $n(n-1)$ . Итак, мы получили, что

$$M(D_e) = \frac{n-1}{n} D_e, \quad (3.11)$$

следовательно,  $D_e$  – смещенная оценка для генеральной дисперсии.

Теорема доказана.

Полученная формула (3.11) для вычисления математического ожидания выборочной дисперсии позволяет указать состоятельную и несмещенную оценку для генеральной дисперсии. Для этого рассмотрим случайную величину

$$S^2 = \frac{n}{n-1} D_e, \quad (3.12)$$

называемую *исправленной дисперсией*. Понятно, что

$$S^2 \xrightarrow{p} D_e,$$

так как  $\frac{n}{n-1} \rightarrow 1$  при  $n \rightarrow \infty$ . С другой стороны,

$$M(S^2) = M\left(\frac{n}{n-1} D_e\right) = \frac{n}{n-1} M(D_e) = \frac{n}{n-1} \cdot \frac{n-1}{n} D_e = D_e.$$

Тем самым доказана

**Теорема 3.4.** Исправленная дисперсия  $S^2$  является состоятельной и несмещенной оценкой для генеральной дисперсии  $D_e$ .

Заметим, что для выборок большого объема множитель  $\frac{n}{n-1}$

близок к 1, поэтому случайные величины  $S^2$  и  $D_e$  мало отличаются друг от друга. Однако для выборок малого объема это отличие может быть существенным.

Возникает вопрос: будет ли несмещенная оценка  $S^2$  эффективной?

Предположим, что случайная величина  $X$  подчиняется нормальному распределению  $N(a, \sigma)$ , а величины  $X_1, X_2, \dots, X_n$ , как обычно, –  $n$  независимых экземпляров независимой величины  $X$ . Тогда минимальная дисперсия несмещенной оценки для дисперсий равна

$$D_{\min} = \frac{2\sigma^4}{n}. \quad (3.13)$$

В п. 4.1 будет показано, что величина  $S^2$  представима в виде

$$S^2 = \frac{\sigma^2}{n-1} \chi_{n-1}^2, \quad (3.14)$$

где  $\chi_{n-1}^2$  – случайная величина, имеющая  $\chi^2$ -распределение с  $n-1$  степенями свободы. Поэтому

$$D(S^2) = \frac{\sigma^4}{(n-1)^2} D(\chi_{n-1}^2) = \frac{2\sigma^4}{n-1}, \quad (3.15)$$

из этого следует

$$D(S^2) = \frac{n}{n-1} D_{\min}. \quad (3.16)$$

Следовательно,  $S^2$ , будучи несмещенной оценкой дисперсии  $D(X)$ , не является эффективной оценкой. Однако при достаточно больших  $n$  увеличение  $D(S^2)$  по сравнению с  $D_{\min}$  пренебрежимо мало.



Заметим, что несмещенная эффективная оценка дисперсии  $D(X)$  нормально распределенной величины  $X = N(a, \sigma)$  имеет вид:

$$S_0^2 = \frac{1}{n} \sum_{i=1}^n (X_i - a)^2.$$

Однако в эту формулу входит математическое ожидание  $a$ , которое, как правило, заранее неизвестно.

### 3.4. Точечная оценка вероятности события

Обозначим через  $p(A)$  неизвестную вероятность события  $A$  в одном испытании. Для оценивания  $p(A)$  проведем  $n$  независимых испытаний, в которых событие  $A$  произошло  $m$  раз. Тогда случайная величина

$$p^* = \frac{m}{n} \quad (3.17)$$

является частностью (относительной частотой) события  $A$ . Свойства этой точечной оценки определяет

**Теорема 3.5.** Относительная частота  $p^* = m/n$  появления события  $A$  в  $n$  испытаниях есть состоятельная, несмещенная и эффективная оценка вероятности  $p(A)$ .

*Доказательство.* Состоятельность оценки  $p^*$  вытекает из теоремы Бернулли, согласно которой для любого  $\varepsilon > 0$  выполняется неравенство

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{m}{n} - P(A)\right| < \varepsilon\right) = 1, \quad (3.18)$$

или в других обозначениях:

$$\frac{m}{n} \xrightarrow{p} p(A).$$

Для доказательства несмещенности этой оценки зафиксируем число испытаний  $n$ . Найдем математическое ожидание частности  $m/n$ , имея в виду, что в условиях испытаний Бернулли величина  $m$

имеет биномиальный закон распределения с характеристиками  $M(m) = np$ ,  $D(m) = np(1-p)$ . Имеем

$$M\left(\frac{m}{n}\right) = \frac{1}{n} M(m) = \frac{1}{n} np = p(A).$$

Следовательно,  $p^* = m/n$  является несмещенной оценкой вероятности  $p(A)$ .

Для доказательства эффективности укажем, что минимум среди дисперсий различных несмещенных оценок вероятности  $p(A)$  равен

$$D_{\min} = \frac{p(1-p)}{n}. \quad (3.19)$$

Определим дисперсию оценки  $p^*$ :

$$D(p^*) = D\left(\frac{m}{n}\right) = \frac{1}{n^2} D(m) = \frac{np(1-p)}{n^2} = \frac{p(1-p)}{n}.$$

Так как  $D(p^*)$  совпадает с минимальной дисперсией  $D_{\min}$ , то частность  $p^*$ , будучи несмещенной оценкой, является также и эффективной.

Теорема доказана.

### 3.5. Метод максимального правдоподобия

В предыдущих пунктах были рассмотрены различные точечные оценки, являющиеся некоторыми функциями от результатов наблюдения. Однако осталось неясным, почему были взяты именно эти функции. Рассмотрим один из методов, позволяющих их получить. Для понимания его сущности обратимся к следующему примеру.

Предположим, что график плотности распределения генеральной совокупности  $X$  имеет вид равнобедренного треугольника  $ABC$ , длина основания и высота которого зафиксированы, а неизвестным параметром  $\theta$  является абсцисса точки  $D$  – середины отрезка  $AB$ . Пусть  $x_1, x_2, \dots, x_n$  – выборка из генеральной совокупности  $X$ . Зададимся вопросом: в какую точку оси абсцисс необходимо поместить точку  $D$ , если в результате опыта получена именно выборка  $x_1, x_2, \dots, x_n$ ? Конечно, никаких ограничений для ее рас-

положения на оси  $x$  нет. Но если мы сдвинем треугольник далеко влево или вправо от элементов выборки, то вероятность получения выборки, попавшей в промежуток  $[L, M]$ , которому принадлежит точка  $D$ , будет равна нулю, так как

$$P(X \in [L, M]) = \int_{[L, M]} p(x) dx = \int_{[L, M]} 0 \cdot dx = 0.$$

Поэтому точка  $D$  должна лежать в "гуще" выборки, т.е. таким образом, чтобы значения ординат  $p(x_i, \theta)$  были в совокупности как можно больше. Тогда становится правдоподобным получение именно выборки  $x_1, x_2, \dots, x_n$ . Данный метод называется *методом максимального правдоподобия*. Итак, параметр  $\theta$ , согласно этому методу, нужно выбирать так, чтобы вероятность получения набора значений  $x_1, x_2, \dots, x_n$  случайной величины  $X$  при этом значении  $\theta$  была наибольшей. Конечно, о вероятности получения данного набора значений мы строго можем говорить лишь в том случае, когда рассматриваемая генеральная совокупность распределена дискретно. Напомним, что для непрерывных случайных величин любые конкретные значения появляются с нулевой вероятностью. Поэтому метод максимального правдоподобия имеет некоторые различия в случае дискретных и непрерывных генеральных совокупностей.

**Дискретная генеральная совокупность.** Пусть  $X$  – дискретная генеральная совокупность, распределение которой зависит от некоторого параметра  $\theta$ , т.е.

$$P(X = y_i) = p_j(\theta),$$

где  $j = 1, \dots, m$ ;  $y_1, \dots, y_m$  – все различные значения, которые может принимать случайная величина  $X$ , а вероятности, с которыми эти значения появляются, зависят от параметра  $\theta$ . Предположим, что  $x_1, x_2, \dots, x_n$  – выборка из генеральной совокупности  $X$ , причем значение  $y_j$  встречается в выборке  $n_j$  раз, т.е.  $n_j$  – частота значения  $y_j$ , и поэтому имеет место равенство

$$\sum_{j=1}^m n_j = n.$$

Учитывая независимость случайных величин  $X_1, \dots, X_n$ , вероятность получения выборки  $x_1, x_2, \dots, x_n$  можно представить как

$$P(X_1 = x_1; \dots; X_n = x_n) = P(X_1 = x_1) \dots P(X_n = x_n).$$

Эта вероятность есть функция от  $x_1, x_2, \dots, x_n$ , которая называется функцией максимального правдоподобия и обозначается  $L(x_1, x_2, \dots, x_n, \theta) = P(X_1 = x_1) \dots P(X_n = x_n)$ .

Учитывая, что значение  $y_i$  встречается в выборке  $n_j$  раз, получаем

$$L(x_1, \dots, x_n, \theta) = p_1^{n_1}(\theta) \dots p_m^{n_m}(\theta).$$

Как уже было сказано, суть метода максимального правдоподобия состоит в том, что в качестве параметра  $\theta$  берется такое значение, которое максимизирует функцию  $L(x_1, \dots, x_n, \theta)$ . Полученное значение, если оно существует, является функцией от  $x_1, x_2, \dots, x_n$ , т.е.  $\theta = \theta_{МП}^*(x_1, x_2, \dots, x_n)$ . Заменяя элементы  $x_1, x_2, \dots, x_n$  случайными величинами  $X_1, \dots, X_n$ , получаем оценку максимального правдоподобия  $\theta_{МП}^*(X_1, X_2, \dots, X_n)$ .

Точка максимума функции  $L(x_1, \dots, x_n, \theta)$  удовлетворяет нелинейному (в общем случае) уравнению

$$\frac{\partial L(x_1, \dots, x_n, \theta)}{\partial \theta} = 0, \quad (3.20)$$

и поэтому конкретное значение оценки  $\theta_{МП}^*(x_1, x_2, \dots, x_n)$  определяют как корень уравнения (3.20).

Функции  $L(x_1, \dots, x_n, \theta)$  и  $\ln L(x_1, \dots, x_n, \theta)$  достигают максимума при одном и том же значении  $\theta$ . Поэтому вместо отыскания максимума функции  $L(x_1, \dots, x_n, \theta)$  находят максимум функции

$\ln L(x_1, \dots, x_n, \theta)$ . Эта функция получила название *логарифмической функции правдоподобия*.

Построение оценки максимального правдоподобия можно разбить на следующие этапы:

Этап 1. Определяют производную логарифмической функции правдоподобия по параметру  $\theta$ .

Этап 2. Приравнявая производную к нулю, находят критическую точку  $\theta_{кр}$  – корень уравнения правдоподобия

$$\frac{\partial L(x_1, \dots, x_n, \theta)}{\partial \theta} = 0.$$

Этап 3. Находят вторую производную  $\frac{\partial^2 \ln L}{\partial \theta^2}$  и ее значение в точке  $\theta_{кр}$ . Если вторая производная в точке  $\theta_{кр}$  меньше нуля, то в точке  $\theta_{кр}$  функция  $L(x_1, \dots, x_n, \theta)$  достигает максимума.

Найденная таким образом  $\theta_{МП}^*$  является функцией случайных величин  $X_1, X_2, \dots, X_n$  и, следовательно, сама является случайной величиной. Конкретное значение оценки  $\theta_{МП}^*$  получается при подстановке в  $\theta_{МП}^*(X_1, \dots, X_n)$  вместо  $X_1, X_2, \dots, X_n$  значений выборки  $x_1, x_2, \dots, x_n$ .

**Непрерывная генеральная совокупность.** Рассмотрим случай, когда генеральная совокупность имеет непрерывный ряд распределения. Функцию максимального правдоподобия определим по правилу

$$L(x_1, \dots, x_n, \theta) = p(x_1, \theta) \cdots p(x_n, \theta),$$

где  $p(x, \theta)$  – плотность распределения генеральной совокупности. Все остальное, изложенное для дискретного случая, переносится на непрерывный.

♦ **Пример 3.1.** Проводится  $n$  независимых опытов, в каждом из которых событие  $A$  повторяется с неизвестной вероятностью  $p$ . Рассмотрим генеральную совокупность  $X$  – количество появлений

события  $A$  в одном опыте. По выборке  $x_1, \dots, x_n$  из генеральной совокупности  $X$  необходимо оценить параметр  $p$ .

*Решение.* Выборка  $x_1, \dots, x_n$  состоит из нулей и единиц, причем  $x_i = 1$ , если в  $i$ -м опыте событие  $A$  произошло, и  $x_i = 0$ , если событие не произошло. Предположим, что  $m$  – частота появления события  $A$  в  $n$  опытах. Тогда выборка  $x_1, \dots, x_n$  содержит  $m$  единиц и  $(n - m)$  нулей. Так как  $P(X = 1) = p, P(X = 0) = 1 - p$ , то

$$L(x_1, \dots, x_n, \theta) = p^m (1 - p)^{n-m}.$$

Найдем точку максимума логарифмической функции максимального правдоподобия

$$\ln L(x_1, \dots, x_n, \theta) = m \ln p + (n - m) \ln(1 - p).$$

Определим из уравнения

$$\frac{\partial \ln L}{\partial p} = 0$$

критическую точку. Имеем

$$\frac{\partial \ln L}{\partial p} = \frac{m}{p} - \frac{n - m}{1 - p}.$$

Решая уравнение

$$\frac{m}{p} - \frac{n - m}{1 - p} = 0,$$

находим  $p_{кр} = \frac{m}{n}$ . Убедимся, что при данном значении параметра  $p_{кр}$  функция  $\ln L$  достигает максимума. Для этого нужно проверить, что

$$\frac{\partial^2 \ln L}{\partial p^2} = \frac{m}{p^2} - \frac{n - m}{(1 - p)^2} < 0.$$

Подставляя в это неравенство вместо  $p$  значение  $p_{кр}$ , убеждаемся в его справедливости. Значит,  $p_{кр} = \frac{m}{n}$  – оценка максимального правдоподобия, т.е.  $p_{МП}^* = \frac{m}{n}$ . Заметим, что полученная оценка –

относительная частота – является состоятельной и несмещенной оценкой для параметра  $p$ . ☹

♦ **Пример 3.2.** Найти оценку максимального правдоподобия для параметра  $\lambda$  распределения Пуассона.

*Решение.* Напомним, что распределение Пуассона имеет вид

$$P(X = m) = \frac{\lambda^m}{m!} e^{-\lambda},$$

где  $m$  принимает любые целые неотрицательные значения. Пусть  $x_1, \dots, x_n$  – выборка из генеральной совокупности  $X$ . Тогда

$$L(x_1, \dots, x_n, \lambda) = \prod_{i=1}^n \frac{\lambda^{x_i}}{x_i!} e^{-\lambda}.$$

Преобразовав произведение, получим

$$L(x_1, \dots, x_n, \lambda) = \frac{\lambda^{x_1 + \dots + x_n}}{x_1! \cdot x_2! \cdot \dots \cdot x_n!} e^{-n\lambda}.$$

Поэтому логарифмическая функция максимального правдоподобия имеет вид:

$$\ln L = -n\lambda + (x_1 + \dots + x_n) \ln \lambda - \ln(x_1! \cdot \dots \cdot x_n!).$$

Находим критическую точку, решая уравнение

$$\frac{\partial \ln L}{\partial \lambda} = 0.$$

Получим

$$-n + \frac{x_1 + \dots + x_n}{\lambda} = 0.$$

Отсюда  $\lambda_{кр} = \frac{x_1 + \dots + x_n}{n}$ . Так как

$$\frac{\partial^2 \ln L}{\partial \lambda^2} = -\frac{x_1 + \dots + x_n}{\lambda^2} < 0$$

при  $\lambda = \lambda_{кр}$ , то найденная критическая точка есть точка максимума. Поэтому оценка максимального правдоподобия для параметра  $\lambda$  является случайной величиной

$$\lambda_{МП}^* = \frac{X_1 + \dots + X_n}{n},$$

т.е.  $\bar{X}_n$ . ☹

♦ **Пример 3.3.** Найти оценку максимального правдоподобия для параметра  $\alpha$  показательного распределения

$$p(x) = \begin{cases} \alpha e^{-\alpha x}, & x > 0; \\ 0, & x \leq 0. \end{cases} \quad (3.21)$$

*Решение.* По выборке  $x_1, \dots, x_n$ , состоящей из положительных чисел, находим

$$L(x_1, \dots, x_n, \alpha) = \prod_{i=1}^n \alpha e^{-\alpha x_i} = \alpha^n e^{-\alpha(x_1 + \dots + x_n)}.$$

Поэтому

$$\ln L = n \ln \alpha - \alpha(x_1 + \dots + x_n).$$

Решая уравнение

$$\frac{\partial \ln L}{\partial \alpha} = 0,$$

находим  $\alpha = \frac{n}{x_1 + \dots + x_n}$ . Так как условие

$$\frac{\partial^2 \ln L}{\partial \alpha^2} = -\frac{n}{\alpha^2} < 0$$

при  $\lambda = \lambda_{кр}$  выполняется, то оценкой максимального правдоподобия для параметра  $\alpha$  является

$$\alpha_{МП}^* = \frac{1}{\bar{X}_n}. \quad \text{☹}$$

♦ **Пример 3.4.** Найти оценки максимального правдоподобия для параметров  $a$  и  $\sigma$  нормально распределенной генеральной совокупности.

*Решение.* Учитывая, что плотность распределения в данном случае

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

получим по выборке  $x_1, \dots, x_n$

$$L(x_1, \dots, x_n, a, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i-a)^2}{2\sigma^2}} = \frac{1}{(\sqrt{2\pi})^n \sigma^n} e^{-\sum_{i=1}^n \frac{(x_i-a)^2}{2\sigma^2}}.$$

Отсюда

$$\ln L = -\frac{n}{2} \ln 2\pi - n \ln \sigma - \sum_{i=1}^n \frac{(x_i - a)^2}{2\sigma^2}.$$

Находим критические точки этой функции, решая систему уравнений

$$\frac{\partial \ln L}{\partial a} = 0; \quad \frac{\partial \ln L}{\partial \sigma} = 0.$$

Вычисляя частные производные, получим

$$\begin{aligned} \frac{\partial \ln L}{\partial a} &= \sum_{i=1}^n \frac{(x_i - a)}{\sigma^2} = 0, \\ \frac{\partial \ln L}{\partial \sigma} &= -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - a)^2 = 0. \end{aligned}$$

Отсюда

$$a_{кр} = \frac{x_1 + \dots + x_n}{n}; \quad (3.22)$$

$$\sigma_{кр}^2 = \frac{\sum_{i=1}^n (x_i - a_{кр})^2}{n}. \quad (3.23)$$

Проверим, что при найденных значениях  $a_{кр}$  и  $\sigma_{кр}$  функция  $\ln L$  принимает максимальное значение. Для этого нужно проверить выполнение неравенств

$$\frac{\partial^2 \ln L}{\partial a^2} < 0, \quad \begin{vmatrix} \frac{\partial^2 \ln L}{\partial a^2} & \frac{\partial^2 \ln L}{\partial a \partial \sigma} \\ \frac{\partial^2 \ln L}{\partial a \partial \sigma} & \frac{\partial^2 \ln L}{\partial \sigma^2} \end{vmatrix} > 0.$$

Вычислим вторые производные:

$$\begin{aligned} \frac{\partial^2 \ln L}{\partial a^2} &= -\frac{n}{\sigma^2} < 0; \\ \frac{\partial^2 \ln L}{\partial a \partial \sigma} &= \frac{\partial^2 \ln L}{\partial \sigma \partial a} = -2 \sum_{i=1}^n \frac{x_i - a}{\sigma^3}; \\ \frac{\partial^2 \ln L}{\partial \sigma^2} &= \frac{n}{\sigma^2} = \frac{n}{\sigma^2} - \frac{3}{\sigma^4} \sum_{i=1}^n (x_i - a)^2. \end{aligned} \quad (3.24)$$

Подставляя значения для  $a_{кр}$  и  $\sigma_{кр}^2$  из (3.22) и (3.23), получаем:

$$\begin{aligned} \frac{\partial^2 \ln L}{\partial \sigma \partial a} &= -\frac{2}{\sigma^3} \left( \sum_{i=1}^n x_i - \sum_{i=1}^n x_i \right) = 0; \\ \frac{\partial^2 \ln L}{\partial \sigma^2} &= \frac{n}{d_{\epsilon}} - \frac{3}{d_{\epsilon}^2} n d_{\epsilon} = -\frac{2n}{d_{\epsilon}}, \end{aligned} \quad (3.25)$$

где  $d_{\epsilon}$  – значения выборочной дисперсии.

Вычисляя определитель в критической точке, получим

$$\begin{vmatrix} \frac{\partial^2 \ln L}{\partial a^2} & \frac{\partial^2 \ln L}{\partial a \partial \sigma} \\ \frac{\partial^2 \ln L}{\partial a \partial \sigma} & \frac{\partial^2 \ln L}{\partial \sigma^2} \end{vmatrix} = \begin{vmatrix} -\frac{n}{d_{\epsilon}} & 0 \\ 0 & -\frac{2n}{d_{\epsilon}} \end{vmatrix} = \frac{2n^2}{d_{\epsilon}^2} > 0.$$

Поэтому при значениях  $a_{кр}$  и  $\sigma_{кр}^2$ , определенных по формулам (3.22) и (3.23), функция  $\ln L$  принимает максимальное значение. Следовательно, оценками максимального правдоподобия будут

$$a_{МП}^* = \bar{X}_{\epsilon}; \quad \sigma_{МП}^* = \sqrt{D_{\epsilon}}. \quad \bullet$$

♦ **Пример 3.5.** Генеральная совокупность распределена равномерно на интервале  $(a, b)$ . По выборке  $x_1, \dots, x_n$  оценить параметры  $a$  и  $b$ .

*Решение.* Найдем оценки максимального правдоподобия для параметров  $a$  и  $b$ . Плотность генеральной совокупности имеет вид:

$$p(x, a, b) = \begin{cases} \frac{1}{b-a}, & x \in (a, b) \\ 0, & x \notin (a, b) \end{cases}. \quad (3.26)$$

Поэтому функция максимального правдоподобия

$$L(x_1, \dots, x_n, a, b) = \prod_{i=1}^n p(x_i, a, b)$$

равна нулю, если хотя бы один сомножитель произведения равен нулю, и больше нуля, если все значения  $x_1, \dots, x_n$  лежат на интервале  $(a, b)$ , т.е.

$$a \leq \min(x_1, \dots, x_n), \quad b \geq \max(x_1, \dots, x_n). \quad (3.27)$$

Тогда  $L(x_1, \dots, x_n, a, b) = \frac{1}{(b-a)^n}$ . Значение этой функции будет максимальным, если величина  $(b-a)$  минимальна. Учитывая (3.27), получим

$$a_{кр} = \min(x_1, \dots, x_n), \quad b_{кр} = \max(x_1, \dots, x_n),$$

т.е.  $a_{мп}^* = \min(X_1, \dots, X_n), \quad b_{мп}^* = \max(X_1, \dots, X_n)$ . ☹

### 3.6. Вычисление точечных оценок в Excel

**Вычисление исправленной дисперсии.** В п. 3.3 показано, что оценка

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_e)^2 \quad (3.28)$$

является несмещенной точечной оценкой для дисперсии случайной величины, и такую оценку часто называют *исправленной дисперсией*.

Для вычисления выборочного значения этой оценки можно использовать статистическую **функцию** Excel **ДИСП**, обращение к которой имеет вид:

$$=ДИСП(арг1; арг2; ...; арг30),$$

где  $арг1; арг2; ...; арг30$  – числа или адреса ячеек, содержащих числовые величины.

♦ **Пример 3.6.** По выборке примера 2.3 вычислить оценку (3.28).

*Решение.* Первоначально, начиная с ячейки А3, введем в столбец А 55 элементов выборки (рис. 3.3). Затем, используя функции КВАДРОТКЛ, ДИСП (как показано на рис. 3.3), вычислим оценку (3.28). Видно ожидаемое совпадение двух вычисленных значений.

☹

	А	В	С	Д	Е
1					
2	Выборочные значения				
3	20,3	=КВАДРОТКЛ(А3:А57)/(55-1)			
4	15,3				
5	14,3		8,760		
6	19,3				
7	10,1		8,760		
8	13,9				
9	19,5		=ДИСП(А3:А57)		
10	17,8				

Рис. 3.3. Фрагмент вычисления исправленной дисперсии

**Вычисление оценок максимального правдоподобия.** В п. 3.5 были рассмотрены оценки, вычисляемые из условия максимума функционала правдоподобия. В приведенных примерах из условий максимума были получены алгебраические уравнения, решения которых определялись достаточно просто.

В общем случае не удастся получить таких простых соотношений и оценки вычисляются непосредственным определением

точек максимума функционала правдоподобия, т.е. необходимо решить оптимизационную задачу.

Для решения такой задачи в Excel есть команда *Поиск решения* пункта меню **Сервис**. Эта команда позволяет решать не только задачи безусловной оптимизации, но и задачи условной оптимизации, т.е. когда ищется максимум функционала с учетом дополнительных ограничений на значения искомых оценок. Например, значение дисперсии  $\sigma^2$  не может быть отрицательным.

Применение команды *Поиск решения* для вычисления оценок максимального правдоподобия покажем на следующем примере.

