

## Установочная лекция по статистике

# Содержание

<b>1</b>	<b>Введение в статистику</b>	<b>2</b>
1.1	Самое-самое введение и мотивировка . . . . .	2
1.2	Немного о выборке . . . . .	6
<b>2</b>	<b>Оценки характеристик распределения</b>	<b>8</b>
2.1	Немного об оценках . . . . .	8
2.2	Эмпирическое распределение . . . . .	9
2.3	Оценки параметров распределения . . . . .	10
2.4	Немного о качестве оценок . . . . .	12
2.5	Параметрические модели . . . . .	13
2.6	А как сравнивать оценки между собой? . . . . .	14
2.7	А существует ли наилучшая оценка? . . . . .	16
2.8	А как вообще получать оценки? . . . . .	17
<b>3</b>	<b>Интервальное оценивание</b>	<b>17</b>
3.1	Немного об интервальном оценивании . . . . .	17
<b>4</b>	<b>Проверка гипотез</b>	<b>19</b>
4.1	А что такое гипотезы? . . . . .	19
4.2	Проверка гипотез . . . . .	19
4.3	Сравнение критериев . . . . .	21

# 1 Введение в статистику

## 1.1 Самое-самое введение и мотивировка

Здравствуйтесь, уважаемые слушатели. В этой лекции мы осветим те базовые понятия статистики, которые будут подробно изучаться в дальнейшем, а также покажем, какие задачи решает математическая статистика.

При решении различных прикладных задач, при наблюдении тех или иных явлений, оказывается, что многие закономерности носят вероятностный характер. Если быть более точным, вероятностная модель этих закономерностей позволяет получить разумные результаты и выводы. Как же адекватно формализовать фразу: «носят вероятностный характер», что это означает? Это означает, что изучаемая нами закономерность (явление) случайна в своих проявлениях и, тем самым, описывается случайной величиной, которая, в свою очередь, имеет вероятностное распределение. Понятие распределения было введено нами в первой части курса, касающейся теории вероятностей, а посему мы знаем, что знание распределения – это очень много, ведь распределение, грубо говоря, говорит нам, что, когда и с какой вероятностью может произойти с величиной  $\xi$ .

Кроме того, можно вычислять такие полезные штуки, как: математическое ожидание  $E\xi$ , дисперсию  $D\xi = E(\xi - E\xi)^2$ , медиану, вероятности попадания в те или иные множества, плотность  $f_\xi$  и многое-многое другое. Иными словами, зная распределение интересующей нас случайной величины, мы сразу знаем очень много, а также можем получать сразу кучи полезной информации.

Давайте приведем пример. Предположим, что случайная величина  $\xi$  – это рост взрослого человека. Ясное дело, что он случаен и меняется от человека к человеку. В то же время оказывается, что он описывается некоторым распределением. Кстати, это распределение хорошо аппроксимируется нормальным распределением (конечно, в разумных пределах, так как нормально распределенная случайная величина может принимать и отрицательные значения, и сколь угодно большие, что для роста, конечно, невозможно). Распределение же характеризуется, например, плотностью  $f_\xi$ . Ну а зная плотность, можно вычислить и средний рост, и среднеквадратическое отклонение, и даже вероятность встретить дядю Стёпу (ну она-то равна нулю, это понятно). Но и тут нужно быть осторожным в интерпретации результатов: параметры этого распределения меняются, как минимум, от страны к стране, от национальности к национальности, ну и так далее. Кроме того, ясное дело, средний рост женщин и мужчин отличается, а значит распределения роста по полу, приближаемые нормальными  $N_{a,\sigma^2}$ , приближаются разными нормальными, то есть нормальными, но с разными параметрами математического ожида-

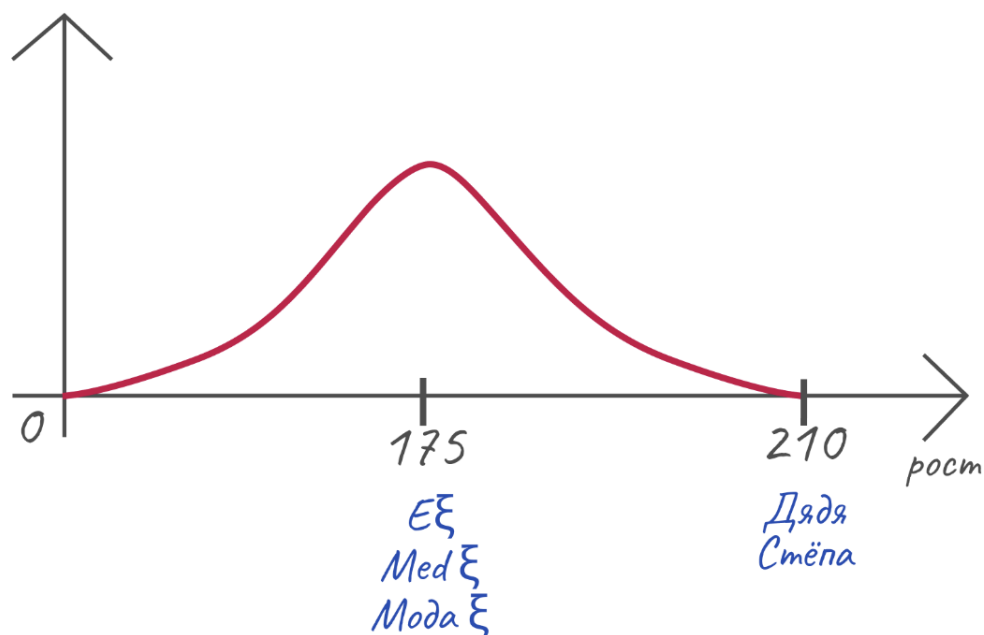


Рис. 1: Нормальное распределение

ния  $a$  и дисперсии  $\sigma^2$ ! Поэтому важно, чтобы наши выводы, сделанные по выборке, касались именно того распределения, из которого эта выборка. Будет очень неправильно, если по выборке из роста мужчин будут сделаны выводы по пошиву одежды для женщин.

Ну а зачем нам знать эти параметры? О чем они нам говорят? В случае нормального распределения мода (локальный максимум плотности), медиана и математическое ожидание совпадают, и все могут быть интерпретированы, как среднее значение случайной величины.  $\sigma^2$  – это дисперсия, а значит  $\sigma$  – среднеквадратическое отклонение. Иными словами, в среднем, с вероятностью примерно 0.68, значения случайной величины находятся в пределах

$$(a - \sigma, a + \sigma) = (E\xi - \sigma_\xi, E\xi + \sigma_\xi).$$

И что же это значит? А это значит, что одежды, с размерами роста из приведенного интервала имеет смысл шить куда больше, чем с размерами вне этого интервала – она просто не будет продаваться в достаточном объеме! А что если взять не  $\sigma$ , а  $2\sigma$ ? Тогда вероятность того, что наша случайная величина попадет в интервал

$$(a - 2\sigma, a + 2\sigma)$$

примерно равна 0.96. Вообще, для нормального распределения есть так называемое правило  $3\sigma$ . Оказывается, что вероятность случайной величине  $\xi$ , имеющей нормальное распределение  $N_{a, \sigma^2}$ , попасть в интервал  $(a - 3\sigma, a + 3\sigma)$  равна 0.997, что помнить вовсе необязательно, но что неплохо осознавать: вероятность чрезвычайно велика!

