

Проверка статистических гипотез.  
Р-значения. Доверительные интервалы

Теория вероятностей  
и математическая статистика / Урок 5



# Проверка статистических гипотез

Статистическая гипотеза — предположение о виде распределения и свойствах случайной величины, которое можно подтвердить или опровергнуть на основании имеющихся данных.

Например,

- 1 Гипотеза: математическое ожидание случайной величины равно 10.
- 2 Гипотеза: случайная величина имеет нормальное распределение.
- 3 Гипотеза: две случайные величины имеют одинаковое математическое ожидание.

Проверяемую гипотезу принято называть **нулевой** и обозначать  $H_0$ .

*Пример.* Имеется станок, изготавливающий шарики для подшипников, который настроен делать шарики с диаметром 1 мм. На основании выборки из значений диаметров таких шариков мы можем проверить, правильно ли станок откалиброван (т.е. делает ли он такие шарики, которые он настроен делать).

В таком случае в качестве нулевой гипотезы  $H_0$  берётся гипотеза о том, что математическое ожидание диаметра шарика равно 1 мм.

Проверяемую гипотезу принято называть **нулевой** и обозначать  $H_0$ .

*Пример.* Имеется станок, изготавливающий шарики для подшипников, который настроен делать шарики с диаметром 1 мм. На основании выборки из значений диаметров таких шариков мы можем проверить, правильно ли станок откалиброван (т.е. делает ли он такие шарики, которые он настроен делать).

В таком случае в качестве нулевой гипотезы  $H_0$  берётся гипотеза о том, что математическое ожидание диаметра шарика равно 1 мм.

Параллельно с нулевой гипотезой рассматривается противоречащая ей гипотеза  $H_1$ , называемая **альтернативной** или **конкурирующей**.

В нашем примере в качестве альтернативной гипотезы будет выступать гипотеза о том, что математическое ожидание диаметра шарика не равно 1 мм.

В зависимости от задачи альтернативные гипотезы бывают левосторонние, правосторонние или двухсторонние.

Например, в примере выше альтернативная гипотеза двухсторонняя, поскольку в соответствии с ней математическое ожидание может быть как больше 1, так и меньше.

Односторонние гипотезы, возникают, например, когда анализируемое значение равно 0, если нулевая гипотеза верна, и больше 0, если гипотеза неверна. Например, такие гипотезы проверяются в F-тестах (подробнее на уроках 7 и 8).

В любом случае альтернативная гипотеза представляет собой дополнение к нулевой (т.е. хотя бы одна из них всегда верна).

Проверяя нулевую гипотезу, мы по выборке считаем некоторое значение (зависит от вида гипотезы) и сравниваем его с теоретическим.

При проверке данной гипотезы наша задача установить, какое отклонение полученного значения от теоретического мы можем принять как допустимое (т. е. **незначимое**), а какое отклонение уже нельзя списать на случайность (**значимое**).

Если отклонение значимо, то гипотеза  $H_0$  отвергается в пользу альтернативной. Если же отклонение не является значимым, то гипотеза  $H_0$  остаётся в силе.

- 1 Формулируются нулевая и альтернативная гипотезы.
- 2 Задаётся некоторая статистика (функция от выборки)  $S(X)$ , которая в условиях справедливости нулевой гипотезы  $H_0$  имеет известное распределение (в частности, известна её функция распределения  $F_S(x) = P(S < x)$ ).
- 3 Фиксируется уровень значимости  $\alpha$  — допустимая для данной задачи вероятность ошибки первого рода (чаще всего 0.01, 0.05 или 0.1).
- 4 Определяется критическая область  $\Omega_\alpha$ , такая, что  $P(S \in \Omega_\alpha | H_0) = \alpha$ .
- 5 Проводится статистический тест: для конкретной выборки  $X$  считается значение  $S(X)$ , и если оно принадлежит  $\Omega_\alpha$ , то заключаем, что данные противоречат гипотезе  $H_0$ , и принимается гипотеза  $H_1$ .



Выбор статистики  $S$  зависит от того, какого рода гипотеза проверяется, какое распределение имеется и что нам известно.

Например, если проверяется гипотеза относительно математического ожидания нормально распределённой случайной величины с *известной* дисперсией, то используется **Z-статистика**:

$$Z = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}},$$

где  $X$  — выборка,  $\bar{X}$  — выборочное среднее,  $\mu$  — утверждаемое гипотезой  $H_0$  значение математического ожидания,  $\sigma$  — известное среднее квадратическое отклонение,  $n$  — число элементов в выборке.