♦ **Пример 3.7.** По выборке примера 2.3 вычислить оценки максимального правдоподобия для математического ожидания  $a$  и дисперсии  $\sigma^2$  из условия максимума функционала правдоподобия вида:

$$-\frac{n}{2}\ln(2\pi) - n\ln(\sigma) - \sum_{i=1}^n \frac{(x_i - a)^2}{2\sigma^2}, \quad (3.29)$$

предполагая при этом, что выборка порождена случайной величиной, подчиняющейся нормальному распределению.

*Решение.* Первоначально, начиная с ячейки A3, введем в столбец A 55 элементов выборки (диапазон A3:A57). Затем в ячейку C8 занесем произвольное значение  $a$  (например, 10), в ячейку D8 – значение  $\sigma$  (например, значение  $4 > 0$ ), в ячейке E8 вычислим  $\sigma^2$ . В ячейках B3:B57 запрограммируем вычисление разностей  $x_i - a$  (рис. 3.4). В ячейке C5 запрограммируем вычисление величины функционала (3.29). В верхней части документа на рис. 3.4 показана запрограммированная формула.

После этих подготовительных операций можно перейти к выполнению команды *Поиск решения*. Для этого необходимо обратиться к пункту основного меню **Сервис** и в появившемся меню щелкнуть мышью на команде *Поиск решения*. Затем в появившемся диалоговом окне выполнить следующие действия (см. рис. 3.4):

- в поле ввода *Установить целевую ячейку*: ввести адрес ячейки, в которой вычисляется значение минимизируемого функционала (в нашем примере C5);

- включить опцию *Равной*: максимальному значению (ищутся значения, при которых функционал достигает максимального значения);

- в поле *Изменяя ячейки*: ввести адреса ячеек, в которых находятся значения искомых оценок (в нашем примере это ячейки C8:D8);

- щелкнув мышью на кнопке *Добавить*, сформировать ограничения на значения искомых оценок (в нашем примере это требование  $\sigma \geq 0.0000001$ , чтобы  $\ln(\sigma)$  не был равен  $-\infty$ ).

	A	B	C	D	E	F
1	=-55/2*LN(2*3,1415)-55*LN(D8)-СУММКВ(B3:B57)/(2*D8^2)					
2	Выборочные значения					
3	20,3	10,3				
4	15,3	5,3				
5	14,3	4,3	-249,03			
6	19,3	9,3				
7	10,1	0,1				
8	13,9	3,9	10,000	4,000	16,000	
9	19,5	9,5	a	σ	σ²	

Величина функционала =A6-C8

**Поиск решения**

Установить целевую ячейку: \$C\$5

Равной: ☒ максимальному значению ☐ значению: 0

☐ минимальному значению

Изменяя ячейки: \$C\$8:\$D\$8

Ограничения: \$D\$8 >= 0,0000001

Кнопки: Выполнить, Закрыть, Предположить, Параметры, Добавить, Изменить, Удалить, Восстановить, Справка

Рис. 3.4. Задание параметров команды *Поиск решения*

После выполнения этих операций щелкнуть на кнопке *Выполнить*. Начинается поиск решения введенной оптимизационной задачи. Спустя некоторое время на экране появится новое диалоговое окно *Результаты поиска решения* (рис. 3.5). Для сохранения найденных значений оценок в соответствующих ячейках необходимо включить опцию *Сохранить найденное решение* и щелкнуть на кнопке ОК.

	A	B	C	D	E	F
1	$=-55/2*\text{LN}(2*3,1415)-55*\text{LN}(D8)-\text{СУММКВ}(B3:B57)/(2*D8^2)$					
2	Выборочные значения					
3	20,3	2,392733				
4	15,3	-2,60727				
5	14,3	-3,60727	-137,22	Величина функционала		
6	19,3	1,392733			=A6-C8	
7	10,1	-7,80727		=D8^2		
8	13,9	-4,00727	17,907	2,933	8,601	
9	19,5	1,592733	$\alpha$	$\sigma$	$\sigma^2$	
10						
11						
12						
13						
14						
15						
16						
17						
18	13,5	-4,40727				

Рис. 3.5. Результаты выполнения команды *Поиск решения*

Из рис. 3.5 видно, что вычисленные значения оценок находятся в ячейках C8, D8 и равны  $\alpha = 17.907$ ,  $\sigma = 2.933$ . Ячейка C5 содержит значение максимизируемого функционала, равное  $-137.22$ . Сравнивая вычисленные значения оценок  $\alpha = 17.907$  и  $\sigma^2 = 8.601$  с

выборочными оценками примера 2.11 (см. рис. 2.7), видим их полное совпадение. ☺

**Задание 3.1.** Предполагая, что выборка примера 2.1 порождена случайной величиной, имеющей показательное распределение (3.21), вычислить оценку максимального правдоподобия для параметра  $\alpha$ , используя команду *Поиск решения*.

*Рекомендация.* Оценку максимального правдоподобия осуществлять из условия максимума функционала

$$n \ln(\alpha) - \alpha \sum_{i=1}^n x_i$$

при ограничении  $\alpha > 0$ . При вызове команды *Поиск решения* использовать пример 3.7. ♥

### Функции Excel для вычисления других точечных оценок.

Для вычисления среднеквадратичных отклонений можно использовать следующие функции Excel.

**Функция СТАНДОТКЛОН** вычисляет

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

Обращение к ней имеет вид:

$$=\text{СТАНДОТКЛОН}(arg1; arg2; \dots; arg30),$$

где  $arg1; arg2; \dots; arg30$  – числовые константы или адреса ячеек, содержащих числовые данные.

**Функция СТАНДОТКЛОНП** вычисляет

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

Обращение к ней имеет вид:

$$=\text{СТАНДОТКЛОНП}(arg1; arg2; \dots; arg30),$$

где  $arg1; arg2; \dots; arg30$  – числовые константы или адреса ячеек, содержащих числовые данные.



**Функция ЭКСЦЕСС** вычисляет оценку

$$\frac{n(n+1)}{(n-1)(n-2)(n-3)} \sum_{i=1}^n \left( \frac{x_i - \bar{x}_e}{d_e} \right)^2 - \frac{3(n-1)^2}{(n-2)(n-3)}$$

для характеристики эксцесс  $\frac{\mu_4}{\sigma^4} - 3$ , которая определяет островершинность или плосковершинность плотности распределения.

Обращение к функции имеет вид:

$$= \text{ЭКСЦЕСС}(arg1; arg2; \dots; arg30),$$

где  $arg1; arg2; \dots; arg30$  – числовые константы или адреса ячеек, содержащих числовые данные.

**Функция МОДА** вычисляет наиболее часто встречающееся значение в заданных аргументах функции, т.е. значение, встречающееся в выборке с максимальной частотой.

Обращение к функции имеет вид:

$$= \text{МОДА}(arg1; arg2; \dots; arg30),$$

где  $arg1; arg2; \dots; arg30$  – числовые константы или адреса ячеек, содержащих числовые данные.

Если в заданных значениях аргументов **нет повторяющихся значений**, то функция возвращает признак ошибки #Н/Д.

**Функция МЕДИАНА** вычисляет значение выборки, приходящееся на середину упорядоченной выборочной совокупности. Если выборка имеет четное число элементов, то значение функции будет равно среднему двух значений, находящихся по середине упорядоченной выборочной совокупности. Например, медиана выборки (200, 236, 250, 305, 337, 220) будет равна  $(236 + 250) / 2 = 243$ .

Обращение к функции имеет вид:

$$= \text{МЕДИАНА}(arg1; arg2; \dots; arg30),$$

где  $arg1; arg2; \dots; arg30$  – числовые константы или адреса ячеек, содержащих числовые данные.

**Функция СКОС** вычисляет оценку

$$\frac{n}{(n-1)(n-2)} \sum_{i=1}^n \frac{(x_i - \bar{x}_e)^3}{d_e^{3/2}}$$

для характеристики асимметрии  $\frac{\mu_3}{\sigma^3}$ , которая для симметричной плотности распределения равна 0.

Обращение к функции имеет вид:

$$= \text{СКОС}(arg1; arg2; \dots; arg30),$$

где  $arg1; arg2; \dots; arg30$  – числовые константы или адреса ячеек, содержащих числовые данные.

**Вычисление описательных статистик.** Описательные статистики можно разделить на следующие группы:

- *характеристики положения* описывают положение данных на числовой оси (среднее, минимальное и максимальное значения, медиана и др.);
- *характеристики разброса* описывают степень разброса данных относительно своего центра (дисперсия, размах выборки, эксцесс, среднее квадратическое отклонение и др.);
- *характеристики асимметрии* определяют симметрию распределения данных относительно своего центра (коэффициент асимметрии, положение медианы относительно среднего и др.);
- *характеристики, описывающие закон распределения* (частоты, относительные частоты, гистограммы и др.).

Основные характеристики положения, разброса и асимметрии можно вычислить, используя режим **Описательная статистика** команды **Пакет анализа**.

Для вызова режима **Описательная статистика** необходимо обратиться к пункту **Сервис**, команде **Пакет анализа**, выбрать в списке режимов **Описательная статистика** и щелкнуть на кнопке ОК. В появившемся диалоговом окне Описательная статистика задать следующие параметры (рис. 3.6):

**Входной интервал:** – адреса ячеек, содержащих элементы выборки.

**Группирование:** – задает способ расположения (по столбцам или по строкам) элементов выборки.

**Метки в первой строке** – включается, если первая строка (столбец) во входном интервале содержит заголовки.

	A	B	C	D	E	F
1						
2	Выборочные значения					
3	20,3					
4	15,3					
5	14,3					

Описательная статистика

Входные данные

Входной интервал: A3:A57

Группирование:

☒ по столбцам
☐ по строкам

☐ Метки в первой строке

Параметры вывода

☒ Выходной интервал: \$D\$5
☐ Новый рабочий лист:
☐ Новая рабочая книга

☒ Итоговая статистика
☒ Уровень надежности: 95 %
☒ К-ый наименьший: 2
☒ К-ый наибольший: 1

OK

Отмена

Справка

Рис. 3.6. Параметры режима *Описательная статистика*

*Выходной интервал:* / *Новый рабочий лист:* / *Новая рабочая книга* – определяет место вывода результатов вычислений. При включении *Выходной интервал:* в поле вводится адрес ячейки, начиная с которой будут выводиться результаты.

*Итоговая статистика:* – включается, если необходимо вывести по одному полю для каждой из вычисленных характеристик.

*Уровень надежности:* – включается, если необходимо вычислить доверительный интервал для математического ожидания с задаваемым (в %) уровнем надежности  $\gamma$ .

*К-й наименьший:* – включается, если необходимо вычислить к-й наименьший (начиная с  $x_{\min}$ ) элемент выборки. При  $k = 1$  вычисляется наименьшее значение.

*К-й наибольший:* – включается, если необходимо вычислить к-й наибольший (начиная с  $x_{\max}$ ) элемент выборки. При  $k = 1$  вычисляется наибольшее значение.

Пример задания параметров приведен на рис. 3.6.

Результаты работы режима *Описательная статистика* выводятся в виде таблицы, в левом столбце которой приводится название вычисленной характеристики (рис. 3.7), позволяющее однозначно трактовать характеристику. Тем не менее, поясним следующие названия характеристик:

- *Интервал* – определяет размах выборки  $x_{\max} - x_{\min}$ ;
- *Сумма* – определяет сумму всех элементов выборки;
- *Счет* – определяет число обработанных элементов выборки;
- *Уровень надежности* – определяет величину  $\Delta_{\bar{x}}$ , от которой зависит доверительный интервал для математического ожидания, имеющий вид

$$[\bar{x}_g - \Delta_{\bar{x}}, \bar{x}_g + \Delta_{\bar{x}}],$$

где  $\bar{x}_g$  – выборочное среднее (подробнее см. п. 4.3).

♦ **Пример 3.8.** По выборке примера 2.3 вычислить описательные статистики, используя режим *Описательная статистика*.

*Решение.* Первоначально, начиная с ячейки A3, введем в столбец A 55 элементов выборки. После этого обратимся к пункту **Сервис**, команде *Пакет анализа*. В списке режимов выберем *Описательная статистика*. В появившемся диалоговом окне включим параметры, показанные на рис. 3.6, и щелкнем ОК. Вычисленные характеристики приведены на рис. 3.7. ☺

<i>Столбец1</i>	
Среднее	<b>17,907</b>
Стандартная ошибка	<b>0,399</b>
Медiana	<b>18,100</b>
Мода	<b>17,800</b>
Стандартное отклонение	<b>2,960</b>
Дисперсия выборки	<b>8,760</b>
Эксцесс	<b>-0,078</b>
Асимметричность	<b>-0,386</b>
Интервал	<b>13,700</b>
Минимум	<b>10,100</b>
Максимум	<b>23,800</b>
Сумма	<b>984,900</b>
Счет	<b>55,000</b>
Наибольший(2)	<b>23,300</b>
Наименьший(1)	<b>10,100</b>
Уровень надежности(95,0%)	<b>0,800</b>

Рис. 3.7. Результаты работы *Описательная статистика*

**Задание 3.2.** Сравните значения характеристик (см. рис. 3.7) со значениями аналогичных характеристик, вычисленных в предыдущих примерах. ♥

## 4. ИНТЕРВАЛЬНЫЕ ОЦЕНКИ НЕИЗВЕСТНЫХ ПАРАМЕТРОВ

### 4.1. Некоторые распределения выборочных характеристик

Генеральные совокупности часто имеют нормальный закон распределения. В этом случае многие выборочные характеристики, в том числе  $\bar{X}_e, D_e, S^2$ , выражаются через небольшое число распределений. Как правило, в математической статистике используются не плотности этих распределений, а некоторые характеристики, представленные таблицами. Чаще всего в качестве такой характеристики выступает квантиль распределения.

Квантилем уровня  $p$  ( $0 < p < 1$ ) или  $p$ -квантилем случайной величины  $X$  называется такое число  $d_p$ , что вероятность  $P(X < d_p)$  равна заданной величине  $p$ .

Из определения следует, что если непрерывная случайная величина  $X$  имеет плотность распределения  $p(x)$ , то квантиль  $d_p$  определяется равенством

$$\int_{-\infty}^{d_p} p(x)dx = p. \quad (4.1)$$

Это означает, что площадь фигуры, ограниченной осью абсцисс, кривой  $f(x)$  и прямой  $x = d_p$ , равна величине  $p$ . На рис. 4.1,а показан квантиль  $d_{0.1}$ , а на рис. 4.1,б – квантиль  $d_{0.9}$ . Площади заштрихованных фигур равны 0.1 и 0.9 соответственно.

Рассмотрим несколько распределений, которым подчиняются выборочные характеристики и которые используются для построения интервальных оценок.

**Распределение  $\chi^2$  (распределение К. Пирсона).** Пусть  $N_1, \dots, N_n$  – независимые нормально распределенные случайные величины с параметрами (0,1). Распределение случайной величины

$$\chi_n^2 = N_1^2 + N_2^2 + N_3^2 + \dots + N_n^2 \quad (4.2)$$

называется *распределением  $\chi^2$  с  $n$  степенями свободы*, а сама величина  $\chi^2$  – случайной величиной  $\chi^2$  с  $n$  степенями свободы.

Заметим, что количество степеней свободы  $n$  является единственным параметром  $\chi^2$ -распределения и значения  $\chi^2$  неотрицательны, т.е.  $P(\chi_n^2 < 0) = 0$ .

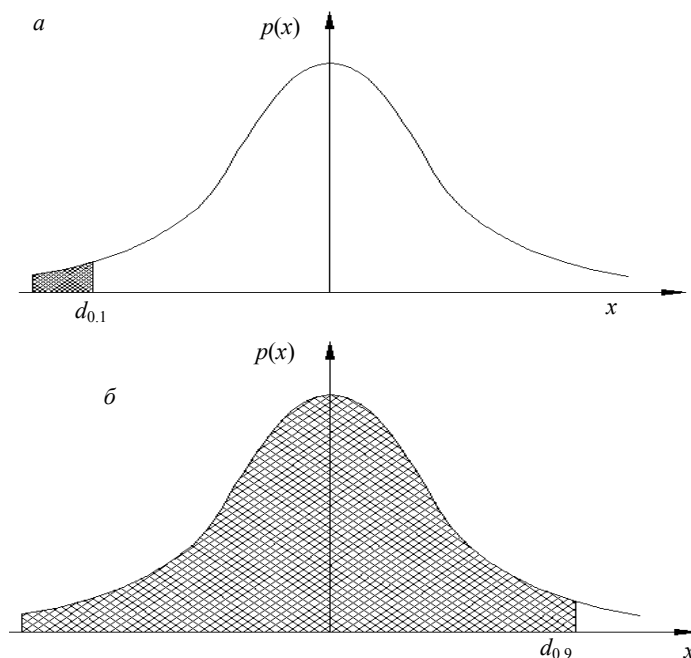


Рис. 4.1. К определению квантилей случайной величины

Определим математическое ожидание величины  $\chi^2$ . По определению (4.2) имеем

$$M(\chi_n^2) = M\left(\sum_{i=1}^n N_i^2\right) = \sum_{i=1}^n M(N_i^2) = \sum_{i=1}^n [D(N_i) + M^2(N_i)],$$

так как  $D(X) = M(X^2) - M^2(X)$ . Но  $D(N_i) = 1, M(N_i) = 0$ , а значит,  $M(\chi_n^2) = n$ . Нетрудно вычислить и дисперсию случайной ве-

личины  $\chi_n^2$ . Так как случайные величины  $N_1^2, \dots, N_n^2$  независимы, то

$$D(\chi_n^2) = nD(N_1^2) = n[M(N_1^4) - M^2(N_1^2)]. \quad (4.3)$$

Плотность распределения случайной величины  $N_1$  равна  $p(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$ , значит,

$$M(N_1^4) = \int_{-\infty}^{\infty} x^4 p(x) dx = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x^4 e^{-\frac{x^2}{2}} dx = 3.$$

Последний интеграл вычисляется методом интегрирования по частям. Далее, так как  $M(N_1^2) = 1$ , то  $D(\chi_n^2) = n(3 - 1) = 2n$ . Таким образом,  $\chi^2$ -распределение с  $n$  степенями свободы имеет следующие числовые характеристики:

$$M[\chi_n^2] = n; \quad D[\chi_n^2] = 2n. \quad (4.4)$$

Согласно центральной предельной теореме, если случайные величины  $N_1^2, N_2^2, \dots, N_n^2$  независимы, одинаково распределены и имеют конечные дисперсии, то последовательность  $\chi_n^2 = N_1^2 + \dots + N_n^2$  асимптотически нормальна. Другими словами, при больших значениях  $n$  распределение случайной величины  $\chi_n^2$  близко к нормальному распределению с параметрами  $a = n, \sigma^2 = 2n$ . Однако при малых значениях  $n$  функция плотности случайной величины  $\chi_n^2$  значительно отличается от кривой Гаусса.

На рис. 4.2 показаны плотности распределения  $p(x)$  случайной величины  $\chi_n^2$  при  $n = 2, n = 6$  и  $n = 20$ . Видно, что при увеличении  $n$  плотность  $p(x)$  "приближается" к плотности нормального распределения.

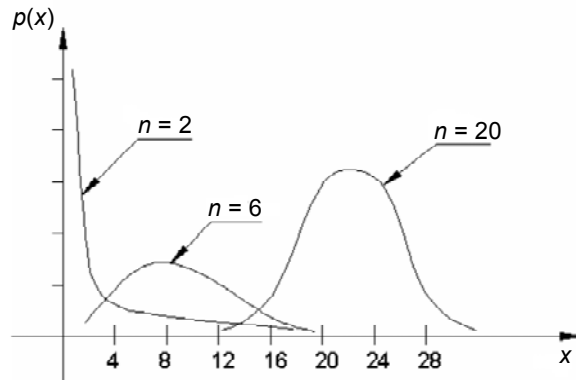


Рис. 4.2. Плотность распределения  $\chi^2$

Обратим внимание на одно замечательное свойство распределения  $\chi_n^2$ . Строго говоря, это свойство можно доказать, используя, например, производящие функции. Свойство состоит в том, что сумма независимых случайных величин  $\chi_n^2 + \chi_m^2$  также распределена по закону  $\chi^2$  с  $(n+m)$  степенями свободы. Объясняется это тем, что случайная величина  $\chi_n^2 + \chi_m^2$  представляется в виде суммы  $(n+m)$  квадратов случайных величин, независимых и нормально распределенных с параметрами  $(0,1)$ .

**Распределение Стьюдента ( $t$ -распределение).** Пусть  $N(0,1)$  – нормально распределенная случайная величина с параметрами  $a=0, \sigma=1$ , а  $\chi_n^2$  – независимая от  $N(0,1)$  случайная величина, подчиняющаяся распределению  $\chi^2$  с  $n$  степенями свободы. Тогда распределение случайной величины

$$T_n = \frac{N(0,1)\sqrt{n}}{\sqrt{\chi_n^2}} \quad (4.5)$$

называется  $t$ -распределением или *распределением Стьюдента*. Сама случайная величина (4.5) называется  $t$ -величиной с  $n$  степенями

свободы. Плотность вероятности случайной величины  $T_n$  имеет вид  $p_n = B_n \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}}$ , где  $B_n$  – некоторая константа, удовле-

творяющая условию нормирования  $\int_{-\infty}^{\infty} p_n(x) dx = 1$ . При больших значениях  $n$  кривая  $p_n(x)$  близка к кривой нормального распределения  $N(0,1)$ . Поэтому в практических расчетах при  $n > 30$  часто считают, что

$$p_n(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}.$$

Заметим, что функция плотности  $p_n(x)$  симметрична относительно оси ординат.

**Распределение Фишера ( $F$ -распределение).** Пусть  $\chi_n^2$  и  $\chi_m^2$  – независимые случайные величины, имеющие  $\chi^2$ -распределение с  $n$  и  $m$  степенями свободы соответственно. Распределение случайной величины

$$F_{n,m} = \frac{\chi_n^2/n}{\chi_m^2/m} \quad (4.6)$$

называется  $F$ -распределением или *распределением Фишера* с  $n$  и  $m$  степенями свободы, а сама величина (4.6) –  $F_{n,m}$  величиной. Так как случайные величины  $\chi_n^2 \geq 0$  и  $\chi_m^2 \geq 0$ , то  $F_{n,m} \geq 0$ .

В дальнейшем мы часто будем ссылаться на следующую теорему о распределении выборочных характеристик  $\bar{X}_e$  и  $D_e$ , доказанную Р. Фишером.

**Теорема 4.1** (о распределении выборочных характеристик). Если генеральная совокупность  $X$  распределена по нормальному закону с параметрами  $a$  и  $\sigma$ , то:

а) случайная величина  $\bar{X}_g$  распределена нормально с параметрами  $(a, \frac{\sigma}{\sqrt{n}})$ ;

б)  $nD_g/\sigma^2$  имеет распределение  $\chi_{n-1}^2$ ;

в) случайные величины  $\bar{X}_g$  и  $D_g$  независимы.

Мы не будем полностью доказывать эту теорему, а ограничимся доказательством утверждения а). Очевидно, что  $\bar{X}_g$  есть линейная комбинация

$$\bar{X}_g = \frac{1}{n} X_1 + \frac{1}{n} X_2 + \dots + \frac{1}{n} X_n$$

независимых, нормально распределенных случайных величин. Как отмечалось в курсе теории вероятностей, в этом случае случайная величина  $\bar{X}_g$  распределена нормально. Легко получить, что

$$M(\bar{X}_g) = M\left(\frac{x_1 + x_2 + \dots + x_n}{n}\right) = \frac{M(x_1) + \dots + M(x_n)}{n} = \frac{na}{n} = a,$$

$$D(\bar{X}_g) = D\left(\frac{x_1 + \dots + x_n}{n}\right) = \frac{D(x_1) + \dots + D(x_n)}{n^2} = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n}.$$

Тем самым первое утверждение теоремы доказано.

Как следует из в), используя случайные величины  $\bar{X}_g$  и  $D_g$ , можно составить случайную величину  $T_{n-1}$ . Действительно, проинормировав  $\bar{X}_g$ , получим  $\frac{(\bar{X}_g - a)\sqrt{n}}{\sigma} = N(0,1)$ . Так как  $\bar{X}_g$  и  $D_g$  независимы, то по (4.5)

$$T_{n-1} = \frac{(\bar{X}_g - a)\sqrt{n}\sqrt{n-1}}{\sigma} : \sqrt{\frac{nD_g}{\sigma^2}} = \frac{(\bar{X}_g - a)\sqrt{n-1}}{\sqrt{D_g}}.$$

Итак, мы получили

**Следствие.** Если условия теоремы о распределении выборочных характеристик выполнены, то случайная величина

$$\frac{(\bar{X}_g - a)\sqrt{n-1}}{\sqrt{D_g}}$$

имеет распределение Стьюдента с  $(n-1)$  степенями свободы.

Напомним, что исправленная дисперсия  $S^2$  определяется как

$$S^2 = \frac{n}{n-1} D_g.$$

Тогда получаем новое

**Следствие.** Если условия теоремы о распределении выборочных характеристик выполнены, то случайная величина

$$\frac{(\bar{X}_g - a)\sqrt{n}}{\sqrt{S^2}}$$

имеет распределение с  $(n-1)$  степенями свободы.

#### 4.2. Понятие интервальной оценки параметра случайной величины

Вычисляя на основании результатов наблюдений точечную оценку  $\theta^*$  неизвестного параметра  $\theta$ , мы понимаем, что величина  $\theta^*$  является (в силу своей случайности) лишь приближенным значением параметра  $\theta$ . При большом числе наблюдений точность приближения бывает достаточной для практических выводов в силу несмещенности, состоятельности и эффективности "хороших" оценок. Для выборок малого объема точечные оценки могут значительно отличаться от оцениваемого параметра и вопрос о точности получаемых оценок становится очень важным. В математической статистике он решается введением интервальных оценок.