Для других распределений ситуация, вообще говоря, не такая. Например, для случайной величины  $\xi$ , имеющей показательное распределение  $\text{Exp}_\lambda$ , мода – это  $x = 0$ , математическое ожидание – это  $E\xi = \frac{1}{\lambda}$  и медиана  $\text{med}\xi = \frac{\ln 2}{\lambda}$ . Все три величины различны. А тогда что и для какой оценки

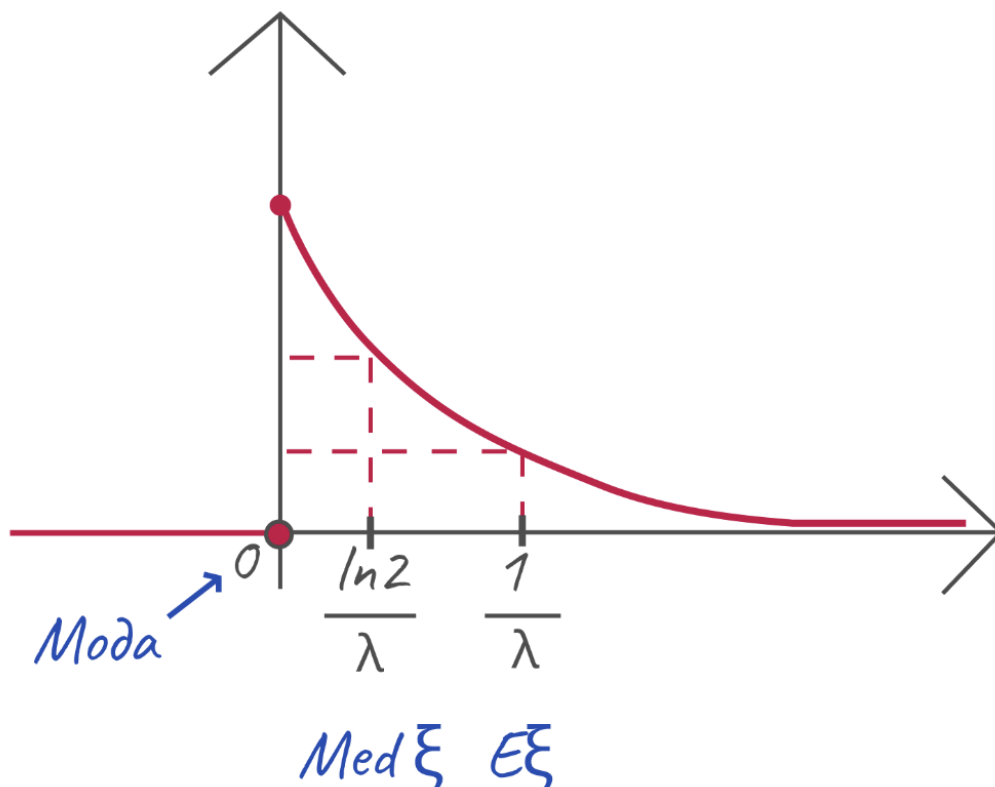


Рис. 2: Показательное распределение

выбирать? Все зависит от задачи. Мы уже вспоминали задачу про зарплаты, напомним ее.

Предположим, что имеется некоторая фирма, в которой работает 100 человек, один из которых начальник. Зарботная плата начальника равна 101000 долларов в месяц, а зарботная плата каждого работника равна 1000 долларов в месяц. Пусть случайная величина  $\xi$  – зарплата работника, тогда ее распределение может быть задано, как

$\xi$	1000	101000
P	$\frac{99}{100}$	$\frac{1}{100}$

Легко понять, что математическое ожидание случайной величины  $\xi$  равно

$$E\xi = 1000 \cdot \frac{99}{100} + 101000 \cdot \frac{1}{100} = 2000,$$

то есть средняя зарплата равна 2000 долларов в месяц. В то же время, медиана равна 1000 и медианная зарплата равна 1000. Ясно, что в этом примере

гораздо более «честной» является медиана, нежели математическое ожидание, из-за такого сильного выброса в заработной плате начальника.

В принципе, можно сказать, что нам хочется найти или оценить плотность  $f_\xi$  случайной величины  $\xi$ , или ее дискретный аналог – ряд распределения, потому что дальше все хорошо: суммы-интегралы дают какие-то числа, которые уже изучены вдоль и поперек, и даже в этом небольшом фрагменте прозвучали не раз. Но плотность нам нужна вовсе не только для этого. По плотности можно понять: а однородно ли распределение?

Не вдаваясь в детали, можно сказать, что однородное распределение – это распределение с одной модой (унимодальные). Давайте все же сформулируем определение моды в одном месте.

**Определение 1.1.1** *Модой случайной величины  $\xi$ , имеющей абсолютно непрерывное распределение с плотностью  $f_\xi$ , называется точка локального максимума  $f_\xi$ .*

*Если же случайная величина  $\xi$  имеет дискретное распределение, заданное таблицей*

$$\begin{array}{c|c|c|c|c|c} \xi & a_1 & a_2 & \dots & a_n & \dots \\ \hline P & p_1 & p_2 & \dots & p_n & \dots \end{array},$$

*то модой случайной величины  $\xi$  называется произвольное значение  $a_i$  такое, что*

$$p_{i-1} \leq p_i, \quad p_{i+1} \leq p_i,$$

*если хотя бы одна из соседствующих вероятностей определена.*

Рисунки иллюстрируют введенные определения. Для дискретного распреде-

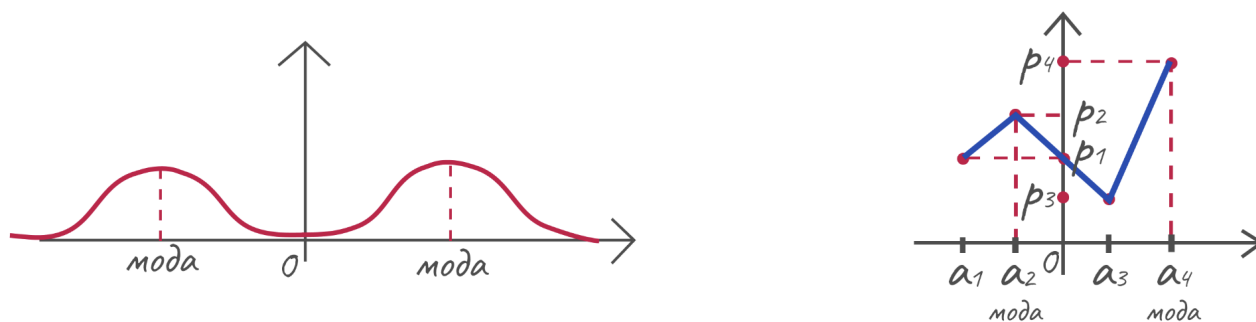


Рис. 3: Моды

ления часто вводят понятие многоугольника распределения (тогда узловые точки  $(a_i, p_i)$  видны куда лучше), он тоже изображен на рисунке. Многоугольник получается просто соединением отрезком соседних точек  $(a_i, p_i)$  и  $(a_{i+1}, p_{i+1})$ .

