

Проверка статистических гипотез • Hypothesis testing

- посмотрим, как отличить статистически достоверное событие от случайного
- узнаем как устроен пайплайн проверки гипотез
- реализуем функции расчета для некоторых тестов

- Статистическая гипотеза – предположение о свойствах генеральной совокупности.
- Всю ГС мы исследовать не можем, значит, мы должны собрать **репрезентативную выборку**, изучить ее, а после проверить гипотезу.

Задача с кофе



Задача с кофе

- В ТЦ стоит один из наших автоматов с кофе.
- Ранее:
 - из тех, кто подходил к нему, с кофе уходил каждый второй.
- Есть гипотеза:
 - У нашего автомата сложный интерфейс, и некоторых людей это сбивает с толку и они уходят
- Изменения:
 - Разработан новый тестовый интерфейс, и поставлен на наш автомат.
 - В случае успеха, можно будет выкатывать на остальные наши автоматы
- После :
 - Из 300 людей, которые подошли к нашему автомату, купили 167

- Подтвердилась ли наша гипотеза? К какому результату интуитивно склоняетесь?
- Желательно не ошибиться с выбором, так как внедрение нового интерфейса на все наши автоматы стоит денег.

Формализуем:

- ГС - все люди, которые подошли бы к нашему автомату
- у нас есть **выборка** $x_1, x_2, x_3, \dots, x_{300}$
- $x_i \sim Be(p)$ (купил/не купил)
- p - неизвестный для нас параметр ГС - доля тех, кто купил бы, подойдя к автомату
- Пример выборки $[1, 1, \dots, 1, 0, 1, 0, 1, 1]$ (167 купили, 133 не купили)

Нулевая гипотеза H_0 – это гипотеза, которой мы придерживаемся, пока наблюдения не заставят признать обратное. Ей всегда сопутствует альтернативная гипотеза H_1 .

- H_0 почти всегда формулируется, как "значимых изменения нет"
- H_1 - "значимые изменения есть"

В нашем случае:

- $H_0 : p = 0.5$ (конверсия в покупку такая же и осталась)
- $H_1 : p > 0.5$ (конверсия увеличилась)

По результатам исследования мы остановимся на одной из гипотез

Ошибка первого и второго рода

- Ошибка первого рода(FP) - это ситуация, когда H_0 отвергается, хотя она, на самом деле, верна
 - α – вероятность ошибки первого рода или уровень значимости
- Ошибка второго рода(FN) - это ситуация, когда H_0 принимается, хотя она неверна
 - β – вероятность ошибки второго рода

Некоторая сложность: при уменьшении ошибки первого рода, увеличивается ошибка второго рода и наоборот.

Ошибка первого и второго рода

- Матрица ошибок (confusion matrix)

Hypothesis testing:

Decision

		H_0 true (Fail to reject)	H_0 false (Rejecting H_0)
Actual	H_0 true	<p>TRUE NEGATIVE</p> <p>Correct decision: Confidence level (prob $1 - \alpha$)</p>	<p>FALSE POSITIVE</p> <p>Type I Error: Significance level/Size (α) (prob α)</p>
	H_0 false	<p>FALSE NEGATIVE</p> <p>Type II Error: fail to reject (prob β)</p>	<p>TRUE POSITIVE</p> <p>Correct decision: Power (prob $1 - \beta$)</p>

- Например, суд выдвигает гипотезу H_0 : подсудимый невиновен. А он, на самом деле, виновен, но суд признает его невиновным за отсутствием улик (презумпция невиновности). То есть суд принимает гипотезу, хотя она неверна.
- "Ложноположительный результат" при медицинских анализах – это ошибка какого рода?