Интервальной оценкой для параметра  $\theta$  называется такой интервал  $(\underline{\theta}^*, \bar{\theta}^*)$  со случайными границами, что

$$P(\underline{\theta}^* < \theta < \bar{\theta}^*) = \gamma. \quad (4.7)$$

Вероятность  $\gamma$  называется *надежностью интервальной оценки* или *доверительной вероятностью*, случайные величины  $\underline{\theta}^*, \bar{\theta}^*$  – *доверительными границами*, а сам интервал  $(\underline{\theta}^*, \bar{\theta}^*)$  иногда называют *доверительным интервалом*. Центром этого интервала является значение точечной оценки  $\theta^*$ .

Надежность  $\gamma$  принято выбирать равной 0.95, 0.99. Тогда событие, состоящее в том, что интервал  $(\underline{\theta}^*, \bar{\theta}^*)$  покроет параметр  $\theta$ , будет практически достоверным.

Общая теория построения интервальных оценок заключается в определении *случайной величины, зависящей от оцениваемого параметра*. Зная распределение этой случайной величины, находят соответствующие доверительные границы и сам доверительный интервал с требуемой точностью. Посмотрим, как эта идея реализуется для различных параметров.

#### 4.3. Интервальные оценки математического ожидания нормального распределения

Пусть генеральная совокупность  $X$  распределена по нормальному закону  $N(a, \sigma)$ , причем параметр  $\sigma$  известен, а параметр  $a$  требуется оценить с надежностью  $\gamma$ . По теореме о распределении

выборочных характеристик случайная величина  $\frac{(\bar{X}_\varepsilon - a)\sqrt{n}}{\sigma}$  рас-

пределена по закону  $N(0,1)$ . На рис. 4.3 изображен график функции плотности этой случайной величины, т.е. кривая

$y = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$ . Выберем число  $x_\gamma$  так, что заштрихованная площадь равна  $\gamma$ , т.е.

$$P(-x_\gamma < \frac{(\bar{X}_\varepsilon - a)\sqrt{n}}{\sigma} < x_\gamma) = \gamma. \quad (4.8)$$

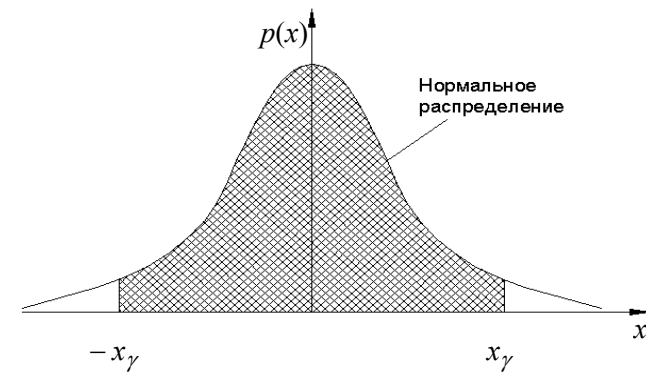


Рис. 4.3. К построению доверительных интервалов

Это значение легко находится с использованием интегральной функции Лапласа  $\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_0^x e^{-\frac{t^2}{2}} dt$ . Действительно,

$$P(-x_\gamma < N(0,1) < x_\gamma) = \Phi(x_\gamma) - \Phi(-x_\gamma) = 2\Phi(x_\gamma) = \gamma. \quad (4.9)$$

Значение  $x_\gamma$ , удовлетворяющее нелинейному уравнению

$$\Phi(x_\gamma) = \frac{\gamma}{2}, \quad (4.10)$$

находится по табл. П1.

Так как  $\sigma > 0$ , то события  $-x_\gamma < \frac{(\bar{X}_\varepsilon - a)\sqrt{n}}{\sigma} < x_\gamma$  и  $\bar{X}_\varepsilon - \frac{x_\gamma \sigma}{\sqrt{n}} < a < \bar{X}_\varepsilon + \frac{x_\gamma \sigma}{\sqrt{n}}$  эквивалентны, а значит, их вероятности равны:

$$P\left(\bar{X}_\varepsilon - \frac{x_\gamma \sigma}{\sqrt{n}} < a < \bar{X}_\varepsilon + \frac{x_\gamma \sigma}{\sqrt{n}}\right) = \gamma. \quad (4.11)$$

Таким образом, для параметра  $a$  мы построили доверительный интервал (интервальную оценку), левая граница которого

$\bar{X}_g - \frac{x_\gamma \sigma}{\sqrt{n}}$ , правая –  $\bar{X}_g + \frac{x_\gamma \sigma}{\sqrt{n}}$ , а точность –  $\delta = \frac{x_\gamma \sigma}{\sqrt{n}}$ . Центр

этого интервала находится в точке с координатой  $\bar{X}_g$ , а длина ин-

тервала  $2 \frac{x_\gamma \sigma}{\sqrt{n}}$ . Если объем выборки неограниченно возрастает, то

интервал стягивается в одну точку  $\bar{X}_g$ , которая является состоя-  
тельной и несмещенной оценкой для параметра  $a$ .

♦ **Пример 4.1.** По выборке объема  $n = 9$  найдено среднее зна-  
чение  $\bar{x}_g = 1.5$ . Считая, что генеральная совокупность распреде-  
лена по нормальному закону с  $\sigma = 2$ , определить интервальную  
оценку для математического ожидания с надежностью  $\gamma = 0.95$ .

*Решение.* Используя табл. П1, находим, что

$$\Phi(x_\gamma) = \frac{0.95}{2} = 0.475$$

при  $x_\gamma = 1.96$ . Тогда  $\delta = 1.96 \cdot \frac{2}{\sqrt{9}} = 1.31$  и доверительный интер-  
вал (4.11) имеет границы  $(\bar{X}_g - 1.31, \bar{X}_g + 1.31)$ . Таким образом, с  
вероятностью 0.95 можно быть уверенным в том, что интервал

$$(\bar{X}_g - 1.31, \bar{X}_g + 1.31) \quad (4.12)$$

накрывает параметр  $a$  или, другими словами, с вероятностью 0.95  
значение  $\bar{X}_g$  дает значение параметра  $a$  с точностью  $\delta = 1.31$ .

Заметим, что эта трактовка *неверна*, если вместо случайной  
величины  $\bar{X}_g$  использовать вычисленное по конкретной выборке  
значение  $\bar{x}_g = 1.5$ . Тогда границы интервала (0.19, 2.81) будут *не*  
*случайными* и возможны два случая:

- точка  $a$  лежит внутри этого интервала, тогда

$$P(0.19 < a < 2.81) = 1;$$

- точка  $a$  не лежит внутри (0.19, 2.81), тогда

$$P(0.19 < a < 2.81) = 0.$$

Поэтому только для интервала (4.12) со случайными границами  
можно утверждать, что

$$P(\bar{X}_g - 1.31 < a < \bar{X}_g + 1.31) = 0.95. \quad \bullet$$

Определим теперь интервальную оценку для неизвестной ге-  
неральной средней  $\bar{x}_g$  нормально распределенной генеральной со-  
вокупности  $X$  в том случае, когда *генеральная дисперсия  $D_g$  неиз-*  
*вестна*, т.е. построим доверительный интервал для параметра  $a$ ,  
если параметр  $\sigma$  неизвестен.

В отличие от предыдущего случая, вместо случайной величи-  
ны  $\frac{(\bar{X}_g - a)\sqrt{n}}{\sigma}$ , распределенной по закону  $N(0,1)$ , рассмотрим

случайную величину  $\frac{(\bar{X}_g - a)\sqrt{n-1}}{\sqrt{D_g}}$ , которая согласно следствию

из теоремы 4.1 распределена по закону Стьюдента  $T_{n-1}$ . При за-  
данном значении  $\gamma$ , пользуясь табл. П2, вычислим значение  
 $t(\gamma, n)$  из условия

$$P\left(-t(\gamma, n) < \frac{(\bar{X}_g - a)\sqrt{n-1}}{\sqrt{D_g}} < t(\gamma, n)\right) = \gamma, \quad (4.13)$$

где  $\gamma$  – надежность интервальной оценки. Заметим, что в табл. П2  
 $n$  означает не число степеней свободы, а объем выборки. Число  
степеней свободы будет равно  $n-1$ .

Замена случайной величины  $\frac{(\bar{X}_g - a)\sqrt{n}}{\sigma}$  на случайную вели-  
чину  $\frac{(\bar{X}_g - a)\sqrt{n-1}}{\sqrt{D_g}}$  вызвана тем, что закон распределения послед-

ней случайной величины известен и в ее запись не входит неиз-  
вестный в данном случае параметр  $\sigma$ . Из условия (4.13) получаем

$$P\left(\bar{X}_g - \frac{t(\gamma, n)\sqrt{D_g}}{\sqrt{n-1}} < a < \bar{X}_g + \frac{t(\gamma, n)\sqrt{D_g}}{\sqrt{n-1}}\right) = \gamma.$$



Таким образом, интервальная оценка надежности  $\gamma$  для неизвестной генеральной средней  $a$  имеет границы

$$\left( \bar{X}_g - \frac{t(\gamma, n)\sqrt{D_g}}{\sqrt{n-1}}, \bar{X}_g + \frac{t(\gamma, n)\sqrt{D_g}}{\sqrt{n-1}} \right).$$

Выразим границы интервала через исправленную дисперсию  $S^2$ .

Так как  $S^2 = \frac{n}{n-1} D_g$ , то  $\frac{\sqrt{D_g}}{\sqrt{n-1}} = \frac{S}{\sqrt{n}}$ . Поэтому

$$\frac{t(\gamma, n)\sqrt{D_g}}{\sqrt{n-1}} = \frac{t(\gamma, n)S}{\sqrt{n}}.$$

Значит, границы доверительного интервала можно записать как

$$\left( \bar{X}_g - \frac{t(\gamma, n)S}{\sqrt{n}}, \bar{X}_g + \frac{t(\gamma, n)S}{\sqrt{n}} \right), \quad (4.14)$$

а точность интервальной оценки определить соотношением

$$\delta = \frac{t(\gamma, n)}{\sqrt{n}} S. \quad (4.15)$$

Как и в предыдущем случае, центр интервала находится в точке  $\bar{X}_g$ , но длина интервала  $2\frac{t(\gamma, n)}{\sqrt{n}}S$  является случайной величиной, принимающей тем меньшие значения, чем больше значение  $n$ . Это объясняется тем, что наличие большей информации  $x_1, \dots, x_n$  о генеральной совокупности  $X$  позволяет сузить интервал.

♦ **Пример 4.2.** По выборке объема  $n = 9$  из нормально распределенной генеральной совокупности найдены значения  $\bar{x}_g = 1.5$  и  $s = 2$ . Построить интервальную оценку для математического ожидания с надежностью  $\gamma = 0.95$ .

*Решение.* Пользуясь табл. П2, находим величину  $t(0.95, 9) = 2.31$ . Тогда точность  $\delta$  определяется соотношением

(см. (4.15)):  $\delta = \frac{t(0.95, 9)S}{\sqrt{n}} = \frac{2.31}{3}S = 0.77S$ , а интервальная оценка

имеет границы  $(\bar{X}_g - 0.77 \cdot S, \bar{X}_g + 0.77 \cdot S)$ , которые зависят от двух случайных величин:  $\bar{X}_g$  и  $S$ . Подставляя вместо  $S$  ее вычисленное значение  $s = 2$ , получаем интервал

$$(\bar{X}_g - 1.54, \bar{X}_g + 1.54).$$

Сравнивая эту оценку с интервальной оценкой примера 4.1 (см. (4.12)), видим, что замена неизвестной величины  $\sigma$  вычисляемой величиной  $s$  приводит к уменьшению точности интервальной оценки и увеличению длины доверительного интервала. Подставив вместо случайной величины  $\bar{X}_g$  ее конкретное значение  $\bar{x}_g = 1.5$ , получаем конкретное значение границ  $(0, 3)$ . ●

#### 4.4. Интервальные оценки дисперсии нормального распределения

Как и при построении интервальных оценок для математического ожидания, в данном случае также необходимо определить случайную величину, распределение которой было известно и включало оцениваемый параметр  $\sigma$ . В соответствии с теоремой 4.1 такой отправной точкой для построения доверительного интервала может быть случайная величина  $\frac{nD_g}{\sigma^2}$ , распределенная по закону  $\chi^2$  с  $(n-1)$  степенями свободы. Заметим, что доверительные интервалы, построенные для параметра  $a$ , вообще говоря, можно было выбрать несимметричными относительно  $\bar{X}_g$  и это не противоречило бы определению интервальной оценки. Но такой выбор интервала, когда в его середине лежит состоятельная и несмещенная оценка параметра, являлся предпочтительным. В данном случае целесообразно выбрать два предела  $\chi_{лев, \gamma}^2$  и  $\chi_{пр, \gamma}^2$  так, что

$$P(\chi_{n-1}^2 < \chi_{лев, \gamma}^2) = P(\chi_{n-1}^2 > \chi_{пр, \gamma}^2) = \frac{\alpha}{2},$$

где  $\alpha = 1 - \gamma$ ,  $\gamma$  – надежность интервальной оценки.

Следовательно,  $\chi_{лев,\gamma}^2$  – квантиль  $\chi_{n-1}^2$ -распределения уровня  $\alpha/2$ ,  $\chi_{np,\gamma}^2$  – уровня  $1-\alpha/2$ . Тогда имеет место равенство  $P\left(\chi_{лев,\gamma}^2 < \frac{nD_6}{\sigma^2} < \chi_{np,\gamma}^2\right) = \gamma$ , а интервал

$$\left(\frac{nD_6}{\chi_{np,\gamma}^2}, \frac{nD_6}{\chi_{лев,\gamma}^2}\right) \quad (4.16)$$

является интервальной оценкой для  $\sigma^2$  надежности  $\gamma$ .

Так как  $D_6 = (n-1)S^2/n$ , то  $nD_6 = (n-1)S^2$  и интервал

$$\left(\frac{n-1}{\chi_{np,\gamma}^2} S^2, \frac{n-1}{\chi_{лев,\gamma}^2} S^2\right) \quad (4.17)$$

является также интервальной оценкой для дисперсии  $\sigma^2$  надежности  $\gamma$ .

Заметим, что границы интервалов (4.16), (4.17) являются случайными величинами (почему?) и с вероятностью  $\gamma$  можно утверждать, что интервалы (4.16), (4.17) накроют неизвестную дисперсию  $\sigma^2$ .

♦ **Пример 4.3.** По выборке объема  $n = 20$  из нормально распределенной генеральной совокупности вычислено значение дисперсии выборки  $d_6 = 1.5$ . Построить интервальную оценку для параметра  $\sigma^2$  надежности  $\gamma = 0.96$ .

*Решение.* Значения  $\chi_{лев,\gamma}^2$ ,  $\chi_{np,\gamma}^2$  находим из условий:

$$P(\chi_{19}^2 < \chi_{лев,\gamma}^2) = 0.02; \quad P(\chi_{19}^2 < \chi_{np,\gamma}^2) = 0.98.$$

Эти условия означают, что  $\chi_{лев,\gamma}^2$  есть квантиль  $\chi^2$ -распределения с 19 степенями свободы уровня 0.02, а  $\chi_{np,\gamma}^2$  – квантиль уровня

0.98. По табл. ПЗ квантилей  $\chi^2$ -распределения находим

$$\chi_{лев,\gamma}^2 = 8.6; \quad \chi_{np,\gamma}^2 = 33.7.$$

Тогда интервальная оценка (4.16) принимает вид

$$(0.59D_6, 2.33D_6).$$

Подставляя вычисленное значение  $d_6 = 1.5$  случайной величины  $D_6$ , получаем

$$0.89 < \sigma^2 < 3.488. \quad \odot$$

#### 4.5. Интервальная оценка вероятности события

В п. 3.4 было показано, что "хорошей" точечной оценкой вероятности  $p$  события является частность  $p^* = m/n$  (см. (3.17)), где  $n$  – общее число независимых испытаний, в каждом из которых событие  $A$  может произойти с вероятностью  $p$ , а  $m$  – число испытаний, в которых произошло событие  $A$ .

Зададимся надежностью интервальной оценки  $\gamma$  и найдем числа  $p_{лев,\gamma}$ ,  $p_{np,\gamma}$  такие, чтобы выполнялось соотношение

$$P(p_{лев,\gamma} < p < p_{np,\gamma}) = \gamma. \quad (4.18)$$

Интервальную оценку построим для двух случаев: когда число испытаний  $n$  сравнительно велико ( $np > 10, n > 30$ ) и для малого числа испытаний.

**Интервальная оценка вероятности при большом числе испытаний.** Если  $np > 10, n > 30$ , то распределение случайной величины  $p^* = \frac{m}{n}$  можно аппроксимировать нормальным распределением  $N(p, \sqrt{pq/n})$ . Следовательно, при этих же условиях распределение величины  $\frac{(p^* - p)}{\sqrt{pq/n}}$  близко к нормальному с нулевым математическим ожиданием и единичной дисперсией, т.е.

$$\frac{p^* - p}{\sqrt{pq/n}} = N(0,1).$$

По аналогии с (4.8) найдем такое число  $x_\gamma$ , для которого справедливо равенство

$$P\left(-x_\gamma < \frac{p^* - p}{\sqrt{pq/n}} < x_\gamma\right) = \gamma. \quad (4.19)$$

Это число является корнем уравнения

$$\Phi(x_\gamma) = \gamma/2,$$

где  $\Phi(x)$  – функция Лапласа, и корень может быть найден с помощью табл. П1.

Неравенство, стоящее в скобках выражения (4.19), разрешим относительно  $p$ . Для этого неравенство перепишем в виде эквивалентного неравенства

$\left|\frac{p^* - p}{\sqrt{pq/n}}\right| < x_\gamma$ . Возведем в квадрат, в результате получим  $(p^* - p)^2 < \frac{p(1-p)}{n} x_\gamma^2$ . Далее, возведя в квадрат

$(p^* - p)$  и перенеся все члены влево, получим

$$\left(1 + \frac{x_\gamma^2}{n}\right)p^2 - \left(2p^* + \frac{x_\gamma^2}{n}\right)p + p^{*2} < 0.$$

Корни  $p_1$  и  $p_2$  квадратного трехчлена, стоящего в правой части неравенства, определяются выражениями

$$p_1 = \frac{p^* + x_\gamma^2/(2n) - x_\gamma \sqrt{p^*(1-p^*)/n + x_\gamma^2/(4n^2)}}{1 + x_\gamma^2/n}; \quad (4.20)$$

$$p_2 = \frac{p^* + x_\gamma^2/(2n) + x_\gamma \sqrt{p^*(1-p^*)/n + x_\gamma^2/(4n^2)}}{1 + x_\gamma^2/n}. \quad (4.21)$$

Корни этого уравнения и являются границами интервальной оценки (4.18)

$$P_{лев,\gamma} = p_1; \quad P_{пр,\gamma} = p_2. \quad (4.22)$$

Если  $n \gg 100$ , то для вычисления  $p_1, p_2$  можно использовать приближенные формулы:

$$p_1 \approx p^* - x_\gamma \sqrt{p^*(1-p^*)/n}; \quad p_2 \approx p^* + x_\gamma \sqrt{p^*(1-p^*)/n}. \quad (4.23)$$

Видно, что границы интервала (4.18) являются случайными величинами и конкретные значения границ получаются в результате подстановки наблюдаемого значения случайной величины  $p^*$ .

♦ **Пример 4.4.** Событие  $A$  в серии из  $n = 100$  испытаний произошло  $m = 78$  раз. Построить интервальную оценку для вероятности  $p$  события с надежностью  $\gamma = 0.9$ .

*Решение.* Значение точечной оценки вероятности  $p$  равно  $p^* = 78/100 = 0.78$ . По табл. П1 определяем  $x_\gamma = 1.64$  и вычисляем по формулам (4.20), (4.21) значения  $p_1, p_2$  при  $p^* = 0.78$ :  $p_1 = 0.705$ ,  $p_2 = 0.848$ . Таким образом, получили реализацию доверительного интервала (0.705, 0.848) для вероятности  $p$  события  $A$ . ●

**Интервальная оценка вероятности при малом числе испытаний.** При малом числе испытаний  $n$  предположение о приближенном распределении случайной величины  $m$  по нормальному закону  $m = N(np, \sqrt{npq})$  становится несправедливым. Для описания распределения величины  $m$  необходимо использовать формулу Бернулли:

$$P(m = x) = C_n^x p^x (1-p)^{n-x}, \quad x = 0, 1, \dots, n.$$

Можно показать, что граничные точки интервальной оценки (4.18) являются решениями следующих нелинейных уравнений:

$$\sum_{x=0}^{m-1} C_n^x p_{лев,\gamma}^x (1 - p_{лев,\gamma})^{n-x} = \frac{1+\gamma}{2}; \quad (4.24)$$

$$\sum_{x=0}^m C_n^x p_{пр,\gamma}^x (1 - p_{пр,\gamma})^{n-x} = \frac{1-\gamma}{2}, \quad (4.25)$$

где  $\gamma$  – надежность интервальной оценки. Вновь заметим, что решения  $p_{лев,\gamma}, p_{пр,\gamma}$  этих уравнений являются случайными величинами (почему?) и только при подстановке конкретного значения  $m$  (количество испытаний, в которых появилось событие  $A$ ) будут получены конкретные значения граничных точек интервальной оценки (4.18).

Корни уравнений (4.24), (4.25) могут быть найдены одним из известных численных методов решения нелинейных уравнений. Кроме этого, существуют специальные таблицы для нахождения  $p_{лев,\gamma}, p_{пр,\gamma}$ , удовлетворяющих уравнениям (4.24), (4.25) по заданным  $n, m - n, \gamma$ . Фрагмент этих таблиц представлен в приложении (табл. П4).

♦ **Пример 4.5.** В пяти испытаниях событие  $A$  произошло три раза. Построить интервальную оценку для вероятности  $p$  события  $A$  с надежностью  $\gamma = 0.95$ .

*Решение.* Из условий примера имеем  $n = 5, m = 3, \gamma = 0.95$ . По табл. П4 находим  $p_{лев,\gamma} = 0.147, p_{пр,\gamma} = 0.947$ , а интервальная оценка определяется как (0.147, 0.947).

Сравнивая интервальные оценки примеров 4.4, 4.5, видим, что длина доверительного интервала для примера 4.5 (равная 0.8) существенно больше длины доверительного интервала примера 4.4 (0.143). Это является следствием разного объема выборок ( $n = 5$  и  $n = 100$ ) и различных дисперсий случайной величины  $p^* = m/n$ .



#### 4.6. Вычисление границ доверительных интервалов в Excel

Границы доверительных интервалов зависят от некоторой величины, которая зависит от распределения точечной оценки и до-

верительной вероятности. Эта величина находится по специальным таблицам. Поэтому часто возникает необходимость интерполяции или экстраполяции табличных данных и, следовательно, требуются дополнительные вычисления. В табличном процессоре Excel определены функции, позволяющие вычислять величины, входящие в интервальные оценки для различных числовых характеристик случайной величины.

**Вычисление величины  $x_\gamma$ ,** входящей в доверительный интервал (4.11):

$$\left[ \bar{X}_e - \frac{x_\gamma \sigma}{\sqrt{n}}, \bar{X}_e + \frac{x_\gamma \sigma}{\sqrt{n}} \right]. \quad (4.26)$$

Величина  $x_\gamma$  является корнем нелинейного уравнения (4.10) и вычисляется с помощью функции НОРМСТОБР:

$$x_\gamma = \text{НОРМСТОБР}((\gamma + 1)/2),$$

где  $\gamma$  – надежность интервальной оценки (4.26).

**Вычисление величины  $x_\gamma \sigma / \sqrt{n}$**  осуществляется с помощью функции ДОВЕРИТ:

$$\Delta_{\bar{X}_e} = x_\gamma \sigma / \sqrt{n} = \text{ДОВЕРИТ}(\alpha; \sigma; n),$$

где  $\alpha = 1 - \gamma$ ,  $\sigma$  – известное среднеквадратичное отклонение,  $n$  – объем выборки. Тогда интервальную оценку (4.26) можно записать в виде  $\left[ \bar{X}_e - \Delta_{\bar{X}_e}, \bar{X}_e + \Delta_{\bar{X}_e} \right]$ .

**Вычисление величины  $t(\gamma, n)$ ,** входящей в доверительный интервал

$$\left[ \bar{X}_e - \frac{t(\gamma, n) \cdot \sqrt{D_e}}{\sqrt{n-1}}, \bar{X}_e + \frac{t(\gamma, n) \cdot \sqrt{D_e}}{\sqrt{n-1}} \right],$$

осуществляют с использованием функции СТЬЮДРАСПОБР, обращение к которой имеет вид:

$$t(\gamma, n) = \text{СТЫЮДРАСПОБР}(\alpha; n),$$

где  $\alpha = 1 - \gamma$ ,  $n$  – **число степеней свободы** (обратите на это внимание).

**Вычисление величин**  $\chi_{лев,\gamma}^2$ ,  $\chi_{пр,\gamma}^2$ , входящих в доверительный интервал (4.17), для дисперсии  $\sigma^2$ :

$$\left[ \frac{n-1}{\chi_{пр,\gamma}^2} S^2, \frac{n-1}{\chi_{лев,\gamma}^2} S^2 \right],$$

где  $S^2$  – исправленная дисперсия. Используется функция ХИ2ОБР:

$$\chi_{лев,\gamma}^2 = \text{ХИ2ОБР}\left(1 - \frac{\alpha}{2}; n\right);$$

$$\chi_{пр,\gamma}^2 = \text{ХИ2ОБР}(\alpha/2; n),$$

где  $\alpha = 1 - \gamma$ ,  $\gamma$  – надежность интервальной оценки.