Легко понять (кроме как, возможно, для биномиального распределения  $\text{Bin}_{n,p}$ ), что почти что все распределения, которые мы рассматривали: вырожденное, Бернулли, биномиальное, Пуассона, показательное, нормальное (короче все, кроме равномерного и левого хвоста показательного) – унимодальные. А что вообще может означать, что распределение не однородно и имеет несколько мод? Часто это означает, что рассматривается смесь унимодальных распределений (линейная комбинация их плотностей). Не углубляясь в детали, поясним, что смесь – это когда с какой-то вероятностью  $p$  проявляется один закон распределения, а с обратной, то есть с вероятностью  $(1 - p)$  – другой закон распределения. У смеси плотность равна линейной комбинации плотностей этих распределений с коэффициентами  $p$  и  $(1 - p)$ . Например, можно смешать крокодилов и черепах и рассматривать распределение смеси их весов. Если есть неоднородность (то есть смесь разного), то такое достаточно бессмысленно исследовать: зачем кому-то нужен средний вес (математическое ожидание) смеси крокодилов и черепах?

В общем, много всего интересного нам хочется узнать про распределение. И тут мы сталкиваемся с главной проблемой: а как это сделать? Нам же не дано распределение, мы наблюдаем только какие-то значения случайной величины. Иными словами, есть проявления некоторой вероятностной закономерности, по которым нужно узнать про генеральное распределение, но формулки-то нет: нет ни плотности, ни ряда распределения. Тут и возникает математическая статистика.

## 1.2 Немного о выборке

Итак, пусть у нас имеется  $n$  проявлений некоторой вероятностной закономерности. Иными словами, мы наблюдаем  $n$  значений некоторой случайной величины  $\xi$ , называемой еще генеральной совокупностью, имеющей какое-то (неизвестное нам) распределение. Тогда перед нами обычный числовой  $n$ -мерный вектор  $X = (x_1, x_2, \dots, x_n)$  – это выборка после эксперимента. Итак,

**Определение 1.2.1** Пусть  $\xi$  – рассматриваемая нами случайная величина. Выборкой (после эксперимента)  $X = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n$  называется  $n$  независимых реализаций случайной величины  $\xi$ . Последнюю часто называют генеральной совокупностью.

Например, пусть  $\xi$  – это рост случайно взятого человека, тогда

$$X = (175, 182, 168, 155, 192)$$

это выборка объема 5 из генеральной совокупности  $\xi$ . Хорошо, допустим теперь мы хотим вычислить среднее, а точнее прикинуть, чему равно математическое ожидание  $\xi$ . Как бы это сделать? Наверное, вы скажете, что логично

посчитать среднее арифметическое, называемое выборочным средним:

$$\bar{X} = \frac{175 + 182 + 168 + 155 + 192}{5} = 174.4.$$

И как, все готово? То есть мы получили математическое ожидание  $\xi$ ? Немного подумав, наверное, можно бросить осторожную фразу вроде: «Видимо, это оценка этого математического ожидания»? И что это значит? Это значит, что, поменяв выборку, мы получим, вообще говоря, другое число. Видимо, это число имеет какое-то отношение к реальному, но какое? Каждый раз меняя «измеряемых» людей, мы получаем другое число. И как тут строить теорию? Опять что-то усреднять? Так можно циклиться до бесконечности.

Важно еще понимать и вот какой момент, называемый неоднородностью. Придя в детский сад и померив людей в нем, получим одно число (вероятно, сравнительно небольшое, даже рост воспитателей и нянечек не поможет), а измерив рост членов баскетбольной команды – разительно отличающееся. Дело в том, что выборки-то берутся из разных распределений! Поэтому мы всегда должны понимать, какое именно распределение мы изучаем по данной выборке. А то получится вроде ситуации с опросом в интернете с вопросом: «Пользуетесь ли вы интернетом»?

Ну ладно, давайте подумаем еще немного. Все-таки  $\xi$  – случайная величина, а значит она определена на каком-то вероятностном пространстве. Значит, меняя элементарные исходы пространства элементарных исходов, эта величина  $\xi$  принимает, вообще говоря, различные значения. У нас элементарные исходы – это люди, потому взяв разных людей, получим, вообще говоря, разные значения роста. Конечно, можно быть совсем занудным: нет людей одинакового роста, все зависит от точности измерений, но не будем об этом, и не будем приводить такие оговорки в дальнейшем. В итоге, разумно привести общее определение выборки.

**Определение 1.2.2** Пусть  $\xi$  – рассматриваемая нами случайная величина. Выборкой  $X = (X_1, X_2, \dots, X_n)$  называется  $n$  независимых случайных величин, имеющих распределение такое же, как и  $\xi$ .

Итак, введенное определение выборки – это, по сути, определение выборки до конкретного эксперимента. В нашем примере, до того, как мы измерим рост конкретных  $n$  людей, выборка – это случайный вектор (вектор из  $n$  независимых случайных величин), а после – это набор из  $n$  чисел  $X = (x_1, x_2, \dots, x_n)$ . Именно на основе этого понятия выборки, как случайного вектора, и строится вся теория математической статистики. В дальнейшем мы не будем конкретизировать то, что мы понимаем под выборкой (числовой вектор или вектор из случайных величин), это должно будет быть понятно из контекста. В принципе, к выборке логично относиться, как к некоторой функции от  $\omega$



до реализации на конкретном элементарном исходе  $\omega$ , и к «значению» этой функции на конкретном  $\omega$  после эксперимента. Ну что, давайте перейдем к задачам математической статистики.

## 2 Оценки характеристик распределения

### 2.1 Немного об оценках

Пусть  $\theta$  – некоторая характеристика распределения случайной величины  $\xi$ , например, пусть  $\theta$  – математическое ожидание  $E\xi$  случайной величины  $\xi$  (о точечном оценивании более строго мы будем говорить позднее). Мы бы хотели оценить этот параметр, используя выборку. Статистика  $\hat{\theta}$  – это некоторая функция от выборки  $X = (X_1, X_2, \dots, X_n)$ , то есть

$$\hat{\theta} = \hat{\theta}(X_1, X_2, \dots, X_n).$$

Например,

$$\hat{\theta} = \bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n} = \frac{1}{n} \sum_{i=1}^n X_i$$

это статистика. Так как у выборки есть два определения, то и статистика может пониматься в двух смыслах. Если выборка рассмотрена на конкретном эксперименте, то  $\hat{\theta}$  – это конкретное число:

$$\hat{\theta} = \hat{\theta}(x_1, x_2, \dots, x_n).$$

В нашем примере на конкретной выборке  $X = (x_1, x_2, \dots, x_n)$  в результате эксперимента получаем

$$\hat{\theta} = \bar{X} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{1}{n} \sum_{i=1}^n x_i.$$

Впрочем, это выражение для нас уже не ново, ведь чуть раньше мы уже считали что-то вроде «среднего роста» по данной выборке. Мы не зря рассматриваем этот пример с выборочным средним, он нам потребуется и далее.