Статистика критерия

Статистика – любая функция, получаемая по выборке. В каком-то смысле это просто посчитанная метрика

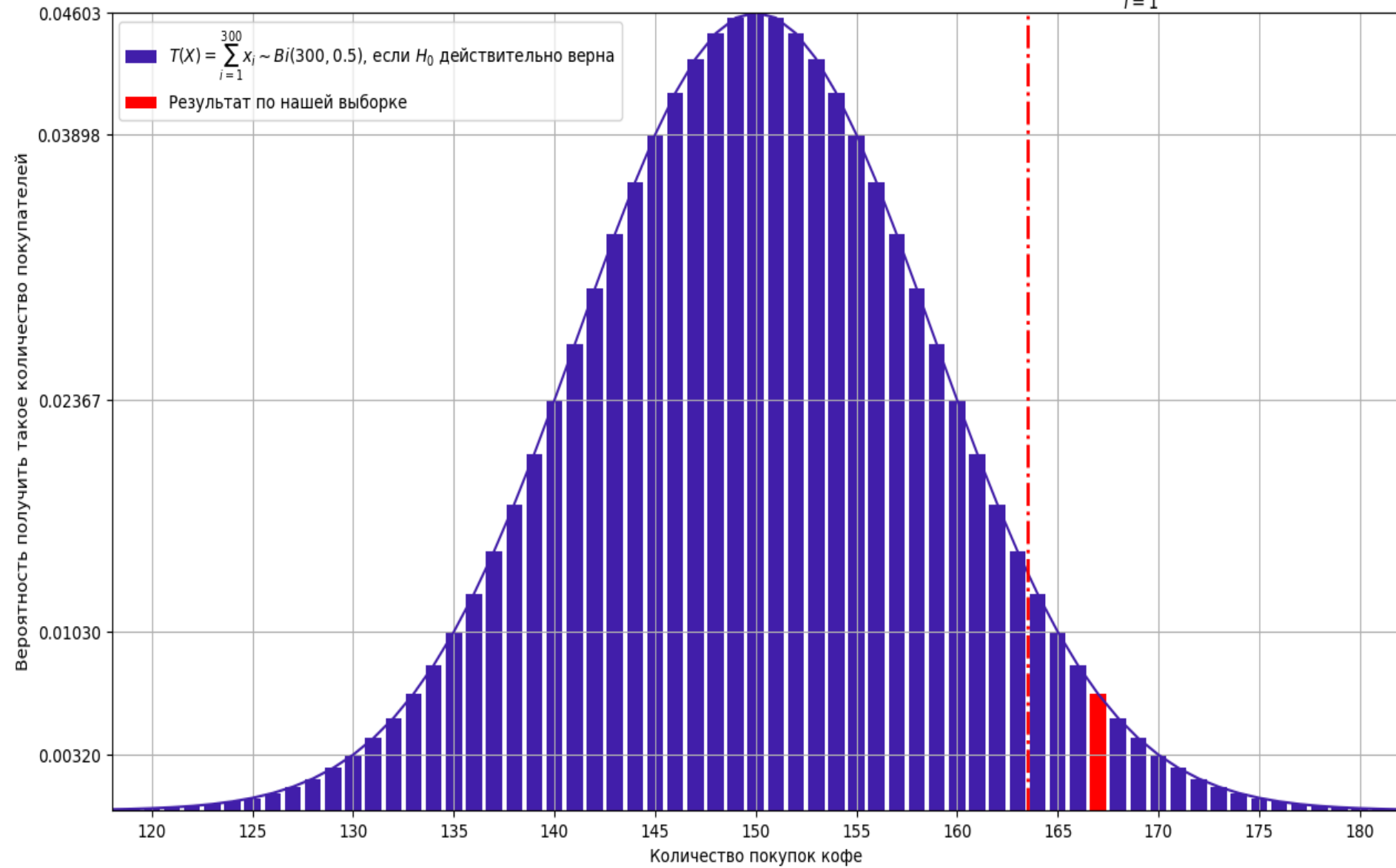
Обозначение: $T(\vec{x})$, где $\vec{x} = (x_1, x_2, x_3, \dots, x_n)$ - выборка

- Статистика агрегирует информацию о выборке.
- Самые частые статистики:
 - Среднее, доля, медиана, количество, квантиль и т.д
- Кастомные статистики
 - Показатель удовлетворенности клиентов (Customer Satisfaction Score, CSS)
 - Показатель устойчивости бизнеса (Customer Loyalty Index, CLI)
 - "Здесь бы могла быть ваша статистика 😊"

- Возьмем в нашем примере с кофе $T(X) = \sum_{i=1}^{300} x_i$, иначе говоря, сколько людей купили у нас кофе.
- Важно понимать, что $T(X)$ тоже является **случайной величиной**, а значит имеет свое распределение, это **ключевой момент** в данной теме.
- Именно знание распределения статистики дает нам понимания, насколько **экстремальное** значение мы вообще получили.
- Например при проверки монетки на честность получить 90 орлов после 100 подбрасываний кажется слишком экстремальным, и скорее она не честная.