В предположении верности гипотезы  $H_0$  статистика  $Z$  имеет *стандартное нормальное распределение* (т.е. нормальное распределение с параметрами  $\mu = 0$ ,  $\sigma = 1$ ).

Если же дисперсия неизвестна, используется t-статистика:

$$t = \frac{\bar{X} - \mu}{\sigma_X / \sqrt{n}},$$

где  $\sigma_X$  — несмещённая оценка среднего квадратического отклонения.

В предположении верности гипотезы  $H_0$  t-статистика имеет распределение Стьюдента с параметром  $df = n - 1$ .

Если же дисперсия неизвестна, используется **t-статистика**:

$$t = \frac{\bar{X} - \mu}{\sigma_X / \sqrt{n}},$$

где  $\sigma_X$  — несмещённая оценка среднего квадратического отклонения.

В предположении верности гипотезы  $H_0$  t-статистика имеет **распределение Стьюдента** с параметром  $df = n - 1$ .

*Замечание.* Если известна дисперсия, то известно и среднее квадратическое отклонение. Наоборот, если известно среднее квадратическое отклонение, то известна и дисперсия.

Описанные выше принципы дают точный результат только для нормального распределения. Однако, мы знаем, что, согласно Центральной предельной теореме, при достаточно больших  $n$  распределение выборочного среднего  $\bar{X}$  близко к нормальному, поэтому в этих случаях также можно использовать t-статистику и распределение Стьюдента.

Однако, если результат теста оказывается пограничным, нужно быть аккуратным при его интерпретации и помнить о том, что ЦПТ работает лишь приблизительно.

Таков подход к проверке гипотез о математическом ожидании. Для проверки других гипотез, как правило, используются другие распределения и другие статистики, а также накладываются свои ограничения.

Например, для проверки гипотез о дисперсии нормального распределения используют *распределение хи-квадрат*. В дальнейшем мы познакомимся и с другими тестами, которые опираются на *распределение Фишера*, *распределение Колмогорова* и др.

Ошибки первого и второго рода возникают в задачах, в которых требуется определить, произошло какое-то событие или нет.

Ошибка первого рода (т.н. **false positive**) соответствует ситуации, когда было определено, что событие произошло, тогда как реально оно не произошло.

Ошибка второго рода (т.н. **false negative**) — обратная ситуация: мы определили, что событие не произошло, а реально оно произошло.

Например, на входе в каждый аэропорт стоит рамка-металлодетектор. Её задача — выявлять людей, которые пытаются пронести в аэропорт что-то опасное.

В этом случае:

- ошибка первого рода: рамка сработала на человеке, который ничего опасного не несёт,
- ошибка второго рода: рамка не сработала на человеке, который несёт что-то опасное.

Например, на входе в каждый аэропорт стоит рамка-металлодетектор. Её задача — выявлять людей, которые пытаются пронести в аэропорт что-то опасное.

В этом случае:

- ошибка первого рода: рамка сработала на человеке, который ничего опасного не несёт,
- ошибка второго рода: рамка не сработала на человеке, который несёт что-то опасное.

Как правило, между вероятностями ошибок первого и второго рода приходится балансировать, т.е. уменьшение одной вероятности приводит к увеличению другой.

Например, в примере с рамкой вероятность ошибки первого рода будет невероятно высокой, поскольку ошибки второго рода в этом случае недопустимы.



В контексте проверки статистических гипотез под **событием** мы понимаем то, что **нулевая гипотеза была отвергнута**.

**Уровень значимости** — это вероятность ошибки первого рода, т.е. вероятность отвергнуть верную нулевую гипотезу.

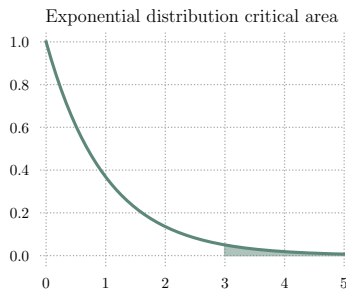
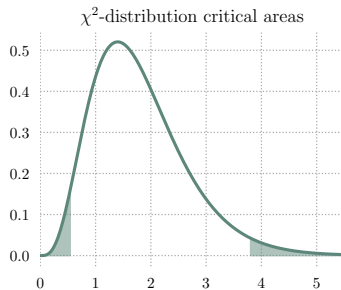
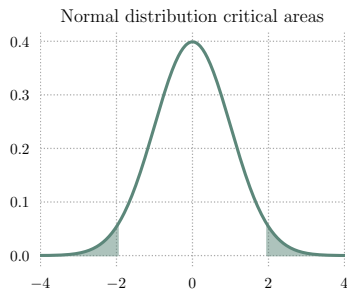
Низкие значения уровня значимости  $\alpha$  делают проводимый тест **более осторожным**: вероятность ошибки первого рода мала, значит, если по результатам теста мы отвергаем нулевую гипотезу, то мы можем быть более уверены в том, что она и правда неверна.