Начали мы с вопроса оценки математического ожидания и, видимо, хочется сказать, что именно выборочное среднее  $\bar{X}$  является хорошей оценкой истинного математического ожидания  $E\xi$  генеральной совокупности  $\xi$ . А что значит, что оценка «хорошая»? И почему вообще такая оценка «естественна»? Откуда она берется?

## 2.2 Эмпирическое распределение

На самом деле, с каждой конкретной выборкой, полученной в результате эксперимента, то есть с выборкой  $X = (x_1, x_2, \dots, x_n)$ , разумно связать новую случайную величину  $\xi^*$ . Какую, спросите вы? Ну, вроде бы ясно какую: такую, которая каждое значение выборки  $x_i$  принимает с вероятностью  $\frac{1}{n}$  — у нас же нет любимчиков в выборке. В итоге, можно написать следующую табличку

$$\begin{array}{c|c|c|c|c} \xi^* & x_1 & x_2 & \dots & x_n \\ \hline \tilde{P} & \frac{1}{n} & \frac{1}{n} & \dots & \frac{1}{n} \end{array}.$$

Кстати, как легко видеть, случайная величина, а точнее ее распределение, конечно, зависит от  $n$ , от объема выборки, но мы не будем наделять ее дополнительным индексом.

Важно отметить, что случайная величина  $\xi^*$ , конечно, задана на другом вероятностном пространстве, нежели  $\xi$ , потому вероятность на ее вероятностном пространстве будем обозначать  $\tilde{P}$ . Например, имея в результате эксперимента выборку

$$X = (1, 2, 1, 3, 1)$$

объема 5, соответствующее распределение новой, разыгранной случайной величины  $\xi^*$  записывается, как

$$\begin{array}{c|c|c|c} \xi^* & 1 & 2 & 3 \\ \hline \tilde{P} & \frac{3}{5} & \frac{1}{5} & \frac{1}{5} \end{array}.$$

Как же вычислить вероятность попадания случайной величины  $\xi^*$  в множество  $A$ ? Да как обычно, ровно как мы и поступаем в дискретном случае: после эксперимента, на конкретной числовой выборке, вероятность попасть в  $A$  вычисляется, как

$$\tilde{P}(\xi^* \in A) = \frac{|\{x_i : x_i \in A\}|}{n},$$

иными словами, вероятность события  $A$  — это отношение количества элементов выборки, лежащих в  $A$ , к объему выборки  $n$ .

Если же выборка «абстрактная», то есть эксперимент не произошел, то  $\tilde{P}(\xi^* \in A)$  — это случайная величина, ведь

$$\tilde{P}(\xi^* \in A) = \frac{|\{X_i : X_i \in A\}|}{n}$$

и справа стоит функция от выборки.

Итак, хотелось бы, чтобы распределение представленной нами случайной величины неплохо приближало распределение настоящей генеральной

совокупности  $\xi$  с ростом  $n$ . И интуиция не подводит, а именно, справедлива теорема, которую мы готовы сформулировать. На самом деле, чтобы не мучить вас всякими  $\sigma$ -алгебрами и на этом этапе не затуманивать суть дела, давайте считать, что  $A$  – произвольное подмножество  $\mathbb{R}$ , и  $A$  – событие.

**Теорема 2.2.1 (ЗБЧ для эмпирического распределения)** *Для каждого множества  $A \subset \mathbb{R}$  имеет место сходимость почти наверное:*

$$\frac{|\{X_i : X_i \in A\}|}{n} \xrightarrow[n \rightarrow +\infty]{\text{п. н.}} P(\xi \in A).$$

**Доказательство.** Доказательство рассмотренной теоремы немедленно следует из закона больших чисел. Ясно, что вероятность  $\xi^*$  попасть в множество  $A$  – это отношение количества элементов  $X_i$ , попавших в  $A$ , к общему числу элементов  $n$ .

$$\tilde{P}(\xi^* \in A) = \frac{|\{X_i : X_i \in A\}|}{n} = \frac{1}{n} \sum_{i=1}^n I(X_i \in A),$$

где  $I$  – индикатор события  $A$ , то есть  $I(X \in A) = 1$ , если  $X \in A$ , и ноль иначе. Все эти индикаторы ( $A$  они, кстати, случайные величины (!)) одинаково распределены, независимы и имеют распределение Бернулли с параметром  $p = P(X_i \in A) = P(\xi \in A)$ . Значит, и математическое ожидание слагаемых равно  $P(\xi \in A)$ , а тогда по ЗБЧ в форме Колмогорова

$$\tilde{P}(\xi^* \in A) = \frac{1}{n} \sum_{i=1}^n I(X_i \in A) \xrightarrow[n \rightarrow +\infty]{\text{п. н.}} P(\xi \in A).$$

□

Итак, наше распределение поточечно почти всюду (или почти наверное) сходится к истинному. Иными словами, при увеличении  $n$  эмпирическое распределение все более и более похоже на истинное.

## 2.3 Оценки параметров распределения

Наверное, теперь ясна основная идея: может быть, в качестве оценки некоторой характеристики генеральной совокупности  $\xi$ , можно использовать соответствующую характеристику  $\xi^*$ , раз распределение последней так неплохо приближает истинное? Попробуем.

Тогда в качестве оценки  $\widehat{E\xi}$  математического ожидания  $\xi$  логично взять соответствующее математическое ожидание  $\tilde{E\xi^*}$ , то есть, по сути, выборочное среднее:

$$\widehat{E\xi} = \tilde{E\xi^*} = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i.$$

В качестве оценки  $\widehat{D\xi}$  дисперсии генеральной совокупности  $\xi$  логично взять дисперсию разыгранной случайной величины:

$$\widehat{D\xi} = \widetilde{D\xi}^* = \widetilde{E}(\xi^* - \widetilde{E\xi}^*)^2 = S^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

Аналогичным образом можно построить оценки и других моментов, но сейчас мы этого делать не будем.

Для построения оценки медианы, нам потребуется понятие вариационного ряда. Упорядочим выборку по порядку, получим новый набор чисел (или случайных величин)  $X_{(i)}$ ,  $i \in \{1, \dots, n\}$ :

$$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)},$$

где, например,  $X_{(1)} = \min(X_1, X_2, \dots, X_n)$ ,  $X_{(n)} = \max(X_1, X_2, \dots, X_n)$ . Полученный таким образом набор чисел и называется вариационным рядом. В вариационном ряде  $i$ -ый член называется  $i$ -ой порядковой статистикой. Рангом элемента выборки  $X_i$  называется его номер в вариационном ряде. Логично взять в качестве оценки медианы случайной величины  $\xi$  либо средний член вариационного ряда, если он определен однозначно (то есть если  $n$  нечетно) и, скажем, полусумму центральных членов вариационного ряда, если  $n$  четно. Иными словами,

$$\widehat{\text{med}} \xi = X_{((n+1)/2)}, \text{ при } n \text{ нечетном,}$$

$$\widehat{\text{med}} \xi = \frac{1}{2} (X_{(n/2)} + X_{(n/2+1)}), \text{ при } n \text{ четном.}$$