Наша задача

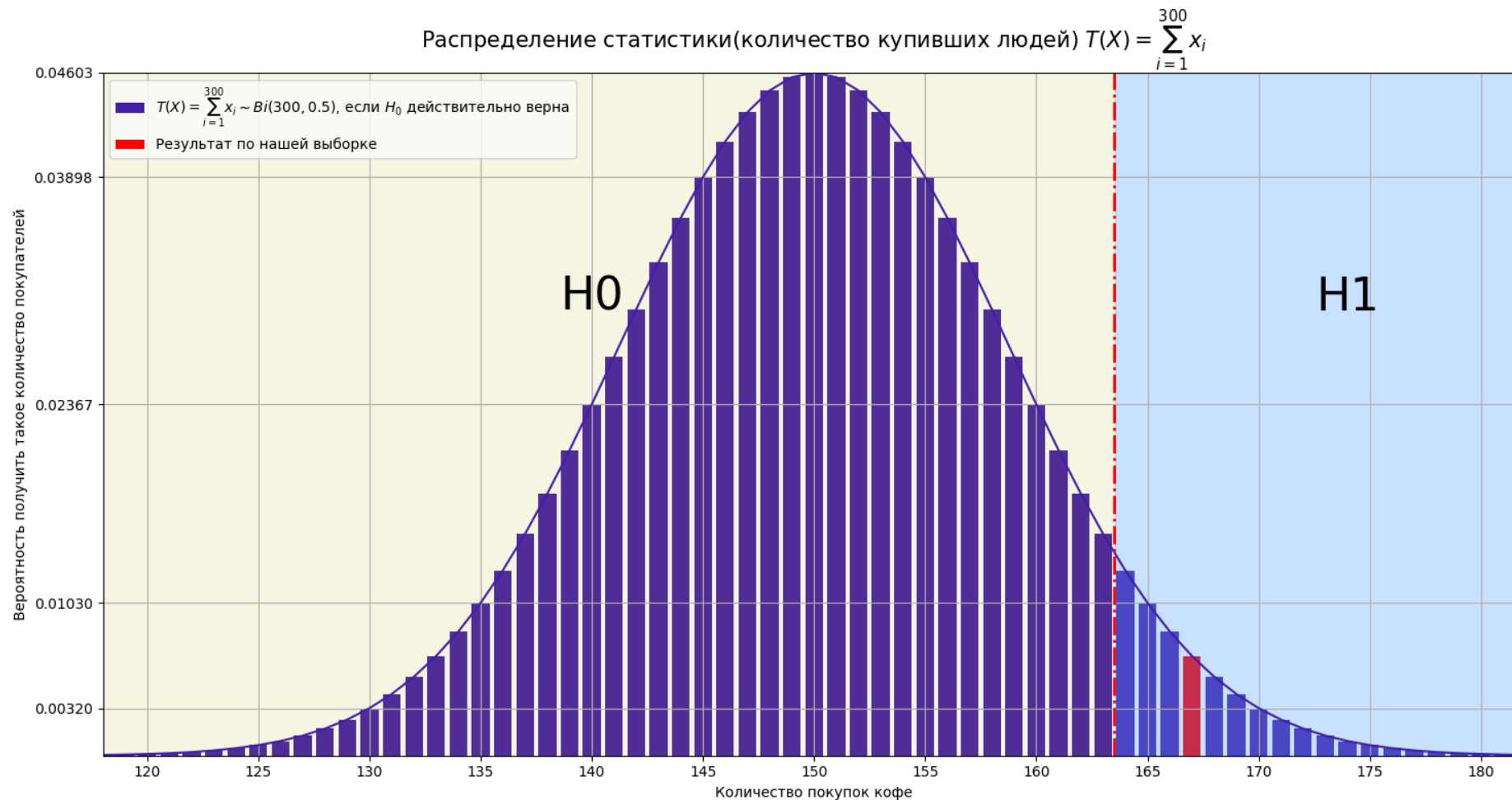
Распределение статистики (количество купивших людей) $T(X) = \sum_{i=1}^{300} x_i$