**Задание 4.1.** Используя функции Excel, вычислите интервальные оценки для примеров 4.1 и 4.2. ♥

**Задание 4.2.** Используя функции Excel, вычислите интервальные оценки для примера 4.3. ♥

## 5. ПРОВЕРКА СТАТИСТИЧЕСКИХ ГИПОТЕЗ

### 5.1. Понятие статистической гипотезы. Основные этапы проверки гипотезы

Прежде чем перейти к математическим формулировкам, рассмотрим один пример.

Результаты многолетних статистических исследований показали, что для населения некоторого региона вероятность предрасположения к данному заболеванию  $R$  равна  $p_0 = 0.1$ . После строительства в этом регионе химического предприятия была проведена выборочная проверка населения. Из 1000 обследованных у 120 человек были обнаружены признаки заболевания  $R$ . Можно ли утверждать: а) полученные данные не противоречат предположению, что строительство не повлияло на уровень заболевания  $R$ , или б) изменение экологической обстановки после строительства комбината повлияло на распространение заболевания  $R$ ? Приведенный пример является типичной задачей проверки статистической гипотезы. Под *статистической гипотезой* понимается всякое высказывание о генеральной совокупности (случайной величине  $X$ ), проверяемое по выборочной совокупности (по результатам наблюдений). В нашем примере высказывание формулируется в терминах вероятности  $p_0$  события  $A = \{\text{наличие у человека заболевания } R\}$ . Не располагая сведениями о всей генеральной совокупности, высказанную гипотезу сопоставляют по определенным правилам с выборочными данными и делают вывод о том, можно принять гипотезу или нет. Эта процедура сопоставления называется *проверкой гипотезы*.

Рассмотрим этапы проверки гипотезы и используемые при этом понятия.

**Э т а п 1.** Располагая выборочными данными и руководствуясь конкретными условиями рассматриваемой задачи, формулируют гипотезу  $H_0$ , которую называют *основной* или *нулевой*, и гипотезу  $H_1$ , *конкурирующую* с гипотезой  $H_0$ . Гипотезу  $H_1$  называют также *альтернативной*.

Термин "конкурирующая" означает, что являются взаимоисключающими следующие два события:

- по выборке принимается решение о справедливости для ге-

неральной совокупности гипотезы  $H_0$ ;

- по выборке принимается решение о справедливости для генеральной совокупности гипотезы  $H_1$ .

Вернемся к нашему примеру. Обозначим через  $A$  событие, состоящее в том, что случайно выбранный человек в данном регионе предрасположен к заболеванию  $R$ . До строительства химического предприятия вероятность события  $A$  была равна 0.1. В качестве гипотезы  $H_0$  рассмотрим гипотезу о том, что после строительства химического предприятия вероятность события  $A$  не изменилась. Таким образом, если  $p_1$  – вероятность события  $A$  после строительства предприятия, то в качестве нулевой (основной) гипотезы принимается

$$H_0 : p_1 = p_0.$$

Учитывая, что: а) строительство комбината вряд ли улучшило экологическую обстановку в регионе; б) при выборке из 1000 человек у 120 человек обнаружено заболевание  $R$ , что соответствует относительной частоте  $p^* = 120/1000 = 0.12 > 0.1$ , в качестве альтернативной гипотезы примем:

$$H_1 : p_1 > p_0.$$

Э т а п 2. Задается вероятность  $\alpha$ , которую называют *уровнем значимости*. Эта вероятность имеет следующий смысл.

Решение о том, можно ли считать высказывание  $H_0$  справедливым для генеральной совокупности, принимается по выборочным данным, т.е. по ограниченному объему информации. Следовательно, это решение может быть ошибочным. При этом может иметь место ошибка двух родов:

- *ошибка первого рода* совершается при отклонении гипотезы  $H_0$  (т.е. принимается альтернативная  $H_1$ ), тогда как на самом деле гипотеза  $H_0$  верна; вероятность такой ошибки обозначим  $P(H_1 / H_0)$ ;
- *ошибка второго рода* совершается при принятии гипотезы  $H_0$ , тогда как на самом деле высказывание  $H_0$  неверно и следовало бы принять гипотезу  $H_1$ ; вероятность ошибки второго рода обозначим как

$$\beta = P(H_0 / H_1). \quad (5.1)$$

Тогда уровень значимости  $\alpha$  определяет ошибку первого рода, т.е.

$$\alpha = P(H_1 / H_0). \quad (5.2)$$

Поэтому вероятность  $\alpha$  задается малым числом, поскольку это вероятность ошибочного высказывания. При этом обычно используются стандартные значения: 0.05; 0.01; 0.005. Например,  $\alpha = 0.05$  означает следующее: если гипотезу  $H_0$  проверять по каждой из 100 выборок одинакового объема, то в среднем в 5 случаях из 100 совершим ошибку первого рода.

Обратим внимание на то, что в результате проверки гипотезы  $H_0$  могут быть приняты *правильные решения* двух следующих видов:

- принимается гипотеза  $H_0$  тогда, когда она верна (т.е.  $H_0$  имеет место в генеральной совокупности); вероятность этого решения равна  $P(H_0 / H_0) = 1 - \alpha$  (почему?);
- не принимается гипотеза  $H_0$  (т.е. принимается гипотеза  $H_1$ ) тогда, когда и на самом деле она неверна (т.е. справедлива гипотеза  $H_1$ ), вероятность этого решения равна (почему?)

$$P(H_1 / H_1) = 1 - \beta. \quad (5.3)$$

Э т а п 3. Определяют величину  $K$  такую, что: а) ее значения зависят от выборочных данных  $x_1, x_2, \dots, x_n$ , т.е.  $K = K(x_1, x_2, \dots, x_n)$ ; б) будучи величиной случайной (в силу случайности выборки  $x_1, \dots, x_n$ ), величина  $K$  подчиняется при выполнении гипотезы  $H_0$  некоторому известному закону распределения; в) ее значения позволяют судить о расхождении гипотезы  $H_0$  с выборочными данными. Величину  $K$  называют *критерием*.

Обратимся к нашему примеру. Пусть  $S_{1000}$  – количество обследуемых, предрасположенных к заболеванию  $R$  в выборке из 1000 человек. Если гипотеза  $H_0$  верна, т.е.  $p_1 = p_0 = 0.1$ , то случайная величина  $S_{1000}$  распределена по биномиальному закону и ее числовые характеристики равны  $M(S_{1000}) = 100$ ,  $D(S_{1000}) = 90$  (почему?). С другой стороны, ее распределение близко к нормальному. Поэтому случайная величина

$$K = \frac{S_{1000} - 100}{9.487} \quad (5.4)$$

распределена по закону, близкому к нормальному  $N(0,1)$ .

Заметим, что если вероятность события  $A$  возросла после строительства химического комбината, то случайная величина  $K$  преимущественно будет принимать положительные значения (почему?) и это может трактоваться в пользу принятия гипотезы  $H_1$ . Видно, что величина (5.4) удовлетворяет требованиям а), б), в) и может быть принята при проверке гипотезы  $H_0: p_1 = p_0$  при альтернативной  $H_1: p_1 > p_0$ .

Э т а п 4. В области всевозможных значений критерия  $K$  выделяют подобласть  $\omega$ , называемую *критической областью*. Значения критерия, попавшие в критическую область, свидетельствуют о существенном расхождении выборки с гипотезой  $H_0$ . Поэтому руководствуются следующим правилом: если вычисленное по выборке значение критерия попадает в критическую область  $\omega$ , то гипотеза  $H_0$  отвергается и принимается альтернативная  $H_1$ . При этом следует помнить, что такое решение может быть ошибочным – на самом деле гипотеза  $H_0$  может быть справедливой. Таким образом, ориентируясь на критическую область, можно совершить ошибку первого рода, вероятность которой задана заранее и равна  $\alpha$ . Отсюда вытекает следующее требование к критической области  $\omega$ :

*Вероятность принятия критерием  $K$  значения из критической области  $\omega$  при справедливости гипотезы  $H_0$  должна быть равна  $\alpha$ , т.е.*

$$P(K \in \omega) = \alpha. \quad (5.5)$$

Однако критическая область определяется равенством (5.5) неоднозначно. Пусть  $p_K(x)$  является плотностью распределения критерия  $K$ . Тогда нетрудно увидеть, что на оси  $X$  существует бесчисленное множество интервалов таких, что площади построенных на них криволинейных трапеций, ограниченных сверху кривой  $p_K(x)$ , равны  $\alpha$ . Поэтому кроме требования (5.5) выдвигается следующее: критическая область  $\omega$  должна быть расположена так, чтобы при заданной вероятности  $\alpha$  – ошибки первого рода вероятность  $\beta$  – ошибки второго рода (см. (5.1)) была минимальной.

Обычно этому требованию удовлетворяют три случая расположения критической области (в зависимости от вида нулевой и альтернативной гипотез, формы и распределения критерия  $K$ ):

- правосторонняя критическая область (рис. 5.1,а), состоящая из интервала  $(x_{np,\alpha}, +\infty)$ , где точка  $x_{np,\alpha}$  определяется из условия

$$P(K > x_{np,\alpha}) = \alpha \quad (5.6)$$

и называется *правосторонней критической точкой*;

- левосторонняя критическая область (см. рис. 5.1,б) состоит из интервала  $(-\infty, x_{лев,\alpha})$ , где  $x_{лев,\alpha}$  определяется из условия

$$P(K < x_{лев,\alpha}) = \alpha \quad (5.7)$$

и называется *левосторонней критической точкой*;

- двусторонняя критическая область (см. рис. 5.1,в), состоящая из двух интервалов:  $(-\infty, x_{лев,\alpha/2})$ ,  $(x_{np,\alpha/2}, +\infty)$ , где точки  $x_{лев,\alpha/2}$ ,  $x_{np,\alpha/2}$  определяются из условий

$$P(K < x_{лев,\alpha/2}) = \alpha/2; \quad P(K > x_{np,\alpha/2}) = \alpha/2. \quad (5.8)$$

Вернемся к нашему примеру. Так как альтернативная гипотеза имеет вид  $H_1: p_1 > p_0$ , то принимается правосторонняя критическая область (см. рис. 5.1,а). Задаваясь  $\alpha = 0.005$ , определяем  $x_{np,\alpha}$  из уравнения (5.6).

При справедливости гипотезы  $H_0$  критерий  $K$ , определяемый выражением (5.4), имеет нормальное распределение  $N(0,1)$ , и, следовательно, по таблице функции Лапласа  $\Phi(x)$  (по табл. П1) необходимо найти такое  $x_{np,\alpha}$ , что  $\Phi(x_{np,\alpha}) = 0.495$ . Это значение равно 2.58. Тогда вероятность того, что критерий  $K$  при справедливости гипотезы  $H_0$  примет значение больше 2.58, равна

$$P(K > 2.58) = P(2.58 < N(0.1) < \infty) = \Phi(\infty) - \Phi(2.58) = 0.005.$$

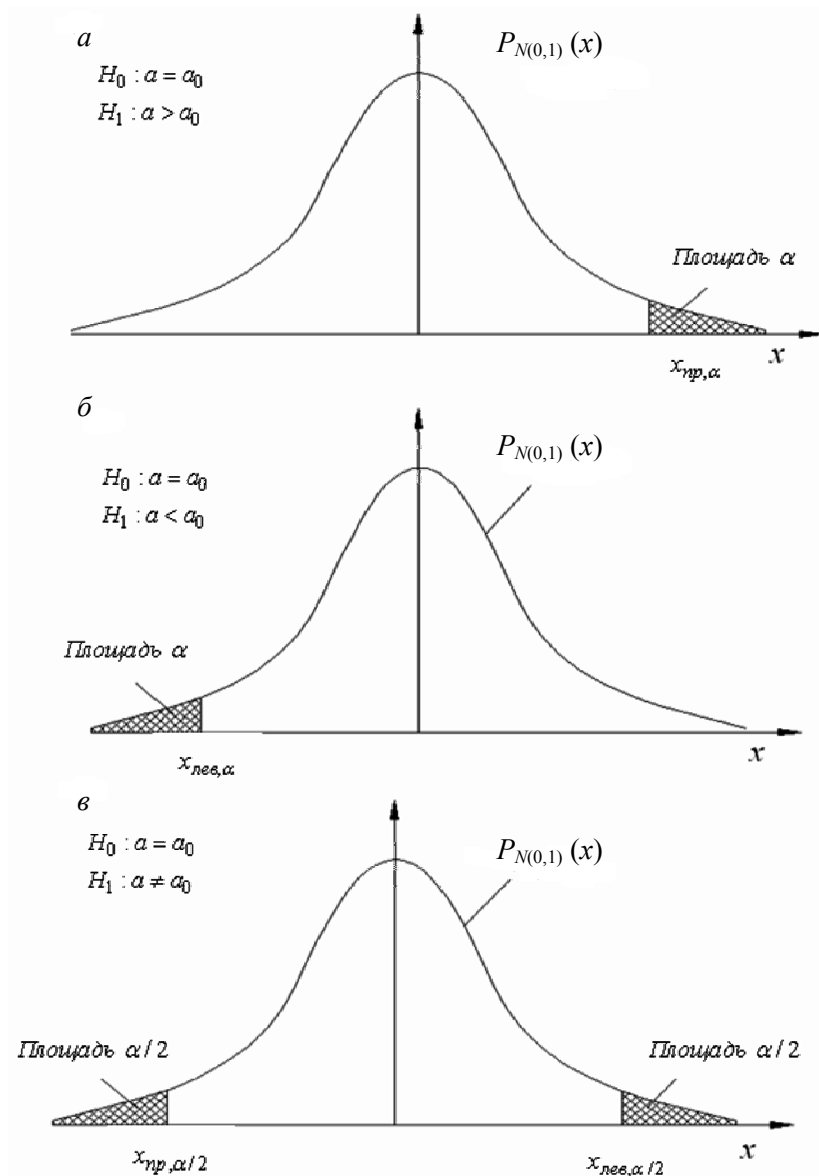


Рис. 5.1. Три вида критических областей при проверке статистических гипотез

Выбор критической области из условия минимума вероятности ошибки второго рода эквивалентен выбору критической области из условия максимума величины

$$m = 1 - \beta,$$

называемой мощностью критерия  $K$  и равной вероятности  $P(H_1 / H_1)$  принятия гипотезы  $H_1$  при справедливости гипотезы  $H_1$ . Поясним понятие мощности критерия следующим примером.

Предположим, что если верна гипотеза  $H_0$ , то критерий  $K$  распределен по нормальному закону  $N(5,3)$  (т.е. математическое ожидание  $a = 5$ , дисперсия  $\sigma^2 = 9$ ), а если верна конкурирующая гипотеза  $H_1$ , то критерий распределен по закону  $N(15,3)$ . Требуется вычислить мощность критерия  $m_1$ , когда в качестве критической рассматривается область больших значений, и мощность  $m_2$ , когда в качестве критической рассматривается область больших по модулю значений. Уровень значимости  $\alpha$  возьмем 0.05. В первом случае границу правосторонней критической области найдем из условия  $P(N(5,3) > x_{np,\alpha}) = 0.05$ , поэтому

$$P(N(5,3) > x_{np,\alpha}) = P(x_{np,\alpha} < N(5,3) < \infty) = \frac{1}{2} - \Phi\left(\frac{x_{np,\alpha} - 5}{3}\right) = 0.05.$$

Значит,  $\Phi\left(\frac{x_{np,\alpha} - 5}{3}\right) = 0.45$ . По таблицам значений функции  $\Phi(x)$

находим, что  $\frac{x_{np,\alpha} - 5}{3} = 1.64$ . Поэтому границы правосторонней

критической области  $x_{np,\alpha} = 9.92$ . Чтобы вычислить ошибку второго рода  $\beta_1$ , нужно найти вероятность попадания критерия в область допустимых значений  $(-\infty, 9.92)$  при условии, что гипотеза  $H_0$  неверна. В этом случае считается справедливой гипотеза  $H_1$ , а критерий будет распределен по закону  $N(15,3)$ . Значит,

$$\begin{aligned} \beta_1 &= P(N(15,3) < 9.92) = 0.5 + \Phi\left(\frac{9.92 - 15}{3}\right) = 0.5 - \Phi(1.69) = \\ &= 0.5 - 0.4545 = 0.0455 \end{aligned}$$

и мощность критерия  $m_1 = 1 - \beta_1 = 0.955$ .



Во втором случае правая граница критической области  $x_{np,\alpha/2}$  вычисляется из условия  $P(N(5,3) > x_{np,\alpha}) = 0.025$ . Поэтому  $\frac{x_{np,\alpha/2} - 5}{3} = 1.96$ . Значит,  $x_{np,\alpha/2} = 10.88$ . Левая граница критической области с точкой  $x_{np,\alpha/2}$  симметрична относительно точки  $x = 5$ , т.е. левая граница  $x_{np,\alpha/2} = 5 - 5.88 = 0.88$ . Тогда вероятность ошибки  $\beta_2$  составит

$$\begin{aligned}\beta_2 &= P(-0.88 < N(15,3) < 10.88) = \Phi\left(\frac{10.88-15}{3}\right) - \Phi\left(\frac{-0.88-15}{3}\right) = \\ &= \Phi(5.29) - \Phi(1.37) = 0.5 - 0.41147 = 0.0853.\end{aligned}$$

Поэтому мощность критерия во втором случае равна  $m_2 = 1 - \beta_2 = 1 - 0.0853 = 0.9147$ . Значит, односторонняя критическая область больших значений является предпочтительной.

**Э т а п 5.** В формулу критерия  $K$ , который является функцией  $n$  случайных величин  $X_1, X_2, \dots, X_n$ , подставляются выборочные значения  $x_1, x_2, \dots, x_n$  и подсчитывается числовое значение критерия  $K_{наб}$ .

Если  $K_{наб}$  попадает в критическую область  $\omega$ , то гипотеза  $H_0$  отвергается и принимается гипотеза  $H_1$ . При этом можно допустить ошибку первого рода с вероятностью  $\alpha$ . Если  $K_{наб}$  не попадает в критическую область, гипотеза  $H_0$  не отвергается. Однако это не означает, что  $H_0$  является единственной подходящей гипотезой: просто  $H_0$  не противоречит результатам наблюдений; возможно, таким же свойством наряду с  $H_0$  могут обладать и другие гипотезы.

Вновь обратимся к нашему примеру. Напомним, что из обследованных 1000 человек признаки заболевания  $R$  были обнаружены у 120 человек, т.е.  $S_{1000} = 120$ . Подставляя это выборочное значение в формулу (5.4), получаем

$$K_{наб} = \frac{120-100}{9.487} = 2.108.$$

Правосторонняя критическая точка ранее была определена как  $x_{np,\alpha/2} = 2.58$ . Так как  $2.108 < 2.58$ , то можно принять гипотезу  $H_0: p_1 = p_0$ , а полученные расхождения между теоретической вероятностью  $p_0 = 0.1$  и наблюдаемой частностью 0.120 считать допустимыми на уровне значимости  $\alpha = 0.005$ .

Если бы количество человек с признаками заболевания  $R$  составило 130 (из 1000 обследованных), то  $K_{наб} = \frac{130-100}{9.487} = 3.162$ .

В этом случае случайная величина  $K$  приняла значение из критической области, т.е. произошло событие  $K > x_{np,\alpha/2}$ , которое практически невозможно, если гипотеза  $H_0$  справедлива. Поэтому следует отвергнуть гипотезу  $H_0$  в пользу альтернативной гипотезы  $H_1: p_1 > p_0$ .

## 5.2. Проверка гипотезы о числовом значении математического ожидания нормального распределения

Полагаем, что  $X$  является случайной величиной, имеющей нормальное распределение с параметрами  $a$  и  $\sigma$ , т.е.  $X = N(a, \sigma)$ , причем числовое значение  $a$  неизвестно.

Дать точный ответ на вопрос, каково численное значение неизвестного параметра  $a$ , по выборочной совокупности, нельзя. Поэтому поступают следующим образом. Полагая, что наблюдения  $X_1, X_2, \dots, X_n$  независимы, вычисляют значение выборочной оценки  $\bar{X}_e$ , которое дает приближенные представления об  $a$ . Затем приступают к проверке гипотез о числовых значениях неизвестного параметра  $a$ .

**Проверка гипотезы о числовом значении математического ожидания при известной дисперсии.** Предполагается, что  $X = N(a, \sigma)$ , причем значение математического ожидания  $a$  неизвестно, а числовое значение дисперсии  $\sigma^2$  известно.

Выдвинем гипотезу  $H_0$  о том, что неизвестный параметр  $a$  равен числу  $a_0$ . Возможны три случая: 1) параметр  $a$  равен числу  $a_1$ ,

которое больше числа  $a_0$  (т.е.  $a > a_0$ ); 2) параметр  $a$  равен числу  $a_1$ , которое не равно  $a_0$  (т.е.  $a \neq a_0$ ); 3) параметр  $a$  равен числу  $a_1$ , которое меньше  $a_0$  (т.е.  $a < a_0$ ). Для случаев 1, 2 рассмотрим этапы проверки гипотезы  $H_0$ , приведенные в п. 5.1.

### Случай 1

Э т а п 1. Сформулируем нулевую гипотезу

$$H_0 : a = a_0 \quad (5.9)$$

и альтернативную

$$H_1 : a = a_1 > a_0. \quad (5.10)$$

Э т а п 2. Зададимся уровнем значимости  $\alpha$ .

Э т а п 3. В качестве критерия возьмем величину

$$K = \frac{\bar{X}_g - a_0}{\sigma/\sqrt{n}}, \quad (5.11)$$

значение которой зависит от выборочных данных (почему?), является случайной величиной и при выполнении гипотезы (5.9) подчиняется нормальному распределению  $N(0,1)$ , т.е.

$$K = \frac{\bar{X}_g - a_0}{\sigma/\sqrt{n}} = N(0,1). \quad (5.12)$$

Э т а п 4. Построим критическую область  $\omega$ , т.е. область таких значений критерия  $K$ , при которых гипотеза  $H_0$  отвергается. Если нулевая и альтернативная гипотезы имеют вид (5.9), (5.10) соответственно, а критерий (5.11) – вид  $K = N(0,1)$ , то критическая область будет правосторонней: ее образует интервал  $(x_{np,\alpha}, +\infty)$ , где  $x_{np,\alpha}$  определяется из условия (5.6), которое с учетом (5.12) записывается как

$$P(N(0,1) > x_{np,\alpha}) = \alpha.$$

Остановимся на методике вычисления  $x_{np,\alpha}$  (которая будет использована в дальнейшем для других критических точек). Вероятность события  $N(0,1) \leq x_{np,\alpha}$  можно представить как

$$\int_{-\infty}^0 p_{N(0,1)}(x)dx + \int_0^{x_{np,\alpha}} p_{N(0,1)}(x)dx = \frac{1}{2} + \Phi(x_{np,\alpha}),$$

где  $p_{N(0,1)}(x)$  – плотность нормального распределения  $N(0,1)$ ;  $\Phi(x)$  – функция Лапласа (см. табл. П1). Следовательно, вероятность противоположного события  $N(0,1) > x_{np,\alpha}$  выражается в виде  $1 - \left[\frac{1}{2} + \Phi(x_{np,\alpha})\right] = \frac{1}{2} - \Phi(x_{np,\alpha})$ , и эта вероятность должна быть равна  $\alpha$ . Таким образом, приходим к уравнению

$$\Phi(x_{np,\alpha}) = \frac{1}{2} - \alpha.$$

Воспользовавшись табл. П1, находим значение  $x_{np,\alpha}$ , удовлетворяющее этому уравнению. Критическая область изображена на рис. 5.1,а.

Э т а п 5. Используя вместо  $X_1, X_2, \dots, X_n$  конкретные числа, находим  $\bar{x}_g$  (см. (2.10)), а затем численное значение  $K_{наб}$  критерия (5.11). Если  $K_{наб} > x_{np,\alpha}$ , то гипотеза  $H_0$  (5.9) отвергается и принимается гипотеза  $H_1$  (5.10). Напомним, что, поступая таким образом, мы можем совершить ошибку первого рода. Вероятность такой ошибки равна  $\alpha$ .

### Случай 2

Э т а п 1. Сформулируем нулевую гипотезу

$$H_0 : a = a_0 \quad (5.13)$$

и альтернативную

$$H_1 : a \neq a_0. \quad (5.14)$$

Э т а п 2. Зададимся уровнем значимости  $\alpha$ .

Э т а п 3. В качестве критерия, как и в случае 1, возьмем величину (5.11), которая при справедливости гипотезы (5.13) удовлетворяет распределению  $N(0,1)$ .

Э т а п 4. Если нулевая и альтернативная гипотезы имеют соответственно вид (5.13), (5.14), а критерий определяется выражением (5.12), то критическая область будет двусторонней: ее образуют интервалы  $(-\infty, x_{лев, \alpha/2})$ ,  $(x_{пр, \alpha/2}, +\infty)$ , где критические точки  $x_{пр, \alpha/2}$ ,  $x_{лев, \alpha/2}$  находятся из условия (5.8), которое, учитывая (5.12), запишется так:

$$P(N(0,1) < x_{лев, \alpha/2}) = \frac{\alpha}{2}; \quad P(N(0,1) > x_{пр, \alpha/2}) = \frac{\alpha}{2}. \quad (5.15)$$

Из рис. 5.1,б видно, что

$$\Phi(x_{пр, \alpha/2}) = \frac{(1-\alpha)}{2}. \quad (5.16)$$

Воспользовавшись табл. П1, находим решение этого уравнения  $x_{пр, \alpha/2}$ . В силу симметричности функции плотности распределения  $N(0,1)$  имеем

$$x_{лев, \alpha/2} = -x_{пр, \alpha/2}.$$

Э т а п 5. Находим числовое значение  $K_{наб}$  критерия (5.11). Если  $K_{наб}$  попадает в интервал  $(-\infty, x_{лев, \alpha/2})$  или  $(x_{пр, \alpha/2}, +\infty)$ , то гипотеза  $H_0$  (5.13) отвергается и принимается альтернативная (5.14). Поступая таким образом, можно с вероятностью  $\alpha$  допустить ошибку первого рода.