Вернемся к нашему синтетическому примеру. Напомним, что выборка из генеральной совокупности  $\xi$  такова

$$X = (1, 2, 1, 3, 1),$$

а распределение случайной величины задается таблицей

$\xi^*$	1	2	3
$\widetilde{P}$	$\frac{3}{5}$	$\frac{1}{5}$	$\frac{1}{5}$

Ясно, что оценкой математического ожидания генеральной совокупности  $\xi$  в нашем случае оказывается

$$\widehat{E\xi} = \widetilde{E\xi}^* = 1 \cdot \frac{3}{5} + 2 \cdot \frac{1}{5} + 3 \cdot \frac{1}{5} = \bar{X} = \frac{1 + 2 + 1 + 3 + 1}{5} = 1.6.$$

Что с дисперсией? Давайте оценим ее. Аналогичным образом получаем, что

$$\widehat{D\xi} = \widetilde{D\xi}^* = S^2 = \frac{16}{25}.$$

И, наконец, для оценивания медианы составим вариационный ряд. Он примет вид

$$1 \leq 1 \leq 1 < 2 < 3.$$

Так как объем выборки равен 5, а это число нечетное, то  $\widehat{\text{med}} \xi = X_{((5+1)/2)} = X_{(3)} = 1$ .

## 2.4 Немного о качестве оценок

Хорошо, мы даже научились получать какие-то оценки, и поняли, из каких соображений их можно получать. А как понять, что они хорошие? Пока что мы смогли обосновать «хорошесть» только эмпирического распределения.

Можно измерить то, насколько ошибка хорошая, используя среднеквадратическую ошибку:

$$\text{MSE} = \mathbb{E} \left( \hat{\theta} - \theta \right)^2.$$

Распишем последнее выражение подробнее, раскрыв скобки и учитывая, что  $\theta \in \mathbb{R}$ :

$$\begin{aligned} \mathbb{E} \left( \hat{\theta} - \theta \right)^2 &= \mathbb{E} \left( \hat{\theta}^2 - 2\hat{\theta}\theta + \theta^2 \right) = \mathbb{E}\hat{\theta}^2 - 2\theta\mathbb{E}\hat{\theta} + \theta^2 = \\ &= \mathbb{E}\hat{\theta}^2 - \left( \mathbb{E}\hat{\theta} \right)^2 + \left( \mathbb{E}\hat{\theta} \right)^2 - 2\theta\mathbb{E}\hat{\theta} + \theta^2 = D\hat{\theta} + \left( \mathbb{E}\hat{\theta} - \theta \right)^2, \end{aligned}$$

откуда

$$\text{MSE} = D\hat{\theta} + \left( \mathbb{E}\hat{\theta} - \theta \right)^2.$$

Ясно, что чтобы оценка была разумной, хотелось бы, чтобы ошибка MSE стремилась к нулю при объеме выборки  $n$ , стремящемся к бесконечности.

Дадим несколько простых определений, более строго к которым подойдем в главе про точечные оценки.

**Определение 2.4.1** Оценка  $\hat{\theta}$  называется несмещенной, если

$$\mathbb{E}\hat{\theta} = \theta.$$

Ясно, что несмещенность, по сути, означает, что «в среднем» значение оценки совпадает с истинным значением параметра. Например, выборочное среднее – это несмещенная оценка математического ожидания  $\xi$ . Действительно,

$$\mathbb{E}\bar{X} = \mathbb{E} \left( \frac{X_1 + X_2 + \dots + X_n}{n} \right) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}X_i = \frac{n\mathbb{E}X_1}{n} = \mathbb{E}X_1 = \mathbb{E}\xi,$$

так как  $X_i$  одинаково распределены. Дадим еще одно определение.

**Определение 2.4.2** Оценка  $\hat{\theta}$  называется асимптотически несмещенной, если

$$\lim_{n \rightarrow +\infty} E\hat{\theta} = \theta.$$

Естественно, любая несмещенная оценка оказывается и асимптотически несмещенной. Кстати, мы уже знакомы со смещенной, но асимптотически несмещенной оценкой – это оценка дисперсии  $S^2$ . Вычисление  $ES^2$  уже не делается в две строчки, поэтому его мы оставим до соответствующей лекции, а пока что просто поверим в ее (оценки) смещенность, но, в то же время, асимптотическую несмещенность.

Ну и зачем нам это все? А выходит, и это совершенно очевидно, что если оценка  $\hat{\theta}$  оказывается несмещенной или асимптотически несмещенной, то из соотношения

$$MSE = D\hat{\theta} + (E\hat{\theta} - \theta)^2.$$

следует, что стремление к нулю  $MSE$  равносильно стремлению к нулю дисперсии оценки  $D\hat{\theta}$  с ростом объема выборки  $n$ .

**Определение 2.4.3** Оценка  $\hat{\theta}$  называется состоятельной в смысле среднеквадратического, если  $MSE$  стремится к нулю с ростом  $n$ .

Давайте проверим, является ли оценка  $\bar{X}$  состоятельной в смысле среднеквадратического? Достаточно вычислить ее дисперсию, предположив, что существует  $D\xi$ :

$$D\bar{X} = D\left(\frac{X_1 + X_2 + \dots + X_n}{n}\right) = \frac{1}{n^2} \sum_{i=1}^n DX_i = \frac{1}{n} DX_1 = \frac{1}{n} D\xi \xrightarrow{n \rightarrow +\infty} 0.$$

Итого, состоятельность в среднеквадратичном доказана.

## 2.5 Параметрические модели

Иногда оказывается, что про распределение случайной величины  $\xi$  что-то известно. Если это учесть, то есть шанс получить что-то более хорошее, нежели работая в режиме «слепых котят». Параметрические модели предполагают знание распределения случайной величины  $\xi$  с точностью до каких-то параметров. Например, может быть известно, что генеральная совокупность имеет нормальное распределение  $N_{a,\sigma^2}$  с неизвестными параметрами  $a$  и  $\sigma^2$ , или равномерное распределение  $U_{a,b}$  с неизвестными параметрами  $a < b$ , или геометрическое распределение  $G_p$  с параметром  $p$ , или, что очень часто, биномиальное распределение  $\text{Bin}_{n,p}$  с параметром  $p$ . И наша задача преобразуется в задачу оценивания конкретных параметров распределения.

Часто параметры модели можно интерпретировать. Так как мы знаем, что у нормального распределения параметр  $a$  отвечает за математическое

ожидание, а параметр  $\sigma^2$  – за дисперсию, то мы уже готовы предложить оценки для этих параметров: выборочное среднее и выборочную дисперсию, а именно

$$\hat{a} = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i,$$

$$\hat{\sigma}^2 = S^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2.$$

Кроме того, мы можем даже найти распределение случайной величины  $\bar{X}$ .

**Лемма 2.5.1** *Выборочное среднее имеет нормальное распределение  $N_{a, \sigma^2/n}$ .*

Почему это так? Да потому, что согласно свойствам нормального распределения, если случайные величины  $X_1$  и  $X_2$  независимы и  $X_1 \sim N_{a_1, \sigma_1^2}$ , а  $X_2 \sim N_{a_2, \sigma_2^2}$ , то  $X_1 + X_2 \sim N_{a_1+a_2, \sigma_1^2+\sigma_2^2}$ . Осталось заметить, что  $X_1 + X_2 + \dots + X_n \sim N_{na, n\sigma^2}$  и воспользоваться свойством линейных преобразований над случайными величинами.