Однако, вместе с тем повышается и вероятность не отвергнуть ложную нулевую гипотезу. Т.е. если нулевая гипотеза не была отвергнута, то далеко не факт, что она и впрямь верна.

При проведении статистического теста мы строим критическую область для статистики  $S(X)$ , построенной ранее. Мы можем это сделать, потому что знаем распределение этой статистики, в частности, её функцию распределения  $F_S(x)$ .

Кроме того, к этому моменту мы уже зафиксировали уровень значимости  $\alpha$ . Это значение является вероятностью попасть в критическую область.

Как правило, используемые в статистике распределения имеют функцию плотности в форме «колокола»: ярко выраженный пик в центре и хвосты по краям.



Некоторые распределения имеют два хвоста, некоторые — один. Кроме того, хвосты могут быть как конечными, так и уходить в бесконечность.

Критические области представляют собой как раз вот эти «хвосты» распределения. Число хвостов критической области определяется числом «сторон» у альтернативной гипотезы.

Например, если альтернативная гипотеза заключается в том, что мат. ожидание *не равно* какому-то числу, то в соответствии с этой гипотезой оно может быть как больше, так и меньше, т.е. альтернативная гипотеза *двухсторонняя*. В этом случае и критическая область будет двухсторонней, т.е. будет иметь два «хвоста».

Как правило, критические области строят следующим образом:

- Левосторонняя область:  $\Omega_\alpha = (-\infty, t_\alpha)$ .
- Правосторонняя область:  $\Omega_\alpha = (t_{1-\alpha}, \infty)$ .
- Двусторонняя область:  $\Omega_\alpha = (-\infty, t_{\alpha/2}) \cup (t_{1-\alpha/2}, \infty)$ .

Здесь  $t_x$  обозначает квантиль порядка  $x$ , т.е.  $F_S(t_x) = x$ .

Итак, осталось лишь провести статистический тест. У нас есть:

- 1 зафиксированная нами статистика  $S(X)$ ,
- 2 построенная нами критическая область  $\Omega_\alpha$ .

Считаем значение статистики от нашей выборки. Если это значение попадает в критическую область, то заключаем, что данные противоречат нулевой гипотезе, и её следует отвергнуть.

Если значение в критическую область не попало, то заключаем, что данные нулевой гипотезе не противоречат.

Итак, осталось лишь провести статистический тест. У нас есть:

- 1 зафиксированная нами статистика  $S(X)$ ,
- 2 построенная нами критическая область  $\Omega_\alpha$ .

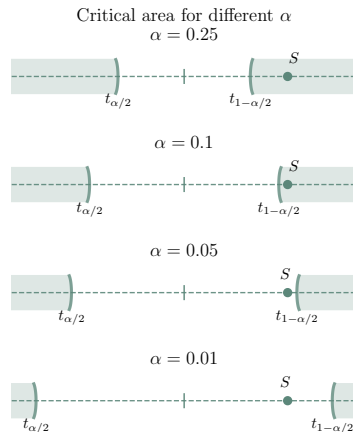
Считаем значение статистики от нашей выборки. Если это значение попадает в критическую область, то заключаем, что данные противоречат нулевой гипотезе, и её следует отвергнуть.

Если значение в критическую область не попало, то заключаем, что данные нулевой гипотезе не противоречат.

*Замечание.* То, что данные не противоречат нулевой гипотезе, не означает, что она верна.

P-значения позволяют получить результат проверки статистических гипотез сразу для многих уровней значимости.

Как мы теперь знаем, уменьшение уровня значимости приводит к расширению границ критической области: чем меньше уровень значимости, тем сложнее попасть в критическую область и, как следствие, отвергнуть нулевую гипотезу. P-значение представляет собой наибольшее значение уровня значимости  $\alpha$ , при котором гипотезу можно принять, т.е. при котором значение статистики, посчитанной по выборке, ещё не попадает в критическую область.





**Р-значение** — это такое значение  $\alpha$ , при котором значение статистики попадает ровно на границу критической области.