♦ **Пример 5.1.** По результатам  $n = 9$  замеров установлено, что среднее время изготовления детали  $\bar{x}_e = 52$  с. Предполагая, что время изготовления подчиняется нормальному распределению с дисперсией  $\sigma^2 = 9$  с<sup>2</sup>, решить на уровне значимости  $\alpha = 0.05$ :

а) можно ли принять 50 с в качестве нормативного времени (математического ожидания) изготовления детали;

б) можно ли принять за норматив 51 с?

*Решение.*

а) по условию задачи нулевая гипотеза  $H_0 : a = 50$  с. Так как  $\bar{x}_e = 52$  с, то в качестве альтернативной возьмем гипотезу  $H_1 : a > 50$  с, т.е. имеем случай 1 (см. (5.9), (5.10)) при  $a_0 = 50$  с. По изложенной схеме получаем  $x_{пр, \alpha} = 1.65$ . Подставляя в (5.11) исходные данные  $\bar{x}_e = 52$  с,  $\sigma = 3$ ,  $n = 9$ , получаем  $K_{наб} = \frac{52-50}{3/\sqrt{9}} = 2$ . Так

как число 2 попадает в критическую область  $(1.65, \infty)$ , то гипотеза  $H_0 : a = 50$  с отвергается и принимается  $H_1 : a > 50$  с;

б) здесь нулевая гипотеза  $H_0 : a = 51$  с, альтернативная  $H_1 : a > 51$  с. Снова имеет место случай 1 при  $a_0 = 51$  с. Так как  $K_{наб} = \frac{51-50}{3/\sqrt{9}} = 1$  не попадает в критическую область, то гипотеза  $H_0 : a = 51$  с не отвергается и в качестве норматива времени изготовления детали берем 51 с. ●

**Проверка гипотезы о числовом значении математического ожидания при неизвестной дисперсии.** В этом случае за основу проверки гипотезы

$$H_0 : a = a_0, \quad (5.17)$$

где  $a_0$  – заранее заданное число, положен критерий

$$K = \frac{\bar{X}_e - a_0}{S/\sqrt{n}}, \quad (5.18)$$

где  $\bar{X}_e$ ,  $S$  – случайные величины, вычисляемые по формулам (2.9) и (3.12). Этот критерий при выполнении гипотезы (5.17) имеет  $t$ -распределение с числом степеней свободы  $k = n - 1$ , т.е.

$$K = \frac{\bar{X}_e - a_0}{S/\sqrt{n}} = T_{n-1}, \quad (5.19)$$

где  $T_{n-1}$  – случайная величина, подчиняющаяся распределению Стьюдента (см. (4.5)).

Задаваясь уровнем значимости  $\alpha$ , построим критическую область для проверки гипотезы (5.17) при следующих альтернативных гипотезах.

### Случай 1

Альтернативная гипотеза

$$H_1 : a > a_0. \quad (5.20)$$

Критическая область является правосторонней: ее образует интервал  $(x_{np,\alpha}, +\infty)$ , где точка  $x_{np,\alpha}$  определяется из условия (5.6), которое с учетом (5.12) можно записать в виде

$$P(T_{n-1} > x_{np,\alpha}) = \alpha.$$

В табл. П2 приведены значения  $t(\gamma, n)$ , определяемые соотношением  $\int_{-t(\gamma,n)}^{t(\gamma,n)} P_T(x) dx = \gamma$ , где  $n$  – объем выборки, а не число степеней свободы. Так как функция плотности  $t$ -распределения симметрична относительно нуля, то искомая точка  $x_{np,\alpha}$  определяется как

$$x_{np,\alpha} = t(1 - 2\alpha, n). \quad (5.21)$$

Подставив в (5.18) конкретные значения  $\bar{X}_g$ ,  $S$ , получаем значение критерия  $K_{наб}$ . Если  $K_{наб} > x_{np,\alpha}$  (т.е. попадает в критическую область), то гипотеза (5.17) отвергается и принимается гипотеза (5.20). При этом возможна ошибка первого рода с вероятностью  $\alpha$ .

### Случай 2

Альтернативная гипотеза

$$H_1 : a \neq a_0. \quad (5.22)$$

Критическая область состоит из двух интервалов  $(-\infty, x_{лев,\alpha/2})$ ,  $(x_{np,\alpha/2}, +\infty)$ , где критические точки  $x_{лев,\alpha/2}$ ,  $x_{np,\alpha/2}$  определяются из условий (5.8), которые с учетом (5.19) можно записать в

виде  $P(T_{n-1} < x_{лев,\alpha/2}) = \alpha/2$ ;  $P(T_{n-1} > x_{np,\alpha/2}) = \alpha/2$ .

Обращаясь к табл. П2, находим

$$x_{лев,\alpha/2} = -t(1 - \alpha, n); \quad x_{np,\alpha/2} = t(1 - \alpha, n). \quad (5.23)$$

Подставляя в (5.18) конкретные значения величин  $\bar{X}_g$ ,  $S$ , получаем значение критерия  $K_{наб}$ . Если  $K_{наб}$  попадает в интервал  $(-\infty, x_{лев,\alpha/2})$  или  $(x_{np,\alpha/2}, +\infty)$ , то гипотеза  $H_0$  (5.17) отвергается и принимается альтернативная гипотеза  $H_1$  (5.22). Если  $K_{наб} \in [x_{лев,\alpha/2}, x_{np,\alpha/2}]$ , то принимается основная гипотеза  $H_0$  (5.17).

### ♦ Пример 5.2. Хронометраж затрат времени на сборку узла

машины  $n = 21$  слесарей показал, что  $\bar{x}_g = 77$  мин, а  $s^2 = 4$  мин<sup>2</sup>. В предположении о нормальности распределения решить вопрос: можно ли на уровне значимости  $\alpha = 0.05$  считать 80 мин нормативом (математическим ожиданием) трудоемкости?

*Решение.* В качестве основной гипотезы принимается  $H_0 : a = 80$  мин, в качестве альтернативной  $H_1 : a \neq 80$  мин, т.е. имеем случай 2, при этом  $a_0 = 80$ . Используя (5.23) и табл. П2 ( $n = 21$ ), находим

$$x_{лев,\alpha/2} = -2.086; \quad x_{np,\alpha/2} = 2.086. \quad (5.24)$$

По формуле (5.18) вычисляем  $K_{наб} = (77 - 80) / (2\sqrt{2}) = -6.708$ . Так как число  $-6.708$  попадает в критическую область (конкретно в интервал  $(-\infty, -2.086)$ ), то гипотеза  $H_0 : a = 80$  мин отвергается. ☹

### 5.3. Проверка гипотезы о числовом значении дисперсии нормального распределения

Полагаем, что  $X$  является случайной величиной, имеющей нормальное распределение  $N(a, \sigma)$ , причем числовое значение дисперсии

$\sigma^2$  неизвестно. Выборочная оценка  $S^2 = \sum_{i=1}^n (X_i - X_e)^2 / (n-1)$  дает приближенное представление о  $\sigma^2$ . Используя эту оценку, проверим гипотезу

$$H_0 : \sigma^2 = \sigma_0^2, \quad (5.25)$$

где  $\sigma_0^2$  – заранее заданное число. В качестве критерия возьмем случайную величину

$$K = \frac{(n-1)S^2}{\sigma_0^2}. \quad (5.26)$$

При выполнении гипотезы (5.25) эта величина подчиняется  $\chi^2$ -распределению с числом степеней свободы  $k = n-1$ , т.е.

$$K = \frac{(n-1)S^2}{\sigma_0^2} = \chi_{n-1}^2. \quad (5.27)$$

Зададимся уровнем значимости  $\alpha$  и перейдем к построению критических областей для проверки гипотезы  $H_0$  (5.25) при следующих двух альтернативных гипотезах  $H_1$ .

#### Случай 1

В качестве альтернативной гипотезы примем

$$H_1 : \sigma^2 > \sigma_0^2. \quad (5.28)$$

Критическая область является правосторонней и определяется интервалом  $(x_{np,\alpha}, +\infty)$ , где критическая точка  $x_{np,\alpha}$  находится из условия (5.6), которое с учетом (5.27) можно записать в виде

$$P(\chi_{n-1}^2 > x_{np,\alpha}) = \alpha.$$

В табл. ПЗ приведены квантили  $\chi^2(\gamma, k)$ , определяемые соотношением

$$P(\chi_k^2 < \chi^2(\gamma, k)) = \gamma = 1 - \alpha.$$

Следовательно, искомая критическая точка  $x_{np,\alpha}$  находится как

$$x_{np,\alpha} = \chi^2(1 - \alpha, n - 1).$$

Подставив в (5.26) конкретные значения  $S^2, \sigma_0^2$ , находим  $K_{наб}$ . Если  $K_{наб} > x_{np,\alpha}$ , то гипотеза  $H_0$  (5.25) отвергается и принимается гипотеза  $H_1$  (5.28).

#### Случай 2

В качестве альтернативной гипотезы примем

$$H_1 : \sigma^2 \neq \sigma_0^2. \quad (5.29)$$

В этом случае критическая область состоит из двух интервалов  $(0, x_{лев,\alpha/2})$  и  $(x_{np,\alpha/2}, +\infty)$ , где критические точки  $x_{лев,\alpha/2}, x_{np,\alpha/2}$  определяются из условий (5.8), которые с учетом (5.27) можно записать в виде

$$P(\chi_{n-1}^2 < x_{лев,\alpha/2}) = \alpha/2; \quad P(\chi_{n-1}^2 > x_{np,\alpha/2}) = \alpha/2.$$

Обращаясь к табл. ПЗ, находим

$$x_{лев,\alpha/2} = \chi^2(\alpha/2, n-1); \quad x_{np,\alpha/2} = \chi^2(1 - \alpha/2, n-1).$$

Если значение  $K_{наб}$ , вычисленное по формуле (5.26), попадает в один из интервалов  $(0, x_{лев,\alpha/2})$  или  $(x_{np,\alpha/2}, \infty)$ , то гипотеза  $H_0$  отвергается и принимается гипотеза  $H_1$  (5.29). В противном случае нет оснований отвергнуть гипотезу  $H_0$  (5.25).

♦ **Пример 5.3.** Точность работы станка-автомата проверяется по дисперсии контролируемого размера изделия. По выборке из 25 деталей вычислена  $s^2 = 0.25$ . При уровне значимости  $\alpha = 0.05$  проверить гипотезу  $H_0 : \sigma^2 = 0.15$ .

*Решение.* За альтернативную примем гипотезу  $H_1: \sigma^2 > 0.15$ , т.е. имеем случай 1. По табл. ПЗ находим  $x_{np,0.05} = \chi^2(0.95, 24) = 36.4$ , следовательно, критическая область  $(36.4, \infty)$ . По формуле (5.26) находим

$$K_{наб} = (25 - 1)0.25 / 0.15 = 40.$$

Так как  $K_{наб}$  попадает в критическую область, гипотезу  $H_0$  отвергаем. ●

#### 5.4. Проверка гипотезы о числовом значении вероятности события

Предположим, что  $A$  – случайное событие, вероятность  $p$  появления которого в единичном испытании неизвестна. Выдвинем гипотезу

$$H_0: p = p_0 \quad (5.30)$$

о том, что вероятность  $p$  равна числу  $p_0$ . В основе проверки этой гипотезы должно лежать сравнение числа  $p_0$  с приближенными значениями вероятности  $p$ , найденными по опытным данным. Хорошим приближением к  $p$  является относительная частота  $\omega = m/n$ , где  $n$  – число независимых испытаний, проводимых в одинаковых условиях,  $m$  – число испытаний (из  $n$  проведенных), в которых произошло событие  $A$ . Поскольку  $A$  – случайное событие, то число  $m$  – случайная величина. Поэтому рассмотрим два случая.

**Случай большого числа наблюдений.** Напомним, что при большом  $n$  распределение величины  $\frac{\omega - p}{\sqrt{p(1-p)/n}}$  можно аппроксимировать нормальным распределением  $N(0,1)$ . Если гипотеза (5.30) справедлива, то распределение критерия

$$\frac{\omega - p_0}{\sqrt{p_0(1-p_0)/n}} \quad (5.31)$$

можно аппроксимировать нормальным распределением  $N(0,1)$ , т.е.

$$\frac{\omega - p_0}{\sqrt{p_0(1-p_0)/n}} = N(0,1). \quad (5.32)$$

Напомним, что при проверке гипотез о численном значении математического ожидания (при известной дисперсии) уже использовался критерий, имеющий нормальное распределение. Поэтому, не останавливаясь на вычислении критических точек, определим только следующие три вида альтернативной гипотезы  $H_1$ .

Альтернативная гипотеза  $H_1$  имеет вид

$$H_1: p > p_0. \quad (5.33)$$

В этом случае критическая область представляет собой отрезок  $(x_{np,\alpha}, +\infty)$  (см. рис.5.1,а). Подставляя в формулу (5.31) значение частности  $\omega$  и заданные числа  $p_0$  и  $n$ , вычисляем значения критерия  $K_{наб}$ . Если  $K_{наб} > x_{np,\alpha}$ , то гипотеза  $H_0$  (5.30) отвергается и принимается гипотеза  $H_1$  (5.33).

Альтернативная гипотеза  $H_1$  имеет вид

$$H_1: p < p_0. \quad (5.34)$$

В этом случае критическая область имеет вид  $(-\infty, x_{лев,\alpha})$  (см. рис. 5.1,б). Если числовое значение  $K_{наб}$  попадает в интервал  $(-\infty, x_{лев,\alpha})$ , то принимается гипотеза  $H_1$  (5.34).

Альтернативная гипотеза  $H_1$  имеет вид

$$H_1: p \neq p_0. \quad (5.35)$$

В этом случае критическая область состоит из двух отрезков  $(-\infty, x_{лев,\alpha/2})$ ,  $(x_{np,\alpha/2}, +\infty)$  (см. рис. 5.1,в). Если числовое значение критерия  $K_{наб}$  попадает в критическую область, принимается гипотеза  $H_1$  (5.35), в противном случае – гипотеза  $H_0$  (5.30).

♦ **Пример 5.4.** Партия принимается, если вероятность того, что изделие окажется бракованным, не превышает  $p_0 = 0.02$ . Среди случайно отобранных  $n = 1000$  деталей оказалось  $m = 40$  бракованных. Можно ли при уровне значимости  $\alpha = 0.01$  принять партию?

*Решение.* Из условий задачи следует, что нулевая гипотеза

$$H_0 : p = 0.02,$$

а альтернативная имеет вид

$$H_1 : p > 0.02.$$

Критическую точку  $x_{np,\alpha}$  находим из уравнения

$$\Phi(x_{np,\alpha}) = \frac{1}{2} - 0.01 = 0.49.$$

По табл. П1 проводим линейную интерпретацию, получаем  $x_{np,\alpha} = 2.33$ . Числовое значение критерия (5.31)

$$K_{наб} = \frac{0.04 - 0.02}{\sqrt{0.02 \cdot 0.98/1000}} = 4.5.$$

Так как это число попадает в критическую область  $(2.33, +\infty)$ , то гипотезу  $H_0 : p = 0.02$  отвергаем и делаем вывод, что при уровне значимости  $\alpha = 0.01$  партию изделий принять нельзя. ☹

**Случай малого числа наблюдений.** При малом числе наблюдений допущение (5.32) несправедливо. В этом случае проверка гипотезы (5.30) проводится следующим образом.

Альтернативная гипотеза  $H_1$  имеет вид

$$H_1 : p > p_0.$$

Задаемся уровнем значимости  $\alpha$ . Полагая  $\gamma = 1 - 2\alpha$  и зная значение  $n, m$ , по табл. П4 находим  $p_1$  (это нижнее число). Если  $p_0 < p_1$ , то принимается гипотеза  $H_1 : p > p_0$ , в противном случае – гипотеза  $H_0 : p = p_0$ .

Альтернативная гипотеза  $H_1$  имеет вид

$$H_1 : p < p_0.$$

Полагая  $\gamma = 1 - 2\alpha$  и зная  $n, m$ , по табл. П4 находим  $p_2$  (верхнее число в таблице). Если  $p_0 > p_2$ , то принимаем гипотезу  $H_1 : p < p_0$ , в противном случае – гипотезу  $H_0 : p = p_0$ .

Альтернативная гипотеза  $H_1$  имеет вид

$$H_1 : p \neq p_0.$$

Полагая  $\gamma = 1 - \alpha$  и зная  $n, m$ , по табл. П4 находим  $p_1, p_2$ . Если  $p_0 < p_1$  или  $p_0 > p_2$ , то принимаем гипотезу  $H_1 : p \neq p_0$ ; если  $p_1 < p_0 < p_2$ , то принимаем гипотезу  $H_0 : p = p_0$ .

♦ **Пример 5.5.** В  $n = 5$  опытах событие  $A$  произошло  $m = 4$  раза. Можно ли принять вероятность  $p$  равной 0.2 при уровне значимости  $\alpha = 0.025$ ?

*Решение.* Основная гипотеза  $H_0$  имеет вид  $H_0 : p = p_0 = 0.2$ . Рассмотрим три случая альтернативной гипотезы.

1.  $H_1 : p > p_0$ . Принимая  $\gamma = 1 - 2\alpha = 0.95$ , по табл. П4 находим  $p_1 = 0.284$ . Так как  $p_0 < p_1$ , то принимаем гипотезу  $H_1$ , т.е. считаем, что  $p > 0.2$ .

2.  $H_1 : p < p_0$ . Для  $\gamma = 1 - 2\alpha = 0.95$  по табл. П4 находим  $p_2 = 0.995$ . Так как  $p_0 < p_2$ , то принимаем гипотезу  $H_0$ , т.е. считаем, что вероятность события  $p = 0.2$ .

3.  $H_1 : p \neq p_0$ . Полагая  $\alpha = 0.05$ , по табл. П4 для  $\gamma = 1 - \alpha = 0.95$  находим  $p_1 = 0.284$  и  $p_2 = 0.995$ . Так как  $p_0 = 0.2$  не попадает в интервал  $(0.284, 0.995)$ , то принимается гипотеза  $H_1 : p \neq 0.2$ . ☹

### 5.5. Проверка гипотезы о равенстве математических ожиданий двух нормальных распределений

Проверка гипотезы о равенстве математических ожиданий двух генеральных совокупностей имеет важное практическое значение. Действительно, иногда оказывается, что средний результат  $\bar{X}_g$  одной серии наблюдений отличается от среднего результата  $\bar{Y}_g$  другой серии. Возникает вопрос: можно ли это различие объяснить случайной ошибкой экспериментов или оно неслучайно? Иначе говоря, можно ли считать, что результаты экспериментов представляют собой выборки из двух генеральных совокупностей с одинаковыми средними. Приведем точную формулировку задачи.

Пусть генеральные совокупности  $X$  и  $Y$  распределены по нормальному закону, причем их средние квадратические отклонения известны и равны соответственно  $\sigma_X$  и  $\sigma_Y$ . Требуется по двум независимым выборкам  $x_1, \dots, x_n$  и  $y_1, \dots, y_m$  из генеральных совокупностей  $X$  и  $Y$  проверить гипотезу о равенстве генеральных средних, т.е. основная гипотеза имеет вид:

$$H_0 : M(X) = M(Y). \quad (5.36)$$

Построим критерий проверки этой гипотезы, основываясь на следующем соображении: так как приближенное представление о математическом ожидании дает выборочная средняя, то в основе проверки гипотезы (5.36) должно лежать сравнение выборочных средних  $\bar{X}_g, \bar{Y}_g$ . Найдем закон распределения разности  $(\bar{X}_g - \bar{Y}_g)$ .

Эта разность является случайной величиной, и если гипотеза  $H_0$  (5.36) верна, то

$$M(\bar{X}_g - \bar{Y}_g) = M\left(\frac{X_1 + \dots + X_n}{n} - \frac{Y_1 + \dots + Y_m}{m}\right) = M(X) - M(Y) = 0.$$

Пользуясь свойствами дисперсии, получим

$$\begin{aligned} D(\bar{X}_g - \bar{Y}_g) &= D\left(\frac{X_1 + \dots + X_n}{n} - \frac{Y_1 + \dots + Y_m}{m}\right) = \\ &= \frac{nD(X)}{n^2} + \frac{mD(Y)}{m^2} = \frac{D(X)}{n} + \frac{D(Y)}{m} = \frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}. \end{aligned} \quad (5.37)$$

Так как случайная величина  $\bar{X}_g - \bar{Y}_g$  является линейной комбинацией независимых нормально распределенных случайных величин  $X_1, \dots, X_n, Y_1, \dots, Y_m$ , то  $\bar{X}_g - \bar{Y}_g$  распределена по нормальному закону с параметрами  $a = 0$ ,  $\sigma^2 = \frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}$ . В качестве критерия выберем пронормированную случайную величину  $\bar{X}_g - \bar{Y}_g$ , т.е.

$$K = \frac{\bar{X}_g - \bar{Y}_g}{\sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}}. \quad (5.38)$$

Таким образом, если гипотеза (5.36) верна, случайная величина  $K$  имеет нормальное распределение  $N(0,1)$ , т.е.

$$K = \frac{\bar{X}_g - \bar{Y}_g}{\sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}} = N(0,1). \quad (5.39)$$

Теперь зададимся уровнем значимости  $\alpha$  и перейдем к построению критических областей и проверке гипотезы (5.36) для двух видов альтернативной гипотезы  $H_1$ . Заметим, что вычисление критических точек критерия, распределенного по нормальному закону  $N(0,1)$ , подробно рассматривалось в п. 5.2. Поэтому здесь ограничимся только определением соответствующих критических областей.

1. Альтернативная гипотеза имеет вид

$$H_1 : M(X) > M(Y). \quad (5.40)$$

В этом случае критическая область есть интервал  $(x_{np, \alpha}, +\infty)$ , где



критическая точка  $x_{np,\alpha}$  определяется из условия  $P(N(0,1) > x_{np,\alpha}) = \alpha$  (см. п. 5.2). Критическая область приведена на рис. 5.1,а. Подставляя в (5.38) числовые значения, найдем значения случайных величин  $\bar{X}_e, \bar{Y}_e$  и значение критерия  $K_{наб}$ . Если  $K_{наб} > x_{np,\alpha}$ , то гипотезу  $H_0$  (5.36) отвергаем и принимаем гипотезу  $H_1$  (5.40). Поступая таким образом, можно допустить ошибку первого рода с вероятностью  $\alpha$ .

♦ **Пример 5.6.** По двум независимым выборкам, извлеченным из нормальных генеральных совокупностей, объемы которых равны  $n = 12$  и  $m = 8$ , найдены средние значения  $\bar{x}_e = 143$ ,  $\bar{y}_e = 122$ . Генеральные дисперсии известны:  $\sigma_X^2 = D(X) = 36$ ,  $\sigma_Y^2 = D(Y) = 8$ . При уровне значимости  $\alpha = 0.005$  проверить гипотезу  $H_0: M(X) = M(Y)$  при конкурирующей гипотезе  $M(X) > M(Y)$ .

*Решение.* Критическую точку  $x_{np,\alpha}$  находим по табл. П1 из условия  $\Phi(x_{np,\alpha}) = \frac{1}{2} - \alpha = 0.495$ . Получаем  $x_{np,\alpha} = 2.58$ . Наблюдаемое значение критерия

$$K_{наб} = \frac{143 - 122}{\sqrt{\frac{36}{12} + \frac{8}{8}}} = \frac{21}{2} = 10.5.$$

Так как  $K_{наб} > 2.58$ , то гипотеза о равенстве генеральных средних отвергается на уровне значимости  $\alpha = 0.005$ . ●

2. Альтернативная гипотеза имеет вид

$$H_1: M(x) \neq M(y). \quad (5.41)$$

В этом случае наибольшая мощность критерия достигается при двусторонней критической области, состоящей из двух интервалов  $(-\infty, x_{лев,\alpha/2})$  и  $(x_{np,\alpha/2}, +\infty)$ . Критические точки определяются из условия (см. п. 5.2)

$$P(N(0,1) < x_{лев,\alpha/2}) = \alpha/2; \quad P(N(0,1) > x_{np,\alpha/2}) = \alpha/2.$$

В силу симметрии плотности распределения  $N(0,1)$  относительно нуля  $x_{лев,\alpha/2} = -x_{np,\alpha/2}$ . Если числовое значение критерия  $K_{наб}$ , вы-

численное по формуле (5.38), попадает в интервал  $(-\infty, x_{лев,\alpha/2})$  или в интервал  $(x_{np,\alpha/2}, +\infty)$ , то принимаем гипотезу  $H_1$  (5.41); если  $x_{лев,\alpha/2} < K_{наб} < x_{np,\alpha/2}$ , то принимаем гипотезу  $H_0$  (5.36).