## 2.6 А как сравнивать оценки между собой?

Давайте теперь вспомним, что, вообще-то, математическое ожидание и медиана у случайной величины, имеющей нормальное распределение  $N_{a, \sigma^2}$ , совпадают. Как думаете, что мы хотим предложить? Конечно, еще одну оценку для математического ожидания такой генеральной совокупности – выборочную медиану. Как же теперь их сравнить? Да все так же, по **MSE**:

$$\text{MSE} = D\hat{\theta} + (E\hat{\theta} - \theta)^2.$$

При этом, если оценки несмещенные (ну-ка быстро вспомните, что это!), то достаточно сравнить их дисперсии: логично считать, что чем дисперсия меньше, тем оценка лучше, да? Осознайте, а для этого вспомните, что такое дисперсия, что она показывает? А теперь подумайте «чисто алгебраически», ведь для несмещенных оценок  $\text{MSE} = D\hat{\theta}$ . **MSE** – это ошибка (не важно какая), а мы знаем, что чем меньше ошибка, тем лучше. Значит, вроде как, чем она меньше, тем лучше.

В несмещенности выборочного среднего  $\bar{X}$  мы уже убедились. Несмещенность выборочной медианы можно объяснить, например, так. Пусть  $X = (X_1, X_2, \dots, X_n)$  – выборка из нормального распределения  $N_{a, \sigma^2}$ . Тогда  $Y = (Y_1, Y_2, \dots, Y_n)$ , где  $Y_i = X_i - a$  – это выборка из нормального распределения  $N_{0, \sigma^2}$ . При этом, согласно свойствам линейных преобразований,  $-Y \sim N_{0, \sigma^2}$ . Но медианы  $Y$  и  $-Y$  отличаются знаком, то есть  $\text{med } Y = -\text{med}(-Y)$ , а

значит ее среднее  $E(\text{med } Y)$ , так как его распределение не поменялось, должно быть равно нулю. Ну а тогда  $E(\text{med } X) = E(a + \text{med } Y) = a$ , и медиана оказывается несмещенной оценкой математического ожидания.

Дисперсию выборочного среднего мы уже вычисляли, она равна  $\frac{\sigma^2}{n}$ . Несколько сложнее вычислить дисперсию выборочной медианы, опустим этот момент, и приведем результат:

$$D(\widehat{\text{med } \xi}) = \frac{\pi}{2} \cdot \frac{\sigma^2}{n}.$$

В итоге, выборочная медиана оказывается хуже в смысле среднеквадратической ошибки **MSE**. И что, занавес? Эх, было бы все так просто. С точки зрения теории – да, но реальность вносит свои коррективы?

Прикладная статистика имеет дело с реальными данными, а данные, конечно же, частенько содержат ошибки и выбросы, попавшие в выборку случайно (как крокодил, случайно затесавшийся среди черепах). А что, это проблема. Представьте, что вы измеряете средний размер черепашек (ну пусть в условных единицах), а в выборку затесался крокодил. Он неплохо так увеличит числитель нашей дроби у выборочного среднего, смотрите. Например, пусть

$$X = (5, 7, 3, 8, 4, 2, 100).$$

В конце – крокодил, а в начале – аквариумные маленькие черепашки. Смотрите, без крокодила выборка бы была объема 6 и выборочное среднее бы давало

$$\frac{5 + 7 + 3 + 8 + 4 + 2}{6} = \frac{29}{6} \approx 4.8.$$

А что с крокодилом? Мы получим

$$\frac{5 + 7 + 3 + 8 + 4 + 2 + 100}{7} \approx 18.4.$$

Как вам черепашки? Это ж пол аквариума, наверное. Конечно, если выброс один, а выборка велика, то его вклад постепенно нивелируется. А что если выбросы появляются с периодичностью? Прибор там барахлит, или крокодилы хорошо притворяются? Тогда вся наука будет долго недоумевать нашим измерениям.

Тут и всплывает понятие устойчивости оценки. Устойчивая (или робастная) оценка — оценка, которая почти не меняется от добавления в данные выброса. И с точки зрения устойчивости выборочная медиана лучше! Ну смотрите, что, в нашем примере выборочная медиана разве что-то бы почувствовала? Да нет, конечно. Без крокодила вариационный ряд бы был

$$2 < 3 < 4 < 5 < 7 < 8$$



и  $\widehat{\text{med}} \xi = \frac{4+5}{2} = 4.5$ . А когда крокодил добавился, вариационный ряд стал

$$2 < 3 < 4 < 5 < 7 < 8 < 100$$

и  $\widehat{\text{med}} \xi = \frac{4+5}{2} = 5$ . Видите, изменение совсем маленькое, по сравнению с выборочным средним.

Вообще, можно сказать следующее: если какая-то оценка выражается через порядковые статистики (или ранги), то она будет иметь повышенную устойчивость.

## 2.7 А существует ли наилучшая оценка?

Хорошо, допустим, мы выбрали из нескольких несмещенных оценок ту, у которой самая маленькая дисперсия. Будет ли она наилучшей в смысле среднеквадратического подхода? Вдруг есть оценка, дисперсия которой еще меньше, просто мы эту оценку еще не нашли? В случае параметрической модели ответ возможен, и его можно получить, используя так называемое неравенство Рао-Крамера. Неравенство Рао-Крамера применимо в случае, когда распределение удовлетворяет так называемым условиям регулярности (о них мы будем говорить в соответствующей лекции) и говорит, что в классе несмещенных оценок выполняется соотношение:

$$D\hat{\theta} \geq \frac{1}{nI(\theta)},$$

где  $I(\theta)$  – так называемая информация Фишера для одного наблюдения (так можно поступать, так как мы имеем дело с повторной независимой выборкой), вычисляемая по распределению генеральной совокупности. В итоге, если нам удастся найти ту оценку, для которой в неравенстве Рао-Крамера возникает равенство, то, в случае регулярного распределения, мы можем быть уверены, что ничего лучше (с точки зрения теории) получиться просто не может. Все остальные оценки будут как минимум не лучше, чем данная. Например, для нормального распределения  $N_{a,\sigma^2}$  при неизвестном  $\theta = a$  и известной дисперсии  $\sigma^2$  информация Фишера равна

$$I(\theta) = \frac{1}{\sigma^2},$$

а тогда в неравенстве Рао-Крамера оценка снизу на дисперсию оказывается такой:

$$D\hat{\theta} \geq \frac{1}{nI(\theta)} = \frac{\sigma^2}{n},$$

откуда следует, что полученная нами оценка  $\bar{X}$  оказывается наилучшей среди несмещенных в смысле среднеквадратического подхода, или, как еще говорят, эффективной, так как, как мы уже неоднократно повторяли, и даже

вычислили,

$$D\bar{X} = \frac{\sigma^2}{n}.$$

## 2.8 А как вообще получать оценки?

Отлично. Мы уже обсудили с вами то, что такое оценка и зачем она нужна, как сравнивать оценки между собой и даже как определить, является ли оценка наилучшей в некотором классе. Но до сих пор за кадром остался такой важный вопрос: а как вообще получать оценки? И, конечно, как получать «хорошие» оценки (теперь-то мы знаем, что такое хорошесть)?