Пусть  $F_S(x)$  — функция распределения рассматриваемой статистики, а  $t_x$  — квантиль порядка  $x$  для этого распределения. **Как считать Р-значение:**

- ① Для правосторонней области  $\Omega_\alpha = (t_{1-\alpha}, \infty)$  имеем условие  $t_{1-\alpha} = S$ , откуда

$$P_r = 1 - F_S(S)$$

- ② Для левосторонней области  $\Omega_\alpha = (-\infty, t_\alpha)$ , условие  $t_\alpha = S$ , откуда

$$P_l = F_S(S)$$

- ③ Для двухсторонней области  $\Omega_\alpha = (-\infty, t_{\alpha/2}) \cup (t_{1-\alpha/2}, \infty)$  нужна комбинация двух:

$$P = 2 \cdot \min(P_l, P_r)$$

На практике использовать P-значения можно так: если оно больше выбранного нами уровня значимости  $\alpha$ , то гипотезу можно принять. В противном случае, гипотезу следует отвергнуть.

Отметим, что у статистических тестов есть ещё одна важная характеристика — **статистическая мощность**. Тогда как уровень значимости связан с ошибками первого рода, статистическая мощность связана с ошибками второго рода.

Как именно они связаны, и что можно делать с помощью статистической мощности — читайте в дополнительных материалах к уроку.

## Доверительные интервалы

Ранее мы познакомились со способами оценивать параметры распределения по выборке. Всё это были **точечные** оценки, т.е. мы оценивали параметр каким-то единственным числом.

Минус таких оценок в том, что у нас нет возможности понять, насколько хорошей такая оценка получилась (т.е. насколько близко мы оказались к реальному значению оцениваемого параметра).

Для исправления этого недостатка используют доверительные интервалы.

Доверительный интервал — это интервал, который с некоторой вероятностью (заданной заранее) содержит значение оцениваемого параметра.

Более строго: пусть задано число  $p$ , называемое уровнем доверия или доверительной вероятностью. Доверительным интервалом для параметра  $\theta$  называется пара статистик  $L$  и  $U$ , таких, что

$$P(L \leq \theta \leq U) = p$$

Доверительные интервалы используют ту же математическую базу, что и статистические тесты.

Например, пусть дана выборка  $X$  из нормально распределённой случайной величины с известной дисперсией  $\sigma^2$ , и требуется построить доверительный интервал для математического ожидания  $\mu$  с доверительной вероятностью  $p$ . Мы знаем, что в этом случае статистика

$$Z = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}}$$

имеет стандартное нормальное распределение.

Обозначим  $\alpha = 1 - p$ . Можно убедиться в том, что

$$P(t_{\alpha/2} \leq Z \leq t_{1-\alpha/2}) = p,$$

где  $t_x$  — квантиль порядка  $x$  для стандартного нормального распределения. Подставляя сюда  $Z$ , получаем

$$P\left(t_{\alpha/2} \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq t_{1-\alpha/2}\right) = p$$

$$P\left(t_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \leq \bar{X} - \mu \leq t_{1-\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}\right) = p$$

$$P\left(\bar{X} + t_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + t_{1-\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}\right) = p$$



В случае **неизвестной дисперсии** мы используем статистику

$$t = \frac{\bar{X} - \mu}{\sigma_X / \sqrt{n}},$$

где  $\sigma_X$  — выборочное среднее квадратическое отклонение. Эта статистика имеет распределение Стьюдента, поэтому


$$P(t_{\alpha/2, n-1} \leq t \leq t_{1-\alpha/2, n-1}) = p,$$

где  $t_{x, n-1}$  — квантиль порядка  $x$  для распределения Стьюдента с параметром  $df = n - 1$ . Аналогичным способом получаем доверительный интервал:

$$P\left(\bar{X} + t_{\alpha/2, n-1} \cdot \frac{\sigma_X}{\sqrt{n}} \leq \mu \leq \bar{X} + t_{1-\alpha/2, n-1} \cdot \frac{\sigma_X}{\sqrt{n}}\right) = p$$

Замечание, сделанное ранее про использование  $t$ -статистики для распределений, отличных от нормального, справедливо и здесь: если размер выборки  $n$  достаточно велик, то доверительные интервалы для математического ожидания можно строить по полученным выше формулам.

Однако, как и ранее, отметим, что при использовании доверительных интервалов в таких случаях надо проявлять осторожность, поскольку строятся они лишь приблизительно.



Спасибо за внимание