### 5.6. Проверка гипотезы о равенстве математических ожиданий двух произвольных распределений по выборкам большого объема

Пусть  $x_1, \dots, x_n$  – выборка из генеральной совокупности  $X$ , а  $y_1, \dots, y_m$  – выборка из генеральной совокупности  $Y$ , причем объемы выборок  $n$  и  $m$  достаточно большие (не менее 30 элементов в каждой). Распределение генеральных совокупностей нам неизвестно, но недостаток этой информации компенсируется большими объемами выборок. Согласно центральной предельной теореме, случайная величина  $\bar{X}_e - \bar{Y}_e$  распределена по закону, близкому к нормальному. Если гипотеза  $H_0: M(X) = M(Y)$  верна, то

$$M(\bar{X}_e - \bar{Y}_e) = 0. \text{ Как и в п. 5.5, } D(\bar{X}_e - \bar{Y}_e) = \frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}, \text{ однако}$$

$\sigma_X^2, \sigma_Y^2$  неизвестны. Но при выборках большого объема случайные величины  $D_{ex}$  (выборочная дисперсия  $X$ ) и  $D_{ey}$  (выборочная дисперсия  $Y$ ) являются достаточно хорошими оценками для  $D(x)$  и  $D(y)$ . Поэтому случайная величина

$$K = \frac{\bar{X}_e - \bar{Y}_e}{\sqrt{\frac{D_{ex}}{n} + \frac{D_{ey}}{m}}} \quad (5.42)$$

распределена по закону, близкому к нормальному  $N(0,1)$ , и может быть принята в качестве критерия. Тогда построение критических областей для двух видов конкурирующих гипотез осуществляется так же, как и в п. 5.5.

♦ **Пример 5.7.** По двум независимым выборкам объемов  $n = 120$ ,  $m = 150$  найдены значения выборочных дисперсий  $d_{ex} = 1.2$  и  $d_{ey} = 4.5$ , а также средние значения  $\bar{x}_e = 30$ ,  $\bar{y}_e = 28.3$ .

При уровне значимости  $\alpha = 0.05$  проверить гипотезу  $H_0 : M(X) = M(Y)$  при конкурирующей  $H_1 : M(X) \neq M(Y)$ .

*Решение.* Вычислим наблюдаемое значение критерия  $K$ :

$$K_{\text{наб}} = \frac{\bar{X}_e - \bar{Y}_e}{\sqrt{\frac{d_{ex}}{n} + \frac{d_{ey}}{m}}} = \frac{30 - 28.3}{\sqrt{\frac{1.2}{120} + \frac{4.5}{150}}} = 8.5.$$

Правую границу  $x_{np, \alpha/2}$  двусторонней критической области  $(x_{np, \alpha/2}, +\infty)$  найдем из условия  $\Phi(x_{np, \alpha/2}) = (1 - \alpha)/2 = 0.475$ . Получаем  $x_{np, \alpha/2} = 1.96$ ,  $x_{лев, \alpha/2} = -1.96$ . Так как  $K_{\text{наб}} > x_{np, \alpha/2}$ , гипотеза о равенстве генеральных средних на уровне значимости  $\alpha = 0.05$  отвергается. ☺

### 5.7. Проверка гипотезы о равенстве математических ожиданий двух нормальных распределений с неизвестными, но равными дисперсиями

Сформулируем задачу. Пусть  $x_1, \dots, x_n$  и  $y_1, \dots, y_m$  – две независимые выборки из нормально распределенных генеральных совокупностей  $X$  и  $Y$  соответственно. Ранее мы рассмотрели случай выборок большого объема и научились проверять гипотезу  $H_0 : M(X) = M(Y)$ . Такую же гипотезу мы можем проверить и в том случае, если выборки имеют малый объем, но  $D(X)$  и  $D(Y)$  известны. Поэтому рассмотрим случай, когда выборки имеют малый объем и их дисперсии  $D(X)$  и  $D(Y)$  неизвестны, но равны.

Таким образом, при следующих предположениях:  
а) случайные величины  $X$  и  $Y$  имеют нормальное распределение и независимы; б)  $D(X) = D(Y) = \sigma^2$ , требуется проверить гипотезу о равенстве математических ожиданий случайных величин  $X$  и  $Y$ , т.е.

$$H_0 : M(X) = M(Y). \quad (5.43)$$

Построим критерий для проверки этой гипотезы. Для этого

рассмотрим случайные величины  $\frac{nD_{ex}}{\sigma^2}$  и  $\frac{mD_{ey}}{\sigma^2}$ . По теореме о распределении выборочных характеристик они имеют распределения  $\chi_{n-1}^2$  и  $\chi_{m-1}^2$  соответственно. Так как рассматриваются независимые выборки, то случайные величины  $\frac{nD_{ex}}{\sigma^2}$  и  $\frac{mD_{ey}}{\sigma^2}$  независимы. Поэтому их сумма имеет распределение  $\chi_{n+m-2}^2$ , т.е.

$$\frac{nD_{ex}}{\sigma^2} + \frac{mD_{ey}}{\sigma^2} = \chi_{n+m-2}^2. \quad (5.44)$$

В силу независимости величин  $X$  и  $Y$  имеем  $D(\bar{X}_e - \bar{Y}_e) = \frac{\sigma^2}{n} + \frac{\sigma^2}{m}$ . Если гипотеза  $H_0$  справедлива, то случайная величина

$$U = \frac{\bar{X}_e - \bar{Y}_e}{\sigma \sqrt{\frac{1}{n} + \frac{1}{m}}} = \frac{\sqrt{nm}}{\sigma \sqrt{n+m}} (\bar{X}_e - \bar{Y}_e) \quad (5.45)$$

имеет нормальное распределение  $N(0,1)$  (убедитесь в этом), т.е.  $U = N(0,1)$ .

Напомним, что случайная величина

$$T_{n+m-2} = \frac{U \sqrt{n+m-2}}{\sqrt{\chi_{n+m-2}^2}}$$

подчиняется распределению Стьюдента с  $n+m-2$  степенями свободы (см. п. 4.1). Подставив вместо  $U$  правую часть выражения (5.45), а вместо  $\chi_{n+m-2}^2$  левую часть (5.44), получим

$$K = \frac{\bar{X}_e - \bar{Y}_e}{\sqrt{nD_{ex} + mD_{ey}}} \times \sqrt{\frac{nm(n+m-2)}{n+m}}. \quad (5.46)$$

Эта случайная величина не содержит неизвестного параметра  $\sigma$  и может быть взята в качестве критерия для проверки гипотезы  $H_0$

(5.43). Если эта гипотеза справедлива, то критерий (5.46) имеет  $t$ -распределение с  $k = n + m - 2$  степенями свободы, т.е.

$$K = T_{n+m-2}. \quad (5.47)$$

Зададимся уровнем значимости  $\alpha$  и перейдем к построению критических областей для трех видов альтернативной гипотезы. Заметим, что ранее рассматривался критерий (5.18), имеющий распределение Стьюдента с  $k = n - 1$  степенями свободы. Сейчас рассмотрим критерий (5.46), имеющий  $t$ -распределение с  $k = n + m - 2$  степенями свободы. Никаких принципиальных различий в алгоритмы построения критических областей это не вносит. Поэтому лишь кратко приведем схемы нахождения критических точек.

1. Альтернативная гипотеза имеет вид

$$H_1 : M(X) > M(Y). \quad (5.48)$$

Критическая область представляет собой интервал  $(x_{np,\alpha}, +\infty)$ , где точка  $x_{np,\alpha}$  находится из условия

$$P(T_{n+m-2} > x_{np,\alpha}) = \alpha.$$

В табл. П2 приведены величины  $t(\gamma, N)$ , определяемые условием  $P(|T_{N-1}| < t(\gamma, N)) = \gamma$ , где  $N$  – объем выборки,  $N - 1$  – число степеней свободы. Поэтому

$$x_{np,\alpha} = t(1 - 2\alpha, n + m - 1). \quad (5.49)$$

Подставив в (5.46) числовые значения, получаем значения критерия  $K_{наб}$ . Если  $K_{наб} > x_{np,\alpha}$ , то принимается гипотеза  $H_1$  (5.48), в противном случае – гипотеза  $H_0$  (5.43).

2. Альтернативная гипотеза имеет вид

$$H_1 : M(X) < M(Y). \quad (5.50)$$

Критическая область – это интервал  $(-\infty, x_{лев,\alpha})$ , где точка  $x_{лев,\alpha}$  определяется из условия  $P(T_{n+m-2} < x_{лев,\alpha}) = \alpha$  и равна

$$x_{лев,\alpha} = -t(1 - 2\alpha, n + m - 1),$$

где  $t(1 - 2\alpha, n + m - 1)$  находится по табл. П2. Если числовое значение  $K_{наб} < x_{лев,\alpha}$ , то принимается гипотеза  $H_1$  (5.50), в противном случае – гипотеза  $H_0$  (5.43).

3. Альтернативная гипотеза имеет вид

$$H_1 : M(X) \neq M(Y). \quad (5.51)$$

В этом случае критическая область состоит из двух интервалов  $(-\infty, x_{лев,\alpha/2})$ ,  $(x_{np,\alpha/2}, +\infty)$ , где критические точки определяются из условий

$$P(T_{n+m-2} < x_{лев,\alpha/2}) = \alpha/2; \quad P(T_{n+m-2} > x_{np,\alpha/2}) = \alpha/2.$$

Используя табл. П2, получаем

$$x_{лев,\alpha/2} = -t(1 - \alpha, n + m - 1); \quad x_{np,\alpha/2} = t(1 - \alpha, n + m - 1).$$

Если числовое значение  $K_{наб}$  попадает в интервал  $(-\infty, x_{лев,\alpha/2})$  или в интервал  $(x_{np,\alpha/2}, +\infty)$ , то принимается гипотеза  $H_1$  (5.51). Если  $K_{наб}$  попадает в интервал  $(x_{лев,\alpha/2}, x_{np,\alpha/2})$ , то принимается гипотеза  $H_0$  (5.43).

♦ **Пример 5.8.** По двум малым выборкам из нормальных генеральных совокупностей  $X$  и  $Y$  найдены средние значения  $\bar{x}_g = 30$ ,  $\bar{y}_g = 39$  и значения исправленных дисперсий  $s_x^2 = 0.8$ ,  $s_y^2 = 0.4$ . Требуется на уровне значимости  $\alpha = 0.05$  проверить гипотезу  $H_0 : M(X) = M(Y)$  при конкурирующей гипотезе  $H_1 : M(X) \neq M(Y)$ . Объемы выборок равны соответственно  $n = 12$ ,  $m = 18$ .

*Решение.* Так как выборки имеют малый объем, то для применения критерия Стьюдента мы должны вначале проверить гипотезу о равенстве генеральных дисперсий  $D(X) = D(Y)$  (см. п. 5.8). Для проверки используем критерий Фишера. В качестве конкурирующей выберем гипотезу  $D(X) > D(Y)$ . Найдем наблюдаемое

значение критерия Фишера:  $K_{наб} = \frac{0.8}{0.4} = 2$ . Граница правосторонней критической области  $x_{np,\alpha} = f_\gamma(11,17) = 2.41$ . Так как  $K_{наб} < x_{np,\alpha}$ , то нет оснований отвергать гипотезу о равенстве дисперсий  $D(X)$  и  $D(Y)$ . Считая их равными, применим критерий (5.46) и вычислим

$$K = \frac{\bar{x}_g - \bar{y}_g}{\sqrt{nd_{gx} + md_{gy}}} \cdot \sqrt{\frac{mn(n+m-2)}{n+m}}.$$

Так как  $S^2 = \frac{n}{n-1} D_g$ , то  $nd_{gx} = (n-1)s_X^2$ ,  $md_{gy} = (m-1)s_Y^2$ . После вычислений получим  $K_{наб} = 3.594$ . Критическая область для критерия является двусторонней. По табл. П2 находим

$$x_{np,\alpha/2} = t(1-\alpha, 29) = 2.048; \quad x_{лев,\alpha/2} = -t(1-\alpha, 29) = -2.048.$$

Так как  $K_{наб} > 2.048$ , то гипотеза о равенстве математических ожиданий  $M(X)$  и  $M(Y)$  отвергается на уровне значимости 0.05. ●

### 5.8. Проверка гипотезы о равенстве дисперсий двух нормальных распределений

В п. 5.7 при проверке гипотезы о равенстве математических ожиданий предполагалось, что дисперсии этих совокупностей одинаковы. Как убедиться в этом, имея лишь значения выборочных дисперсий? Задача проверки гипотезы о равенстве дисперсий имеет и самостоятельный интерес. Так как дисперсия, например, характеризует точность работы прибора или технологического процесса, то, убедившись в равенстве дисперсий, можно говорить об одинаковой точности прибора или технологического процесса.

Пусть  $X$  и  $Y$  – две случайные величины, имеющие нормальные распределения и неизвестные дисперсии  $\sigma_X^2$  и  $\sigma_Y^2$ . Требуется проверить гипотезу

$$H_0 : \sigma_X^2 = \sigma_Y^2. \quad (5.52)$$

Построим критерий для проверки этой гипотезы. Для этого рассмотрим исправленные дисперсии:

$$S_X^2 = \frac{\sum_{i=1}^n (X_i - \bar{X}_g)^2}{n-1}, \quad S_Y^2 = \frac{\sum_{j=1}^m (Y_j - \bar{Y}_g)^2}{m-1}.$$

Как известно (см. п. 3.3), эти величины могут быть приняты за приближенные значения  $\sigma_X^2$  и  $\sigma_Y^2$ . Имеют место следующие распределения (см. теорему 4.1):

$$\frac{(n-1)S_X^2}{\sigma_X^2} = \chi_{n-1}^2; \quad \frac{(m-1)S_Y^2}{\sigma_Y^2} = \chi_{m-1}^2.$$

Поэтому в соответствии с определением  $F$ -распределения (см. п. 4.1) отношение  $\frac{\chi_l^2/l}{\chi_k^2/k}$  или отношение  $\frac{(n-1)S_X^2}{\sigma_X^2(n-1)} / \frac{(m-1)S_Y^2}{\sigma_Y^2(m-1)}$  будет иметь распределение Фишера с  $l = n-1$  и  $k = m-1$  степенями свободы, т.е.

$$\frac{S_X^2}{\sigma_X^2} / \frac{S_Y^2}{\sigma_Y^2} = F_{n-1, m-1}. \quad (5.53)$$

Если гипотеза (5.52) верна, то из (5.53) непосредственно получаем критерий

$$K = \frac{\max(S_X^2, S_Y^2)}{\min(S_X^2, S_Y^2)}, \quad (5.54)$$

который подчиняется распределению Фишера с  $l$  и  $k$  степенями свободы, т.е.

$$K = F_{l,k}. \quad (5.55)$$

Предположим, что выборка с большей исправленной дисперсией имеет объем  $n_1$ , с меньшей –  $m_1$ . В этом случае

$$l = n_1 - 1; k = m_1 - 1.$$

Зададим уровень значимости  $\alpha$  и перейдем к построению критических областей и проверке гипотезы (5.52) для двух следующих видов альтернативной гипотезы.

1. Альтернативная гипотеза имеет вид

$$H_1: \sigma_X^2 > \sigma_Y^2. \quad (5.56)$$

В этом случае критическая область представляет собой интервал  $(x_{np,\alpha}, +\infty)$ , где точка  $x_{np,\alpha}$  определяется из условия

$$P(F_{l,k} > x_{np,\alpha}) = \alpha.$$

Исходя из этого условия, найдем  $x_{np,\alpha}$ . В табл. П5 приведены значения  $f_\gamma(l, k)$ , удовлетворяющие условию

$$P(F_{l,k} < f_\gamma(l, k)) = \gamma = 1 - \alpha.$$

Тогда, задавая  $\gamma = 1 - \alpha$ , приходим к соотношению

$$x_{np,\alpha} = f_\gamma(l, k). \quad (5.57)$$

Перейдем к проверке гипотезы  $H_0$ . В соответствии с выражениями

$$s_X^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_e)^2, \quad s_Y^2 = \frac{1}{m-1} \sum_{j=1}^m (y_j - \bar{y}_e)^2,$$

где  $x_i, y_j$  – выборочные значения,  $\bar{x}_e, \bar{y}_e$  – значения выборочных средних, находим  $s_X^2, s_Y^2$ . Подставляя эти значения в (5.54), вычисляем числовое значение критерия  $K_{наб}$ . Если  $K_{наб} > \bar{x}_{np,\alpha}$ , то гипотеза  $H_0$  (5.52) отвергается и принимается гипотеза  $H_1$ . При этом можно совершить ошибку первого рода с вероятностью  $\alpha$ . Если  $K_{наб} < x_{np,\alpha}$ , то принимается гипотеза  $H_0$ .

♦ **Пример 5.9.** По двум независимым выборкам объемов  $n=9, m=13$ , извлеченным из нормальных генеральных совокупностей, найдены исправленные дисперсии  $s_X^2=12, s_Y^2=6$ . При уровне значимости  $\alpha=0.05$  проверить нулевую гипотезу  $H_0: \sigma_X^2 = \sigma_Y^2$  при альтернативной  $H_1: \sigma_X^2 > \sigma_Y^2$ .

*Решение.* Вычислим значение критерия по формуле (5.54):  $K_{наб} = 12/6 = 2$ . В соответствии с соотношением (5.57) находим точку

$$x_{np,\alpha} = f_{0.95}(8, 12) = 2.85 \quad (l = n_1 - 1 = 9 - 1 = 8; k = m_1 - 1 = 13 - 1 = 12).$$

Так как  $K_{наб} < 2.85$ , то принимается гипотеза  $H_0: \sigma_X^2 = \sigma_Y^2$ . ●

2. Альтернативная гипотеза  $H_1$  имеет вид

$$H_1: \sigma_X^2 \neq \sigma_Y^2. \quad (5.58)$$

В этом случае критическая область состоит из двух интервалов  $(0, x_{лев,\alpha/2})$ ,  $(x_{np,\alpha/2}, +\infty)$ , где точки  $x_{лев,\alpha/2}$  и  $x_{np,\alpha/2}$  определяются следующими соотношениями (докажите это):

$$x_{лев,\alpha/2} = \frac{1}{f_{1-\alpha/2}(l, k)}; \quad x_{np,\alpha/2} = f_{1-\alpha/2}(l, k), \quad (5.59)$$

в которых, как и прежде, значения  $f_\gamma(l, k)$  находятся по табл. П5.

При попадании числового значения  $K_{наб}$  (5.54) в интервал  $(0, x_{лев,\alpha/2})$  или  $(x_{np,\alpha/2}, +\infty)$  принимается гипотеза  $H_1$  (5.58); если  $K_{наб}$  попадает в интервал  $[x_{лев,\alpha/2}, x_{np,\alpha/2}]$ , то принимается гипотеза  $H_0$  (5.52).

♦ **Пример 5.10.** По двум независимым выборкам, объемы которых  $n=13, m=15$ , извлеченным из нормальных генеральных совокупностей, найдены исправленные выборочные дисперсии  $s_X^2=1.05, s_Y^2=0.35$ . При уровне значимости  $\alpha=0.10$  проверить гипотезу  $H_0: \sigma_X^2 = \sigma_Y^2$  при конкурирующей гипотезе  $H_1: \sigma_X^2 \neq \sigma_Y^2$ .

*Решение.* Вычислим  $K_{наб} = s_X^2/s_Y^2 = 1.05/0.35 = 3$ . Количество степеней свободы  $l=13-1=12; k=15-1=14$ . По табл. П5 для  $\gamma=1-\alpha/2=0.95, l=12, k=14$  находим  $f_{0.95}(12, 14) = 2.53$ . Тогда, используя (5.59), получаем

$$x_{лев,\alpha/2} = 1/2.53 = 0.395; \quad x_{np,\alpha/2} = 2.53.$$

Так как  $K_{наб} = 3 > 2.53$ , то гипотеза  $H_0: \sigma_X^2 = \sigma_Y^2$  отвергается и принимается гипотеза  $H_1: \sigma_X^2 \neq \sigma_Y^2$ . ●

В заключение сделаем следующее замечание. Выше, в п. 5.2, 5.3, 5.5, 5.7, предполагалась нормальность распределения исследуемых случайных величин  $X$  и  $Y$ . Однако приведенные критерии

весьма устойчивы (особенно при больших объемах выборок) к отклонению от нормального распределения. Данный факт позволяет надеяться на успешное использование этих критериев для проверки гипотез в случаях, когда нет уверенности в нормальном распределении случайных величин  $X$  и  $Y$ .

### 5.9. Проверка гипотезы о законе распределения с применением критерия согласия Пирсона

В предыдущих пунктах этой главы рассматривались гипотезы, относящиеся к отдельным параметрам распределения случайных величин, при этом предполагался известный вид самого распределения.

При обработке статистических данных большого объема часто возникает ситуация, когда закон распределения генеральной совокупности не известен заранее. Однако сравнение гистограммы с известными кривыми функций плотностей позволяет выдвинуть гипотезу о виде распределения генеральной совокупности. Так, например, если гистограмма имеет один явно выраженный пик (рис. 5.2,а), то можно предположить, что исследуемая генеральная совокупность распределена по нормальному закону  $N(a, \sigma)$ , т.е. имеет плотность

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-a)^2}{2\sigma^2}}.$$

Если гистограмма представляет собой "убывающие ступеньки прямоугольников" (см. рис. 5.2,б), то генеральная совокупность может быть распределена по показательному закону:

$$p(x) = \begin{cases} 0, & x < x_0; \\ \lambda e^{-\lambda(x-x_0)}, & x \geq x_0. \end{cases}$$

Для гистограммы, представленной на рис. 5.2,в, естественно выдвинуть гипотезу о равномерном распределении генеральной совокупности.

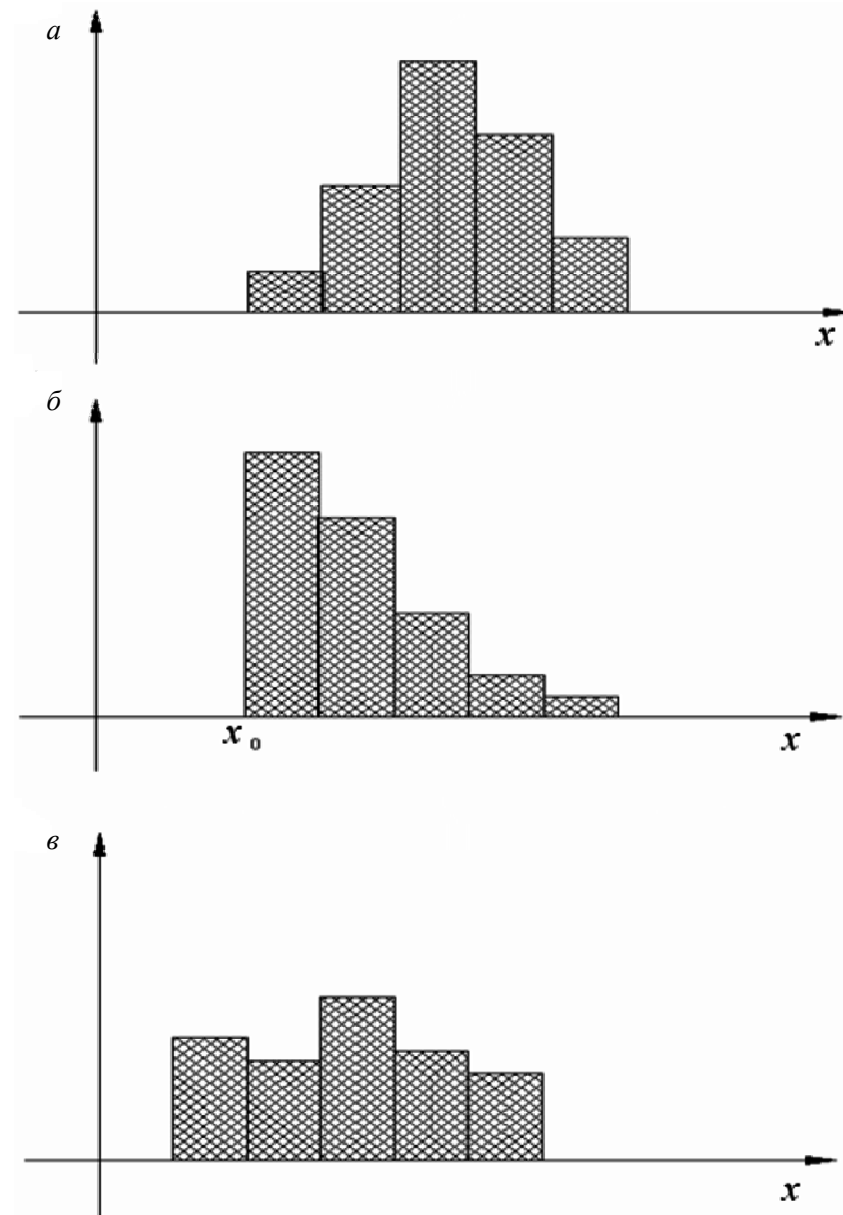


Рис. 5.2. К проверке гипотезы о законе распределения

Возникает вопрос о критерии проверки по выборочным данным гипотезы о том, что случайная величина  $X$  подчиняется распределению с плотностью  $y = p(x)$ . Такие критерии называются *критериями согласия*. Рассмотрим лишь один критерий согласия, использующий  $\chi^2$ -распределение и получивший название критерия согласия Пирсона (или критерия  $\chi^2$ ). Выдвигая гипотезу о виде распределения генеральной совокупности, мы должны различать два случая. В первом из них вид функции плотности определен в гипотезе полностью. Например, мы выдвигаем гипотезу о том, что генеральная совокупность распределена по нормальному закону с параметрами  $a = 0$  и  $\sigma = 1$ . Такие гипотезы называются *простыми*. Если же гипотеза состоит лишь в том, что функция плотности  $p(x)$  принадлежит к некоторому семейству функций, то такая гипотеза называется *сложной*. Например, можно выдвинуть гипотезу о том, что генеральные совокупности распределены по показательному закону, не оговаривая значений параметров  $\lambda$  и  $x_0$ . Такая гипотеза будет сложной.