Оказывается, в параметрической модели есть метод, который позволяет получить самую лучшую оценку (конечно, при объеме выборки, стремящемся к бесконечности). Метод этот называется методом максимального правдоподобия. При выполнении условий регулярности, оценки метода максимального правдоподобия оказываются асимптотически несмещенными, то есть

$$\lim_{n \rightarrow +\infty} E\hat{\theta} = \theta,$$

асимптотически эффективными (то есть неравенство Рао-Крамера в пределе становится равенством):

$$\lim_{n \rightarrow +\infty} (D\hat{\theta} \cdot nI(\theta)) = 1$$

и асимптотически нормальными (новое для нас понятие), которое означает, что

$$Y_n = \frac{\hat{\theta} - \theta}{\sqrt{D\hat{\theta}}} \xrightarrow[n \rightarrow +\infty]{d} Y \sim N_{0,1}.$$

Иными словами, так называемая центрированная  $(\hat{\theta} - \theta)$  и нормированная случайная величина сходится по распределению к случайной величине, имеющей стандартное нормальное распределение.

## 3 Интервальное оценивание

### 3.1 Немного об интервальном оценивании

То, что мы с вами изучали так подробно – это точечное оценивание. Мы оценивали какую-то числовую характеристику, или числовой параметр распределения, числом. Хорошо ли это? Смотря как посмотреть. Дело в том, что мы, конечно, не знаем истинного значения параметра (или характеристики), иначе зачем его оценивать, а потому возникает вот какой вопрос: число 5

– это хорошая оценка, скажем, математического ожидания, или нет? Теоретически выборочное среднее – это хорошая оценка, но мы-то находимся в рамках конкретного эксперимента!

Все зависит от того, насколько эта оценка реально близка к истинному математическому ожиданию. Точечное оценивание, в свою очередь, не дает нам об этом никакой информации: у нас нет информации о погрешности  $|\hat{\theta} - \theta|$ . Да, мы знаем, что если оценка хорошая, то она, конечно, приближается к истинному значению параметра (вспомним **MSE**), но ведь это приближение идет в среднем, с ростом  $n$ , точнее на бесконечности. А мы на практике, конечно, не можем рассматривать выборки бесконечного размера, можем только конечного, и тут же возникает проблема: а какое значение  $n$  достаточно, при каком значении  $n$  оценка «близка» к истинному значению параметра?

Тут вступает в бой так называемое интервальное, или доверительное оценивание. Оно дает возможность написать интервал, в который, с заданной наперед вероятностью, попадет значение истинного параметра. Не стоит наивно ожидать, что значение вероятности берется равным единице, и дело в шляпе: в этом случае наш доверительный интервал будет очень большим, а пользы не будет никакой (подробнее об этом мы поговорим в соответствующей лекции).

Зададим большую вероятность, например, 0.95, называемую уровнем доверия (доверительным уровнем, уровнем надежности).

**Определение 3.1.1** *Доверительный интервал, построенный по выборке  $X = (X_1, X_2, \dots, X_n)$  уровня доверия 0.95 – это интервал с концами  $\theta^\pm(X_1, X_2, \dots, X_n)$ , для которого*

$$P(\theta^- < \theta < \theta^+) \geq 0.95.$$

Концы доверительного интервала  $\theta^\pm$  – это случайные величины, поэтому они случайны.

Пусть опять выборка берется из нормального распределения  $N_{a, \sigma^2}$  с неизвестным параметром  $\theta = a$  и известной дисперсией  $\sigma^2$ . Тогда, как мы знаем, выборочное среднее  $\bar{X}$  имеет распределение  $N_{a, \sigma^2/n}$  и легко получить (и это мы обязательно получим), что

$$P\left(\bar{X} - c_{0.95} \frac{\sigma}{\sqrt{n}} < \theta < \bar{X} + c_{0.95} \frac{\sigma}{\sqrt{n}}\right) = 0.95,$$

где число  $c_{0.95} \approx 1.96$  (что это за числа мы тоже детально обсудим в дальнейшем). Что видно из этого интервала? Во-первых, его середина – это выборочное среднее. Во-вторых, с ростом объема выборки  $n$ , длина интервала уменьшается и стремится к нулю со скоростью порядка  $n^{-1/2}$ . Так, например, если  $n$  увеличить в 100 раз, то длина интервала уменьшится в 10 раз.

## 4 Проверка гипотез

### 4.1 А что такое гипотезы?

Кроме оценивания характеристик распределения, часто бывает важно проверить какую-то гипотезу, то есть какое-то предположение об этом распределении.

Например, мы измеряем объем продаж какого-то товара и хотим проверить гипотезу, что продажи в этом году, в среднем, больше, чем в прошлом году. Просто сравнивать средние больше-меньше нельзя, так как могло случайно получиться больше. Гипотеза же, в свою очередь, не про выборочные средние, а про математические ожидания неизвестного генерального распределения, которые нам неизвестны.

Можно привести и другой пример. Есть монетка, с помощью которой решается «кто сегодня готовит ужин»? Нам очень важно, чтобы монетка давала справедливое решение, т.е. с вероятностью 0.5 выпадала решка и с вероятностью 0.5 – орел. Как это проверить? Да просто бросить монетку много раз и посчитать, сколько раз что выпало. Вообще говоря, мы получим выборку из распределения Бернулли с вероятностью успеха  $p$ . Гипотеза, которая нас интересует, звучит как  $p = 0.5$ .

Итак, гипотеза – это некоторое утверждение про распределение генеральной совокупности  $\xi$ .

### 4.2 Проверка гипотез

Мы будем обычно проверять гипотезу  $H_0$ , которая, как правило, говорит, что какой-то эффект не наблюдается. Например, что продажи не изменились, или монетка правильная.

Если все случайно, невозможно делать безошибочные выводы. Именно поэтому выбирается небольшая вероятность  $\alpha$ , с которой мы согласны ошибочно сказать, что верная нулевая гипотеза не верна. Она называется уровнем значимости (или так называемой вероятностью ошибки первого рода).

Критерий  $\delta$  – это правило, зависящее от выборки и  $\alpha$ , по которому мы говорим, что нулевая гипотеза верна или не верна. Обычно содержательные выводы делаются в случае, если нулевая гипотеза отвергается. Как правило, в нулевой гипотезе содержится утверждение, что некоторого эффекта нет. Соответственно, если гипотеза отвергается, то делается вывод, что эффект есть. Но можем ли мы делать такой вывод?

Кроме нулевой гипотезы  $H_0$  есть еще альтернативная гипотеза  $H_1$ , которая конкретизирует, что мы хотим понимать под фразой: «нулевая гипотеза неверна». Например, если нулевая гипотеза говорит о том, что новое лекарство лучше старого на 5%, то альтернативная гипотеза может утверждать,

что лекарство либо хуже, либо лучше, но не больше, чем на 4%. Вариант отклонения от нулевой гипотезы на, например, 0.1%, нас не волнует.