Остановимся вначале на простой гипотезе, предполагая, что генеральная совокупность распределена непрерывно. В качестве нулевой гипотезы принимается предположение, что неизвестная плотность распределения  $p_X(x)$  исследуемой случайной величины  $X$  совпадает с предполагаемой плотностью  $p(x)$ , т.е.

$$H_0 : p_X(x) = p(x). \quad (5.60)$$

В качестве предполагаемой (теоретической) плотности могут быть рассмотрены различные плотности (нормальная, показательная и т.д.). Выберем наименьшее и наибольшее значения в данной выборке:  $a = \min\{x_1, \dots, x_n\}$ ,  $b = \max\{x_1, \dots, x_n\}$ . Промежуток  $[a, b]$  разобьем на  $l$  промежутков равной длины  $h = \frac{b-a}{l}$ . Границы этих промежутков обозначим  $z_0 = a, z_1, \dots, z_l = b$ , где  $z_{i+1} = z_i + h$  при  $i = 0, \dots, l-1$ . Считаем, что гипотеза верна. Вычислим частоту  $m_i (i = 1, \dots, l)$  попадания элементов генеральной совокупности на каждый промежуток. Понятно, что

$m_1 + m_2 + \dots + m_l = n$ . Сдвинем границу левого интервала на  $-\infty$ , а правого на  $+\infty$ , т.е. вместо первого интервала  $(z_0, z_1)$  рассмотрим интервал  $(-\infty, z_1)$ , а вместо последнего  $(z_{l-1}, z_l)$  — интервал  $(z_{l-1}, \infty)$ . Вычислим вероятность попадания случайной величины  $X$  на каждый из полученных промежутков  $\Delta_1, \dots, \Delta_l$ , воспользовавшись известной формулой:

$$p_i = \int_{\Delta_i} p(x) dx, \quad i = 1, 2, \dots, l.$$

Заметим, что первый и последний из интегралов являются несобственными. Полученные вероятности  $p_1, \dots, p_l$  должны удовлетворять условию  $p_1 + p_2 + \dots + p_l = 1$ .

Рассмотрим  $n$  опытов, каждый из которых состоит в выборе случайного значения величины  $X$  и события  $A_i = \{\text{значение попало в интервал } \Delta_i\}$ . Событие  $A_i$  в каждом опыте происходит с вероятностью  $p_i$ . Поэтому ожидаемое количество появлений события  $A$  в  $n$  опытах равно  $np_i$  (математическое ожидание биномиального распределения). Понятно, что если гипотеза верна, то между фактическими частотами  $m_i$  и теоретическими  $np_i$  попаданий на  $i$ -й интервал не должно быть "больших" расхождений, т.е. величины  $np_1, \dots, np_l$  и числа  $m_1, \dots, m_l$  должны быть соответственно близки друг к другу. В качестве меры расхождения между ними используем сумму квадратов взвешенных расхождений:

$$Y_i = \frac{m_i - np_i}{\sqrt{np_i}}.$$

Случайная величина  $\sum_{i=1}^l Y_i^2 = \sum_{i=1}^l \frac{(m_i - np_i)^2}{np_i}$  при большом объеме выборки  $n$  имеет распределение, близкое к  $\chi^2$  с  $(l-1)$  степенями свободы. Поэтому эта случайная величина принимается за критерий

$$K = \sum_{i=1}^l \frac{(m_i - np_i)^2}{np_i}. \quad (5.61)$$

Если гипотеза  $H_0$  (5.60) справедлива, то критерий  $K$  имеет  $\chi^2$ -распределение с  $k = l - 1$  степенями свободы, т.е.

$$K = \sum_{i=1}^l \frac{(m_i - np_i)^2}{np_i} = \chi_k^2. \quad (5.62)$$

Далее задаемся уровнем значимости  $\alpha$  и, зная распределение критерия  $K$ , строим правостороннюю критическую область. Это будет область вида  $(x_{np,\alpha}, +\infty)$ . Критическая точка  $x_{np,\alpha}$  находится из условия  $P(\chi_k^2 > x_{np,\alpha}) = \alpha$ . В табл. ПЗ приведены значения  $\chi_\gamma^2$ , удовлетворяющие условию  $P(\chi_k^2 < \chi_\gamma^2) = \gamma$ . Следовательно,

$$x_{np,\alpha} = \chi^2(1 - \alpha, l - 1). \quad (5.63)$$

Если числовое значение критерия  $K_{наб}$ , вычисляемое по формуле (5.61), попадает в критическую область  $(x_{np,\alpha}, \infty)$ , то делается вывод о неправомерности гипотезы  $H_0$  (5.60). Следует помнить, что этот вывод может быть ошибочным (т.е. генеральная совокупность имеет плотность распределения  $p(x)$  с вероятностью  $\alpha$  (ошибка первого рода)).

Отметим одну рекомендацию для выбора длины интервала  $h$ .

Чтобы случайная величина  $\sum_{i=1}^l \frac{(m_i - np_i)^2}{np_i}$  была приемлемо близка

к распределению  $\chi_{l-1}^2$ , достаточным для практических расчетов является выполнение условия  $np_i \geq 10$  для всех  $i$ . В том случае, когда для некоторого  $i$  имеет место  $np_i < 10$ , рекомендуется объединить несколько интервалов, пока данное условие не будет выполнено.

♦ **Пример 5.11.** По выборке объема  $n = 144$  составлен группированный статистический ряд:

$X$	0–1	1–2	2–3	3–4	4–5	5–6	6–7	7–8
$m_i$	16	17	19	16	24	19	17	16

Проверить на уровне значимости  $\alpha = 0.05$  гипотезу о равномерности распределения генеральной совокупности на отрезке  $[0, 8]$ .

*Решение.* Нулевая гипотеза имеет вид

$$H_0 : p_X(x) = p(x) = \begin{cases} \frac{1}{8-0}, & 0 \leq x \leq 8; \\ 0, & \text{для остальных } x. \end{cases} \quad (5.64)$$

Вычислим вероятность попадания случайной величины  $X$  в каждый интервал:

$$p_i = \int_{i-1}^i \frac{1}{8} dx = \frac{1}{8}(i - i + 1) = \frac{1}{8}, \quad i = 1, 2, \dots, 8.$$

Поэтому  $np_i = \frac{1}{8}144 = 18$  при любом  $i$ . Так как  $np_i \geq 10$ , то нет необходимости объединять несколько интервалов. Результаты дальнейших вычислений сведены в табл. 5.1.

Таблица 5.1

Номер интервала	$m_i$	$np_i$	$m_i - np_i$	$\frac{(m_i - np_i)^2}{np_i}$
1	16	18	-2	0.22
2	17	18	-1	0.06
3	19	18	1	0.06
4	16	18	-2	0.22
5	24	18	6	2.00
6	19	18	1	0.06
7	17	18	-1	0.06
8	16	18	-2	0.22
$\Sigma$	144	144	0	2.9



Таким образом, числовое значение  $K_{наб} = 2.9$ . Для заданного уровня значимости  $\alpha = 0.05$  находим  $\gamma = 1 - \alpha = 0.95$ ,  $\chi^2 = (0.95, 7) = 14.1$ . Так как  $K_{наб} < \chi_{пр, \alpha}$ , то гипотеза  $H_0$  (5.60) принимается. ☺

Обычной является ситуация, когда предполагается лишь, что распределение генеральной совокупности принадлежит некоторому классу распределений. Например, генеральная совокупность распределена нормально. В этой гипотезе не оговорены значения параметров  $a$  и  $\sigma$ . Отличие в применении критерия  $\chi^2$  в этом случае от ранее рассмотренного состоит в том, что нет возможности сразу вычислить значения вероятностей. Поэтому вначале находят оценки неизвестных параметров. Например, для оценки параметра  $a$ , как известно, можно использовать случайную величину  $\bar{X}_g$  и заменить  $a$  ее значением, т.е.  $a = \bar{x}_g$ .

В качестве оценки параметра  $\sigma^2$  можно выбрать исправленную дисперсию  $S^2$  и заменить  $\sigma^2$  ее значением  $s^2$ . Таким образом,

$$p(x) = \frac{1}{\sqrt{2\pi}s} e^{-\frac{(x-\bar{x}_g)^2}{2s^2}}.$$

В качестве критерия также принимается случайная величина (5.61). Если гипотеза  $H_0$  справедлива, то критерий имеет  $\chi^2$ -распределение с  $k$  степенями свободы. Однако количество степеней свободы критерия подсчитывается по формуле  $l - r - 1$ , где  $r$  – количество параметров, оцененных по выборке. В рассмотренном примере  $r = 2$ , так как по выборке были оценены два параметра  $a$  и  $\sigma$ . В этом же примере вероятность  $p_i$  попадания случайной величины  $X$  в интервал  $[z_{i-1}, z_i]$  находится с помощью функции Лапласа

$$p_i = P(z_{i-1} < N(\bar{x}_g, s) < z_i) = \Phi\left(\frac{z_i - \bar{x}_g}{s}\right) - \Phi\left(\frac{z_{i-1} - \bar{x}_g}{s}\right).$$

♦ **Пример 5.12.** Группированный статистический ряд частот занесен в графы 2 и 3 табл. 5.2. По выборке объема  $n = 200$  най-

дено  $\bar{x}_g$ ,  $s^2 = 94.26$ . При уровне значимости  $\alpha = 0.02$  проверить гипотезу о нормальности распределения генеральной совокупности.

Таблица 5.2

Но- мер ин- тер- вала	Границы интер- валов	$m_i$	$\frac{z_{i-1} - \bar{x}_g}{s}$	$\Phi\left(\frac{z_{i-1} - \bar{x}_g}{s}\right)$	$p_i$	$np_i$	$\frac{(m_i - np_i)^2}{np_i}$
1	2	3	4	5	6	7	8
1	[-20,15]	7	-1.99	-0.4767	0.023	4.66	1.18
2	[-15,10]	11	-1.47	-0.4292	0.047	9.50	0.24
3	[-10,-5]	15	-0.96	-0.331	0.098	19.54	1.05
4	[-5,0]	24	-0.44	-0.1700	0.162	32.30	2.13
5	[0,5]	49	0.07	0.0279	0.198	39.58	2.24
6	[5,10]	41	0.59	0.222	0.194	38.90	0.11
7	[10,15]	26	1.10	0.364	0.142	28.38	0.20
8	[15,20]	17	1.62	0.4474	0.083	16.62	0.01
9	[20,25]	7	2.13	0.4834	0.053	10.52	0.03
10	[25,30]	3	$+\infty$	0.5			
$\Sigma$		200			1	200.0	7.19

*Решение.* Так как  $p_i = \Phi\left(\frac{z_i - \bar{x}_g}{s}\right) - \Phi\left(\frac{z_{i-1} - \bar{x}_g}{s}\right)$ , то в графе 4

вычислены значения  $\frac{z_{i-1} - \bar{x}_g}{s}$ . При этом левая граница первого ин-

тервала заменена на  $-\infty$ , а правая граница последнего интервала заменена на  $+\infty$ . В графе 5 вычислены значения  $\frac{z_{i-1} - \bar{x}}{s}$ , в графе 6 – вероятности  $p_i$ , в графе 7 – математические ожидания  $np_i$ , а в графе 8 – взвешенные отклонения  $\frac{(m_i - np_i)^2}{np_i}$ . Так как для 9-го и 10-го интервалов  $np_9 = 7.2 < 10$  и  $np_{10} = 3.32 < 10$ , то эти интервалы объединяем. Для полученного интервала  $np = 10.52 > 10$  (см. графу 7). Числовое значение критерия  $K_{наб} = 7.19$  (см. итог графы 8). По табл. ПЗ при  $\gamma = 1 - \alpha = 0.98$  и  $k = 9 - 2 - 1 = 6$  находим  $\chi^2(0.98) = 15.0$ ,  $x_{np,\alpha} = 15.0$ . Так как  $K_{наб} < 15.0$ , то гипотеза  $H_0$  о нормальности распределения генеральной совокупности принимается на уровне значимости  $\alpha = 0.02$ . ●

#### 5.10. Проверка гипотезы о независимости двух генеральных совокупностей с применением критерия $\chi^2$

Пусть  $(X, Y)$  – двумерная генеральная совокупность, причем все значения случайной величины  $X$  исчерпываются числами  $a_1, \dots, a_l$ , а все значения случайной величины  $Y$  – числами  $b_1, \dots, b_s$ . Выборка объема  $n$  в этом случае состоит из пар  $(x_1, y_1), \dots, (x_n, y_n)$ , где  $x_i$  и  $y_i$  – соответствующие значения случайных величин  $X$  и  $Y$ . Заполним таблицу, называемую корреляционной, в первой строке которой перечислим все различные значения случайной величины  $Y$ , в первом столбце – все различные значения случайной величины  $X$ , а на пересечении  $i$ -й строки и  $j$ -го столбца поместим число  $n_{ij}$  – количество пар  $(a_i, b_j)$ , встречающихся в выборке. Сумму элементов  $\sum_{j=1}^s n_{ij}$   $i$ -й строки обозна-

чим  $n_{i\bullet}$ . Аналогично  $\sum_{i=1}^l n_{ij} = n_{\bullet j}$ . Ясно, что

$$\sum_{i=1}^l \sum_{j=1}^s n_{ij} = \sum_{j=1}^s \sum_{i=1}^l n_{ij} = \sum_{j=1}^s n_{\bullet j} = \sum_{i=1}^l n_{i\bullet} = n.$$

Если числа  $n_{ij}$  концентрируются вдоль диагонали, идущей из левого верхнего угла к правому нижнему, то между величинами  $X$  и  $Y$  можно предположить тесную прямую связь.

Если числа  $n_{ij}$  сосредоточены вдоль другой диагонали, то между случайными величинами  $X$  и  $Y$  вероятна обратная связь, т.е. с ростом  $X$  значения  $Y$  убывают. Если числа  $n_{ij}$  распределены по большинству ячеек таблицы, то между  $X$  и  $Y$  скорее всего нет связи.

Предположим, что анализ корреляционной таблицы позволил нам выдвинуть гипотезы: основную  $H_0$  – случайные величины  $X$  и  $Y$  независимы и альтернативную  $H_1$  – случайные величины  $X$  и  $Y$  зависимы. Используем критерий  $\chi^2$  для проверки этих гипотез. Если гипотеза  $H_0$  верна, то

$$P(X = a_i, Y = b_j) = P(X = a_i) \cdot P(Y = b_j).$$

Корреляционная таблица

$X \backslash Y$	$b_1$	$b_2$	...	$b_s$
$a_1$	$n_{11}$	$n_{12}$	...	$n_{1s}$
$a_2$	$n_{21}$	$n_{22}$	...	$n_{2s}$
...	...	...	...	...
$a_l$	$n_{l1}$	$n_{l2}$	...	$n_{ls}$

Пусть значение  $X = a_i$  встречается среди чисел  $x_1, \dots, x_n$   $n_{i\bullet}$  раз. Тогда относительная частота события  $\{X = a_i\}$  равна  $n_{i\bullet}/n$ . Она является состоятельной и несмещенной оценкой параметра

$p_i = P(X = a_i)$ . Аналогично  $n_{\bullet j}/n$  – состоятельная и несмещенная оценка вероятности  $p'_j = P(Y = b_j)$ . Если гипотеза  $H_0$  верна, то ожидаемое количество попаданий в клетку  $(i, j)$  можно найти по формуле  $n'_{ij} = n \cdot \frac{n_{i\bullet} n_{\bullet j}}{n^2} = \frac{n_{i\bullet} n_{\bullet j}}{n}$  (как математическое ожидание случайной величины, распределенной по биномиальному закону с параметрами  $n$  и  $p = \frac{n_{i\bullet} n_{\bullet j}}{n^2}$ ) и числа  $n_{ij}$ ,  $\frac{n_{i\bullet} n_{\bullet j}}{n}$  близки друг к другу в совокупности. В качестве критерия примем случайную величину

$$K = \sum_{i=1}^l \sum_{j=1}^s \frac{\left( n_{ij} - \frac{n_{i\bullet} n_{\bullet j}}{n} \right)^2}{n_{i\bullet} n_{\bullet j} / n}. \quad (5.65)$$

Если гипотеза  $H_0$  справедлива, то эта случайная величина имеет  $\chi^2$ -распределение с  $k = (l-1)(s-1)$  степенями свободы, т.е.

$$K = \chi^2_{(l-1)(s-1)}. \quad (5.66)$$

Критическая область представляет собой отрезок  $(x_{np}, +\infty)$ , где точка  $x_{np, \alpha}$  определяется соотношением  $x_{np, \alpha} = \chi^2(1 - \alpha, (l-1)(s-1))$ .

Если числовое значение критерия  $K_{наб}$ , найденное по формуле (5.65), попадает в критическую область, т.е.  $K_{наб} > x_{np, \alpha}$ , то нулевая гипотеза о независимости  $X$  и  $Y$  отвергается.

Заметим, что вместо ограничения  $np_i \geq 10$ , указанного в п. 5.9, здесь желательно выполнение условия  $\frac{n_{i\bullet} n_{\bullet j}}{n} \geq 4$ . Если это условие не выполняется, то соответствующие строки и столбцы должны быть объединены с соседними.

♦ **Пример 5.13.** Комплекующие изделия одного наименования поступают с трех предприятий: 1, 2, 3. Результаты проверки изделий приведены в табл. 5.3.

Таблица 5.3

Результаты проверки изделий	Поставщик			Всего
	1	2	3	
Годные	29	38	53	120
Негодные	1	2	7	10
Всего	30	40	60	130

Можно ли считать, что качество изделий не зависит от поставщика? Уровень значимости принять равным 0.05.

*Решение.* Находим наблюдаемое значение критерия:

$$K_{наб} = \frac{\left( 29 - \frac{120 \cdot 30}{130} \right)^2}{\frac{120 \cdot 30}{130}} + \frac{\left( 38 - \frac{120 \cdot 40}{130} \right)^2}{\frac{120 \cdot 40}{130}} + \frac{\left( 53 - \frac{120 \cdot 60}{130} \right)^2}{\frac{120 \cdot 60}{130}} + \frac{\left( 1 - \frac{10 \cdot 30}{130} \right)^2}{\frac{10 \cdot 30}{130}} + \frac{\left( 2 - \frac{10 \cdot 40}{130} \right)^2}{\frac{10 \cdot 40}{130}} + \frac{\left( 7 - \frac{10 \cdot 60}{130} \right)^2}{\frac{10 \cdot 60}{130}} = 2.55. \quad (5.67)$$

По табл. ПЗ для числа степеней свободы  $k = (l-1)(s-1) = (2-1)(3-1) = 2$  и  $\alpha = 0.05$  находим  $\chi^2(0.95, 2) = 6$ ,  $x_{np, \alpha} = 6$ . Так как  $K_{наб} < 6$ , то можно принять гипотезу  $H_0$  о независимости качества изделий от поставщика. ☹

### 5.11. Проверка статистических гипотез в Excel

В табличном процессоре Excel определены несколько функций и режимов работы *Пакета анализа*, которые можно использовать для проверки различных статистических гипотез.

**Проверка гипотезы о числовом значении математического ожидания нормального распределения при известной дисперсии.** В качестве нулевой гипотезы  $H_0$  принимается (5.13), в качестве альтернативной  $H_1$  – (5.14). Уровень значимости  $\alpha$  принимается равным 0.05.

Используется функция ZTEST, обращение к которой имеет вид:

$$=ZTEST(\text{массив}; a_0; \sigma),$$

где *массив* – адреса ячеек, содержащих выборочные данные случайной величины, математическое ожидание которой сравнивается с заданной величиной  $a_0$ ;

$a_0$  – задаваемое значение математического ожидания;

$\sigma$  – задаваемое среднее квадратичное отклонение случайной величины (если этот параметр опущен, то используется выборочная дисперсия, вычисленная по той же выборке).

Результатом работы функции является корень  $x_{np,0.05/2}$  уравнения (5.8), т.е.

$$x_{np,0.05/2} = ZTEST(\text{массив}; a_0; \sigma).$$

Величины  $x_{np,0.05/2}$ ,  $x_{лев,0.05/2} = -x_{np,0.05/2}$  определяют критические области  $(-\infty, x_{лев,0.05/2}]$ ,  $[x_{np,0.05/2}, \infty)$ .

**Проверка гипотезы о равенстве математических ожиданий двух нормальных распределений с известными дисперсиями.** Изучаются две нормально распределенные случайные величины  $X \sim N(a_X, \sigma_X)$ ,  $Y \sim N(a_Y, \sigma_Y)$ . Числовые значения дисперсий  $\sigma_X^2$ ,  $\sigma_Y^2$  известны. Проверяется основная гипотеза  $H_0$  (5.41) –  $H_0: M(X) = M(Y)$ .

Для проверки этой гипотезы используется режим работы **Двухвыборочный z-тест для средних**. Для вызова этого режима необходимо обратиться к пункту **Сервис** строки меню Excel, команде **Пакет анализа**. Затем в появившемся списке режимов вы-

брать данный режим и щелкнуть ОК. В диалоговом окне (рис. 5.3) задаются следующие параметры:

Рис. 5.3. Задание параметров режима **Двухвыборочный z-тест для средних**

*Интервал переменной 1:* – адреса ячеек, содержащих выборочные значения случайной величины  $X$ .

*Интервал переменной 2:* – адреса ячеек, содержащих выборочные значения случайной величины  $Y$ .

*Гипотетическая средняя разность:* – задает число, равное предполагаемой разности математических ожиданий  $a_X - a_Y$  (при проверке гипотезы о равенстве математических ожиданий задается 0).

*Дисперсия переменной 1 (известная):* – вводится известное значение  $\sigma_X^2$ .

*Дисперсия переменной 2 (известная):* – вводится известное значение  $\sigma_Y^2$ .

*Метки* – включается, если первая строка содержит заголовки столбцов.

*Альфа*: – задается уровень значимости.

*Выходной интервал*: / *Новый рабочий лист*: / *Новая рабочая книга* – указывается, куда выводятся результаты вычислений. При включении *Выходной интервал*: вводится адрес ячейки, начиная с которой выводятся результаты, которые оформлены в виде таблицы (пример такой таблицы приведен на рис. 5.4).

♦ **Пример 5.14.** Выборочные данные о диаметре валиков (мм), изготовленных автоматом 1 и автоматом 2, приведены в столбцах А, В документа Excel (рис. 5.5). Предварительным анализом установлено, что размер валиков, изготовленных каждым автоматом, имеет нормальное распределение с дисперсиями  $\sigma_X^2 = 5 \text{ мм}^2$  (автомат 1) и  $\sigma_Y^2 = 7 \text{ мм}^2$  (автомат 2).

Необходимо проверить нулевую гипотезу  $H_0: a_X = a_Y$  при альтернативной гипотезе  $H_1: a_X \neq a_Y$ .

Двухвыборочный z-тест для средних		
	<i>Автомат 1</i>	<i>Автомат 2</i>
Среднее	181,979	185,03
Известная дисперсия	5,000	7,00
Наблюдения	14,000	9,00
Гипотетическая разность средних	0,000	
z	-2,867	
P(Z<=z) одностороннее	0,002	
z критическое одностороннее	1,645	
P(Z<=z) двухстороннее	0,004	
z критическое двухстороннее	1,960	

Рис. 5.4. Результаты работы режима *Двухвыборочный z-тест для средних*

	А	В
1	Автомат 1	Автомат 2
2	182,3	185,3
3	183,0	185,6
4	181,8	184,8
5	181,4	186,2
6	181,8	185,8
7	181,6	184,0
8	183,2	185,2
9	182,4	184,2
10	182,5	184,2
11	179,7	
12	179,9	
13	181,9	
14	182,8	
15	183,4	
16	Среднее	Среднее
17	<b>182,0</b>	<b>185,0</b>

Рис. 5.5. Исходные данные к примеру 5.14

**Проверка гипотезы о равенстве математических ожиданий двух нормальных распределений с неизвестными, но равными дисперсиями.** Изучаются две нормально распределенные случайные величины  $X \sim N(a_X, \sigma_X)$  и  $Y \sim N(a_Y, \sigma_Y)$ . Дисперсии равны, но не известны, т.е.  $\sigma_X^2 = \sigma_Y^2 = \sigma^2$ . Необходимо проверить статистическую гипотезу  $H_0: a_X = a_Y$  при альтернативной гипотезе  $H_1: a_X \neq a_Y$ .

Для проверки этой гипотезы используется режим *Двухвыборочный t-тест с одинаковыми дисперсиями*. Для вызова режима необходимо обратиться к пункту *Сервис* строки меню Excel, команде *Пакет анализа*. Затем в появившемся списке режимов выбрать данный режим и щелкнуть ОК. В появившемся диалоговом окне этого режима задаются следующие параметры (рис. 5.6):

*Решение.* Обратимся к режиму *Двухвыборочный z-тест для средних* и в появившемся диалоговом окне зададим необходимые параметры (см. рис. 5.3), а затем щелкнем на ОК. Результаты работы режима показаны на рис. 5.4. Величина  $z$  является расчетным значением критерия (5.39)  $K_{наб} = z = -2.867$ . Это значение попадает в критическую область  $|K_{наб}| > |z_{кр}| = 1.96$ . Поэтому нулевая гипотеза с уровнем значимости  $\alpha = 0.05$  отвергается и принимается альтернативная гипотеза  $a_X \neq a_Y$ . ☹

**Двухвыборочный t-тест с одинаковыми дисперсиями**

**Входные данные**

Интервал переменной 1:

Интервал переменной 2:

Гипотетическая средняя разность:

☒ Метки

Альфа:

**Параметры вывода**

☒ Выходной интервал:

☐ Новый рабочий лист:

☐ Новая рабочая книга

OK, Отмена, Справка

Рис. 5.6. Задание параметров режима  
*Двухвыборочный t-тест с одинаковыми дисперсиями*

*Интервал переменной 1:* – адреса ячеек, содержащих выборочные значения случайной величины  $X$ .

*Интервал переменной 2:* – адреса ячеек, содержащих выборочные значения случайной величины  $Y$ .

*Гипотетическая средняя разность:* – задает число, равное предполагаемой разности математических ожиданий  $a_X - a_Y$  (при проверке гипотезы  $a_X = a_Y$  задается 0).

*Метки* – включается, если первая строка содержит заголовки столбцов.

*Альфа:* – задает уровень значимости  $\alpha$ .

*Выходной интервал:* / *Новый рабочий лист:* / *Новая рабочая книга* – указывается, куда выводятся результаты вычислений. При включении *Выходной интервал:* вводится адрес ячейки, начиная с которой выводятся результаты, представленные в виде таблицы (пример такой таблицы приведен на рис. 5.7).

Двухвыборочный t-тест с одинаковыми дисперсиями		
	Старая технология	Новая технология
Среднее	307,11	304,92
Дисперсия	1,61	2,24
Наблюдения	9,00	13,00
Объединенная дисперсия	1,99	
Гипотетическая разность средних	0,00	
df	20,00	
t-статистика	3,58	
P(T<=t) одностороннее	0,00094	
t критическое одностороннее	1,72	
P(T<=t) двухстороннее	0,00189	
t критическое двухстороннее	2,09	

Рис. 5.7. Результаты работы режима  
*Двухвыборочный t-тест с одинаковыми дисперсиями*

♦ **Пример 5.15.** Выборочные данные о расходе сырья при производстве продукции по старой и новой технологии приведены в столбцах А, В документа Excel (рис. 5.8). Предполагая, что расход сырья по старой и новой технологии распределен по нормальному закону и имеет одинаковую дисперсию, проверить статистическую гипотезу  $a_X = a_Y$  при уровне значимости  $\alpha = 0.05$ .

	А	В
	Старая технология	Новая технология
1		
2	308	308
3	308	304
4	307	306
5	308	306
6	304	306
7	307	304
8	307	304
9	308	306
10	307	306
11		304
12		303
13		304
14		303

Рис. 5.8. Исходные данные к примеру 5.15

**Проверка гипотезы о равенстве дисперсий двух нормальных распределений.** В качестве границ критической области выступают квантили  $f_\gamma(l, k)$  распределения Фишера (см. (5.57) или (5.59)). Для вычисления этих квантилей используется функция ФРАСПОБР, обращение к которой имеет вид:

$$=\text{ФРАСПОБР}(\text{вероятность}; \text{степень1}; \text{степень2}),$$

где *вероятность* – уровень значимости  $\alpha$  при построении правосторонней критической области; *степень1* – число степеней свободы  $l$ ; *степень2* – число степеней свободы  $k$ .

Граница  $x_{np, \alpha}$  правосторонней критической области (см. (5.57)) вычисляется с помощью выражения

$$x_{np, \alpha} = \text{ФРАСПОБР}(\alpha; l; k).$$

*Решение.* Обратимся к режиму **Двухвыборочный  $t$ -тест с одинаковыми дисперсиями**. В появившемся диалоговом окне зададим необходимые параметры (см. рис. 5.6), а затем щелкнем ОК. Результаты работы режима показаны на рис. 5.7 ( $t$ -статистика является наблюдаемым значением критерия (5.46):  $K_{наб} = 3.58$ ). Это значение попадает в критическую область  $(-\infty, -2.09] \cup [2.09, \infty)$ . Действительно,  $|K_{наб}| > |t_{кр}| = 2.09$ . Следовательно, нулевая гипотеза  $a_X = a_Y$  с уровнем значимости 0.05 отвергается и принимается альтернативная гипотеза  $a_X \neq a_Y$ . ☹

Граница  $x_{np, \alpha/2}$  при построении двухсторонней критической области вычисляется с помощью выражения

$$x_{np, \alpha/2} = \text{ФРАСПОБР}(\alpha/2; l; k).$$

Проверить гипотезу о равенстве дисперсий двух случайных величин  $X \square N(a_X, \sigma_X)$ ,  $Y \square N(a_Y, \sigma_Y)$  можно с использованием режима **Двухвыборочный  $F$ -тест для дисперсии**. Для вызова режима необходимо обратиться к пункту **Сервис** строки меню Excel, команде **Пакет анализа**. Затем в появившемся списке режимов выбрать данный режим и щелкнуть ОК. В появившемся диалоговом окне этого режима задаются следующие параметры (рис. 5.9):

Рис. 5.9. Задание параметров режима **Двухвыборочный  $F$ -тест для дисперсии**

*Интервал переменной 1:* – адреса ячеек, содержащих выборочные значения случайной величины  $X$ .

*Интервал переменной 2:* – адреса ячеек, содержащих выборочные значения случайной величины  $Y$ .

*Метки* – включается, если первая строка содержит заголовки столбцов.

*Альфа:* – задает уровень значимости  $\alpha$ .



Выходной интервал: / Новый рабочий лист: / Новая рабочая книга – указывается, куда выводятся результаты вычислений. При включении *Выходной интервал*: вводится адрес ячейки, начиная с которой выводятся результаты, представленные в виде таблицы (пример такой таблицы приведен на рис. 5.10).

Двухвыборочный F-тест для дисперсии		
	Старая технология	Новая технология
Среднее	307,111	304,923
Дисперсия	1,611	2,244
Наблюдения	9,000	13,000
df	8,000	12,000
F	0,718	
P(F<=f) одностороннее	0,326	
F критическое одностороннее	0,305	

Рис. 5.10. Результаты работы режима *Двухвыборочный F-тест для дисперсии*

♦ **Пример 5.16.** Выборочные данные о расходе сырья при производстве продукции по старой и новой технологии приведены в столбцах А, В документа Excel (см. рис. 5.8). Предполагая, что расход сырья по старой и новой технологии распределен по нормальному закону, нужно проверить статистическую гипотезу  $\sigma_X^2 = \sigma_Y^2$  при уровне значимости  $\alpha = 0.05$ .

*Решение.* Обратимся к режиму *Двухвыборочный F-тест для дисперсии*. В появившемся диалоговом окне зададим необходимые параметры (см. рис. 5.9), а затем щелкнем ОК. Результаты работы режима показаны на рис. 5.10. Так как  $S_X^2 < S_Y^2$ , то в качестве альтернативной гипотезы  $H_1$  принимаем  $\sigma_X^2 < \sigma_Y^2$  и строим левостороннюю критическую область  $(0, x_{лев, \alpha})$ . Граница  $x_{лев, \alpha} = 0.305$ , а наблюдаемое значение  $K_{наб}$  (5.54) равно 0.73 и не попадает в критическую область. Следовательно, можно принять гипотезу о равенстве дисперсий  $\sigma_X^2 = \sigma_Y^2$  с уровнем значимости  $\alpha = 0.05$ . ☺

## 6. ВОПРОСЫ ДЛЯ САМОКОНТРОЛЯ

### Основные понятия математической статистики

1. Что называется генеральной совокупностью?
2. Что называется выборкой (выборочной совокупностью)?
3. Что называется объемом выборки и выборочными характеристиками?
4. Как определяются повторная выборка и бесповторная выборка?
5. Как определяется простая статистическая совокупность?
6. Как определяется вариационный ряд?
7. Как определяется статистический ряд для дискретной случайной величины?
8. Как производится группирование статистических данных для непрерывной случайной величины?
9. Как строится гистограмма?
10. Какой смысл имеет гистограмма?
11. Какой вид имеет статистическая (эмпирическая) функция распределения?
12. Какие вычисления осуществляет функция Excel ЧАСТОТА?
13. Как построить в Excel гистограмму?
14. Какие вычисления осуществляет функция Excel СЧЁТ?
15. Какая функция Excel вычисляет выборочную дисперсию?

### Статистическое оценивание. Точечная оценка

1. Что такое статистическая оценка и какова ее основная особенность?
2. Какая оценка называется точечной?
3. Как определяется несмещенная оценка и смещенная оценка?
4. Как определяется состоятельная оценка?
5. Как находится точечная оценка математического ожидания?
6. Как формулируются теоремы о несмещенности и состоятельности точечной оценки математического ожидания?
7. Как находится точечная оценка дисперсии случайной величины?
8. Как формулируется теорема о смещенности выборочной дисперсии?
9. Что такое исправленная выборочная дисперсия и исправленное выборочное среднее квадратическое отклонение?
10. Какая функция Excel вычисляет исправленную дисперсию  $S^2$ ?



11. Как осуществить вычисление оценок максимального правдоподобия в табличном процессоре Excel?
12. Что осуществляет команда *Поиск решения*?
13. Можно ли задать априорные ограничения на значения оценок максимального правдоподобия, вычисляемые в табличном процессоре Excel? Если да, то как это осуществить?
14. Назовите функции Excel, осуществляющие вычисление точечных оценок по заданной выборке.
15. Что такое описательные статистики и на какие группы они делятся?
16. Как вычислить описательные статистики в табличном процессоре Excel?

### ***Интервальные оценки неизвестных параметров***

1. Какая оценка называется интервальной?
2. Что называется доверительным интервалом, доверительными границами и доверительной вероятностью?
3. В чем заключается смысл интервальной оценки?
4. Какое распределение используют при интервальном оценивании математического ожидания нормально распределенной случайной величины при известной дисперсии?
5. Какое распределение используют при интервальном оценивании дисперсии нормально распределенной случайной величины?
6. Какое распределение используют при интервальном оценивании математического ожидания нормально распределенной случайной величины при неизвестной дисперсии?
7. Какую величину вычисляет функция Excel ДОВЕРИТ?
8. Какие вычисления осуществляет функция Excel ХИ2ОБР?
9. Какие вычисления осуществляет функция Excel СТЬЮДРАСПОБР?

### ***Проверка статистических гипотез***

1. Что называется критерием, уровнем значимости, критической областью и областью допустимых значений критерия?
2. Что такое ошибки первого и второго рода?
3. Что называется мощностью критерия?
4. Сформулируйте этапы проверки статистических гипотез.
5. Как проверить гипотезу о виде распределения генеральной совокупности?

6. Как проверить гипотезу о равенстве генеральных средних в различных случаях?
7. Как проверить гипотезу о равенстве генеральных дисперсий?
8. Как проверить гипотезу о некоррелированности двух генеральных совокупностей?
9. Проверку какой гипотезы осуществляет функция Excel ZTEST?
10. Как выполнить проверку в табличном процессоре Excel гипотезы о равенстве математических ожиданий при известных дисперсиях?
11. Как выполнить проверку в табличном процессоре Excel гипотезы о равенстве математических ожиданий при неизвестных, но равных дисперсиях?
12. Как выполнить проверку в табличном процессоре Excel гипотезы о равенстве дисперсий двух нормальных распределений?

## **ЗАКЛЮЧЕНИЕ**

В данном учебном пособии были изложены основные методы математической статистики, позволяющие сделать выводы о статистических закономерностях, которым подчиняется изучаемое явление.

Наличие в учебном пособии большого числа рассмотренных типовых примеров позволяет не только лучше усвоить теоретические положения математической статистики, но и успешно использовать методы математической статистики для решения практических задач, возникающих в математико-статистических исследованиях. Включение в учебное пособие фрагментов документов табличного процессора Excel, в которых реализуются алгоритмы решения задач математической статистики, существенно повысит эффективность использования методов математической статистики на практике.

## БИБЛИОГРАФИЧЕСКИЙ СПИСОК

1. **Тимошенко Е. И.** Теория вероятностей : учеб. пособие / Е. И. Тимошенко, Ю. Е. Воскобойников. – Новосибирск : НГАСУ, 2003. – 88 с. (есть электронная версия: <http://www.ngasu.nsk.su/prikl/terver.html>).
2. **Гмурман В. Е.** Теория вероятностей и математическая статистика : учеб. для вузов / В. Е. Гмурман. – 6-е изд., стер. – М. : Высш. шк., 1997. – 479 с.
3. **Смирнов Н. В.** Курс теории вероятностей и математической статистики для технических приложений / Н. В. Смирнов, И. В. Дунин-Барковский. – 3-е изд., стер. – М. : Наука, 1969. – 511 с.
4. **Калинина В. Н.** Математическая статистика : учеб. для техникумов / В. Н. Калинина, В. Ф. Панкин. – М. : Высш. шк., 1994. – 336 с.
5. **Вентцель Е. С.** Теория вероятностей : учеб. для вузов / Е. С. Вентцель. – 5-е изд., стер. – М. : Высш. шк., 1998. – 576 с.
6. **Гмурман В. Е.** Руководство к решению задач по теории вероятностей и математической статистике : учеб. пособие для вузов / В. Е. Гмурман. – 5-е изд., стер. – М. : Высш. шк., 2000. – 400 с.

## ПРИЛОЖЕНИЕ

Таблица П1

Значения функции  $\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_0^x \exp(-z^2/2) dz$

$x$	$\Phi(x)$	$x$	$\Phi(x)$	$x$	$\Phi(x)$	$x$	$\Phi(x)$
0.00	0.0000	0.33	0.1293	0.66	0.2454	0.99	0.3389
0.01	0.0040	0.34	0.1331	0.67	0.2486	1.00	0.3413
0.02	0.0080	0.35	0.1368	0.68	0.2517	1.01	0.3438
0.03	0.0120	0.36	0.1406	0.69	0.2549	1.02	0.3461
0.04	0.0160	0.37	0.1443	0.70	0.2580	1.03	0.3485
0.05	0.0199	0.38	0.1480	0.71	0.2611	1.04	0.3508
0.06	0.0239	0.39	0.1517	0.72	0.2642	1.05	0.3531
0.07	0.0279	0.40	0.1554	0.73	0.2673	1.06	0.3554
0.08	0.0319	0.41	0.1591	0.74	0.2703	1.07	0.3577
0.09	0.0359	0.42	0.1628	0.75	0.2734	1.08	0.3599
0.10	0.0398	0.43	0.1664	0.76	0.2764	1.09	0.3621
0.11	0.0438	0.44	0.1700	0.77	0.2794	1.10	0.3643
0.12	0.0478	0.45	0.1736	0.78	0.2823	1.11	0.3665
0.13	0.0517	0.46	0.1772	0.79	0.2852	1.12	0.3686
0.14	0.0557	0.47	0.1808	0.80	0.2881	1.13	0.3708
0.15	0.0596	0.48	0.1844	0.81	0.2910	1.14	0.3729
0.16	0.0636	0.49	0.1879	0.82	0.2939	1.15	0.3749
0.17	0.0675	0.50	0.1915	0.83	0.2967	1.16	0.3770
0.18	0.0714	0.51	0.1950	0.84	0.2995	1.17	0.3790
0.19	0.0753	0.52	0.1985	0.85	0.3023	1.18	0.3810
0.20	0.0793	0.53	0.2019	0.86	0.3051	1.19	0.3830
0.21	0.0832	0.54	0.2054	0.87	0.3078	1.20	0.3849
0.22	0.0871	0.55	0.2088	0.88	0.3106	1.21	0.3869
0.23	0.0910	0.56	0.2123	0.89	0.3133	1.22	0.3883
0.24	0.0948	0.57	0.2157	0.90	0.3159	1.23	0.3907
0.25	0.0987	0.58	0.2190	0.91	0.3186	1.24	0.3925
0.26	0.1026	0.59	0.2224	0.92	0.3212	1.25	0.3944
0.27	0.1064	0.60	0.2257	0.93	0.3238	1.26	0.3962
0.28	0.1103	0.61	0.2291	0.94	0.3264	1.27	0.3980
0.29	0.1141	0.62	0.2324	0.95	0.3289	1.28	0.3997
0.30	0.1179	0.63	0.2357	0.96	0.3315	1.29	0.4015
0.32	0.1225	0.65	0.2422	0.98	0.3365	1.31	0.4049

Окончание табл. П1

$x$	$\Phi(x)$	$x$	$\Phi(x)$	$x$	$\Phi(x)$	$x$	$\Phi(x)$
1.32	0.4066	1.69	0.4545	2.12	0.4830	2.86	0.4979
1.33	0.4082	1.70	0.4554	2.14	0.4838	2.88	0.4980
1.34	0.4099	1.71	0.4564	2.16	0.4846	2.90	0.4981
1.35	0.4115	1.72	0.4573	2.18	0.4854	2.92	0.4982
1.36	0.4131	1.73	0.4582	2.20	0.4861	2.94	0.4984
1.37	0.4137	1.74	0.4591	2.22	0.4868	2.96	0.4985
1.38	0.4162	1.75	0.4599	2.24	0.4875	2.98	0.4986
1.39	0.4177	1.76	0.4608	2.26	0.4881	3.00	0.49865
1.40	0.4192	1.77	0.4616	2.28	0.4887	3.20	0.49931
1.41	0.4207	1.78	0.4625	2.30	0.4893	3.40	0.49966
1.42	0.4222	1.79	0.4633	2.32	0.4898	3.60	0.499841
1.43	0.4236	1.80	0.4641	2.34	0.4904	3.80	0.499928
1.44	0.4251	1.81	0.4649	2.36	0.4909	4.00	0.499968
1.45	0.4265	1.82	0.4656	2.38	0.4913	4.50	0.499997
1.46	0.4279	1.83	0.4664	2.40	0.4918	5.00	0.499997
1.47	0.4292	1.84	0.4671	2.42	0.4922		
1.48	0.4306	1.84	0.4678	2.44	0.4927		
1.49	0.4319	1.86	0.4686	2.46	0.4931		
1.50	0.4332	1.87	0.4693	2.48	0.4934		
1.51	0.4345	1.88	0.4699	2.50	0.4938		
1.52	0.4357	1.89	0.4706	2.52	0.4938		
1.53	0.4370	1.90	0.4713	2.54	0.4945		
1.54	0.4382	1.91	0.4719	2.56	0.4948		
1.55	0.4394	1.92	0.4726	2.58	0.4951		
1.56	0.4406	1.93	0.4732	2.60	0.4953		
1.57	0.4418	1.94	0.4738	2.62	0.4956		
1.58	0.4429	1.95	0.4744	2.64	0.4959		
1.59	0.4441	1.96	0.4750	2.66	0.4961		
1.60	0.4452	1.97	0.4756	2.68	0.4961		
1.61	0.4463	1.98	0.4761	2.70	0.4963		
1.62	0.4474	1.99	0.4767	2.72	0.4965		
1.63	0.4484	2.00	0.4772	2.74	0.4967		
1.64	0.4495	2.02	0.4783	2.76	0.4971		
1.65	0.4505	2.04	0.4793	2.78	0.4973		
1.66	0.4515	2.06	0.4803	2.80	0.4974		
1.68	0.4535	2.10	0.4821	2.84	0.4977		

Таблица П2

Таблица значений  $t(\gamma, n)$ , определяемых выражением $P(|T_n| < t(\gamma, n)) = \gamma$ , где  $n$  – объем выборки

$\gamma \backslash n$	0.95	0.99	0.999	$\gamma \backslash n$	0.95	0.99	0.999
<b>5</b>	2.78	4.6	8.61	<b>20</b>	2.093	2.861	3.883
<b>6</b>	2.57	4.03	6.86	<b>25</b>	2.064	2.797	3.745
<b>7</b>	2.45	3.71	5.96	<b>30</b>	2.045	2.756	3.659
<b>8</b>	2.37	3.50	5.41	<b>35</b>	2.032	2.720	3.600
<b>9</b>	2.31	3.36	5.04	<b>40</b>	2.023	2.0708	3.558
<b>10</b>	2.26	3.25	4.78	<b>45</b>	2.016	2.692	3.527
<b>11</b>	2.23	3.17	4.59	<b>50</b>	2.009	2.679	3.502
<b>12</b>	2.20	3.11	4.44	<b>60</b>	2.001	2.662	3.464
<b>13</b>	2.18	3.06	4.32	<b>70</b>	1.996	2.649	3.439
<b>14</b>	2.16	3.01	4.22	<b>80</b>	1.991	2.640	3.418
<b>15</b>	2.15	2.98	4.14	<b>90</b>	1.987	2.633	3.403
<b>16</b>	2.13	2.95	4.07	<b>100</b>	1.984	2.627	3.3392
<b>17</b>	2.12	2.92	4.02	<b>120</b>	1.980	2.617	3.374
<b>18</b>	2.11	2.90	3.97	$\infty$	1.960	2.576	3.291
<b>19</b>	2.10	2.88	3.92				

Таблица ПЗ

Таблица значений квантилей  $\chi_k^2$ -распределения,

определяемых соотношением

$$P(\chi_k^2 < \chi^2(\gamma, k)) = \gamma$$

$k \backslash \gamma$	0.02	0.05	0.1	0.9	0.95	0.98
1	0.006	0.0039	0.016	2.7	3.8	5.4
2	0.040	0.103	0.211	4.6	6.0	7.8
3	0.185	0.352	0.584	6.3	7.8	9.8
4	0.43	0.71	1.06	7.8	9.5	11.7
5	0.75	1.14	1.61	9.2	11.1	13.4
6	1.13	1.63	2.20	10.6	12.6	15.0
7	1.56	2.17	2.83	12.0	14.1	16.6
8	2.03	2.73	3.49	13.4	15.5	18.2
9	2.53	3.32	4.17	14.7	16.9	19.7
10	3.06	3.94	4.86	16.0	18.3	21.2
12	4.2	5.2	6.3	18.5	21.0	24.1
14	5.4	6.6	7.8	21.1	23.7	26.9
16	6.6	8.0	9.3	23.5	26.3	29.6
18	7.9	9.4	10.9	26.0	28.9	32.3
20	9.2	10.9	12.4	28.4	31.4	35.0
22	10.6	12.3	14.0	30.8	33.9	37.7
24	12.0	13.8	15.7	33.2	36.4	40.3
26	13.4	15.4	17.3	35.6	38.9	42.9
28	14.8	16.9	18.9	37.9	41.3	45.4
30	16.3	18.5	20.6	40.3	43.8	48.0

Таблица П4

Доверительные границы  $p_2$  и  $p_1$   
 для вероятности  $p$  при  $\gamma = 0.95$  (значения  $p_2$  приведены  
 в верхней строке,  $p_1$  – в нижней)

$n-m \backslash m$	1	2	3	4	5	6	7	8
0	0.975	0.842	0.708	0.602	0.522	0.459	0.410	0.369
	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
1	0.987	0.906	0.806	0.716	0.641	0.579	0.527	0.463
	0.013	0.008	0.006	0.005	0.004	0.004	0.003	0.003
2	0.992	0.932	0.853	0.727	0.710	0.651	0.600	0.556
	0.094	0.068	0.053	0.037	0.037	0.032	0.028	0.028
3	0.094	0.947	0.882	0.816	0.755	0.701	0.652	0.610
	0.194	0.147	0.118	0.099	0.085	0.075	0.067	0.060
4	0.995	0.957	0.901	0.843	0.788	0.738	0.692	0.651
	0.284	0.223	0.184	0.157	0.137	0.122	0.109	0.099

Таблица П5

Квантили  $f_\gamma(l, k)$  распределения Фишера,

определяемые уравнением

$$P(F(l, k) < f_\gamma(l, k)) = \gamma = 0.95$$

( $l$  – степени свободы для большей дисперсии,

$k$  – для меньшей дисперсии)

$\begin{smallmatrix} l \\ k \end{smallmatrix}$	1	2	3	4	6	8	12	24
1	161.4	199.5	215.7	224.6	234.0	238.9	243.9	249.0
2	18.51	19.00	19.16	19.25	19.33	19.37	19.41	19.45
3	10.13	9.55	9.28	9.21	8.84	8.82	8.74	8.64
4	7.71	6.94	5.59	6.39	6.16	6.04	5.91	5.77
5	6.61	5.79	5.41	5.19	4.95	4.82	4.68	4.53
6	5.99	5.14	4.76	4.53	4.88	4.15	4.00	3.84
7	5.59	4.74	4.35	4.12	3.87	3.73	3.57	3.41
8	5.32	4.46	4.07	3.84	3.58	3.44	3.28	3.12
9	5.12	4.26	3.86	3.63	3.37	3.23	3.07	2.90
10	4.96	4.10	3.71	3.48	3.22	3.07	2.91	2.74
12	4.75	3.88	3.49	3.26	3.00	2.85	2.69	2.50
14	4.60	3.74	3.34	3.11	2.85	2.70	2.53	2.35
16	4.49	3.63	3.24	3.01	2.74	2.59	2.42	2.24
18	4.41	3.55	3.16	2.93	2.66	2.51	2.34	2.15
20	4.35	3.49	3.10	2.87	2.60	2.45	2.28	2.08
22	4.30	3.44	3.05	2.82	2.55	2.40	2.23	2.03
24	4.26	3.40	3.01	2.78	2.51	2.36	2.18	1.98
26	4.22	3.37	2.98	2.74	2.47	2.32	2.15	1.95
28	4.20	3.34	2.95	2.71	2.44	2.29	2.12	1.91
30	4.17	3.32	2.92	2.69	2.42	2.27	2.09	1.89
40	4.08	3.23	2.84	2.61	2.34	2.18	2.00	1.79
60	4.00	3.15	2.76	2.52	2.25	2.10	1.92	1.70