Если выбрана альтернативная гипотеза, то возникает так называемая ошибка второго рода, которая означает, что верна альтернативная гипотеза, но мы не отвергли нулевую гипотезу. Вероятность ошибки второго рода может быть большой (мы ее не выбираем, в отличие от вероятности ошибки 1 рода), поэтому мы не говорим фразу: «гипотезы принимаем», а говорим - не отвергаем.

Мощность критерия, равная единице минус вероятность ошибки второго рода, то есть вероятность отвергнуть гипотезу  $H_0$ , если верна  $H_1$ , характеризует умение критерия обнаруживать отличие  $H_1$  от нулевой гипотезы. Она обладает следующими свойствами: во-первых, чем больше объем выборки  $n$ , тем мощность больше; во-вторых, чем больше уровень значимости  $\alpha$ , тем мощность больше; в третьих, чем дальше альтернативная гипотеза от основной гипотезы, тем мощность больше. В частности, отсюда следует, что чтобы увеличить мощность, выбирается уровень значимости максимальным из допустимых. Понятие допустимости зависит от того, каковы риски в случае принятия неверного решения. Например, если мы потеряем 100 рублей в случае ошибки 1 рода, мы можем взять  $\alpha = 0.1$ , а если 1000000 рублей, то  $\alpha$  должно быть существенно меньше.

Так как уровень значимости  $\alpha$  определяется не внутри процедуры проверки гипотезы, а задается исследователем, то часто результатом критерия является пороговое значение  $p$ , которое называется  $p$ -значением. Если уровень значимости  $\alpha$  больше  $p$ , то гипотеза  $H_0$  отвергается, а если меньше – нет оснований отвергать гипотезу. Поэтому  $p$ -значение – это минимальный уровень значимости, при котором гипотеза отвергается (и максимальный, при котором не отвергается). Если  $p$  большое, например, 0.7, то при всех разумных уровнях значимости нулевая гипотеза не отвергается. А если, например,  $p = 0.0001$ , то гипотеза будет отвергаться. Для промежуточных  $p$ , например, в районе 0.01 – 0.1, результат проверки гипотезы зависит от конкретного уровня значимости.

Обычно мы не можем одновременно уменьшить  $\alpha$  и увеличить мощность (то есть уменьшить ошибку второго рода). Давайте разберем такой шуточный пример. Пусть у нас есть выборка объема 1 из нормального распределения  $N_{a,1}$ . Рассмотрим две гипотезы:  $H_0 : a = 0$  и  $H_1 : a = 1$  и критерий

$$\delta(X_1) = \begin{cases} H_0, & X_1 \leq d \\ H_1, & X_1 > d \end{cases}.$$

Что такое ошибка первого рода? Это вероятность не отвергнуть гипотезу  $H_1$  в случае, когда верна гипотеза  $H_0$ , то есть  $\alpha = P_{H_0}(\delta = H_1) = P_{H_0}(X_1 > d)$ .

Что такое ошибка второго рода? Это вероятность не отвергнуть гипотезу  $H_0$  в случае, когда верна гипотеза  $H_1$ , то есть  $P_{H_1}(\delta = H_0) = P_{H_1}(X_1 \leq d)$ . На рисунке четко видно, что уменьшая вероятность ошибки первого рода ( $\alpha$ ), автоматически возрастает вероятность ошибки второго рода, тем самым уменьшается мощность критерия, и наоборот. Давайте покажем то, что кри-

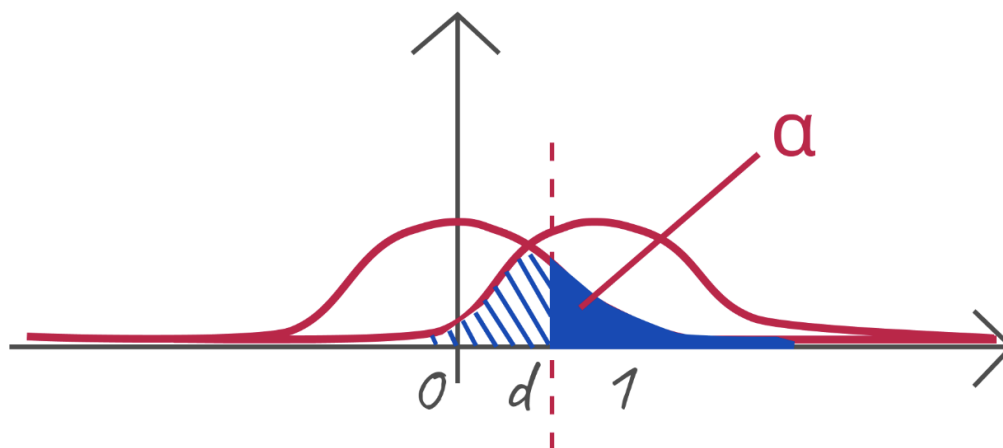


Рис. 4: Ошибки первого и второго рода

терии по-разному регулируют ошибки первого и второго рода. Напомним, что ошибку первого рода мы фиксируем, а ошибка второго рода – как получится, так получится. Тем самым, может быть не все равно, что рассмотреть в качестве нулевой гипотезы, а что в качестве альтернативы.

Например, мы собираемся лететь на самолете и есть две гипотезы – самолет исправный и самолет неисправный. При проверки гипотезы возникают две ошибки – отвергнуть гипотезу, что самолет исправный, хотя самолет исправный (соответственно, зря сдать билет и зря потерять деньги) и отвергнуть гипотезу, что самолет неисправный, хотя самолет неисправный (соответственно, полететь на самолете и попасть в катастрофу).

В каком случае ошибку важнее контролировать? Конечно, во втором. Поэтому, в идеале, надо, чтобы нулевая гипотезы была о том, что самолет неисправен. В этом случае мы полетим на неисправной самолете с выбранной нами (и, конечно, очень маленькой) вероятностью. Ну а если не сможем контролировать ошибку второго рода, то просто есть шанс, что потеряем в деньгах. Так часто бывает, что по смыслу ошибок они друг от друга отличаются по критичности последствий. Соответственно, лучше, если более критичное последствие произойдет в случае ошибки первого рода, выбираемой нами самими.

### 4.3 Сравнение критериев

Как ранее оценок для одного параметра могло быть много, так и критериев для проверки одной и той же гипотезы может быть много. Как их

сравнивать?

Сначала нужно зафиксировать альтернативу так, чтобы она включала туда те отличия от нулевой гипотезы, которые нам важно обнаружить.

1. Критерии можно сравнивать по мощности против выбранной альтернативы, т.е. вероятности отвергнуть гипотезу, если верна альтернатива.
2. Также, как и оценки, важным свойством критериев является их устойчивость к выбросам. Поэтому критерии, которые используют не сами значения выборки, а только их ранги, являются более устойчивыми.

Например, пусть  $H_0 : E\xi = 0$ , где  $\xi$  — то, на сколько возросли продажи товара. Если у нас учет хороший, то мы можем использовать так называемый  $t$ -тест. А если плохой и предполагаем, что там могут быть ошибки, то нужен более устойчивый тест. Подробнее об этих вещах мы будем говорить в соответствующих лекциях.