1. Сформулировать основную и альтернативную гипотезы, задать уровень значимости α
2. Найти критические значения статистики для соответствующего уровня значимости
3. Вычислить значение статистики и определить, попало ли оно в критическую область
4. Сделать вывод: если значение попало в критическую область - отвергнуть нулевую гипотезу, в противном случае принять

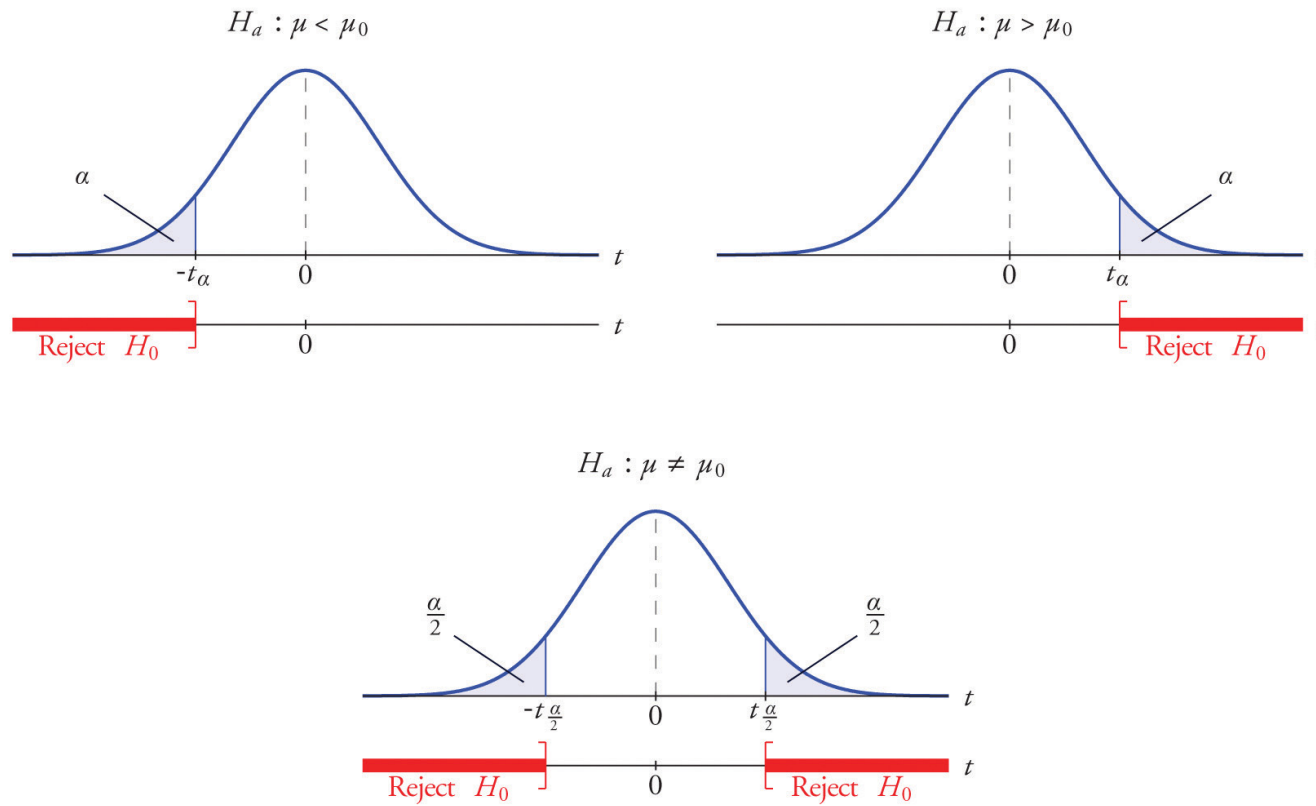
Наша задача

- $\alpha = 0.05$ - уровень значимости или ошибка первого рода



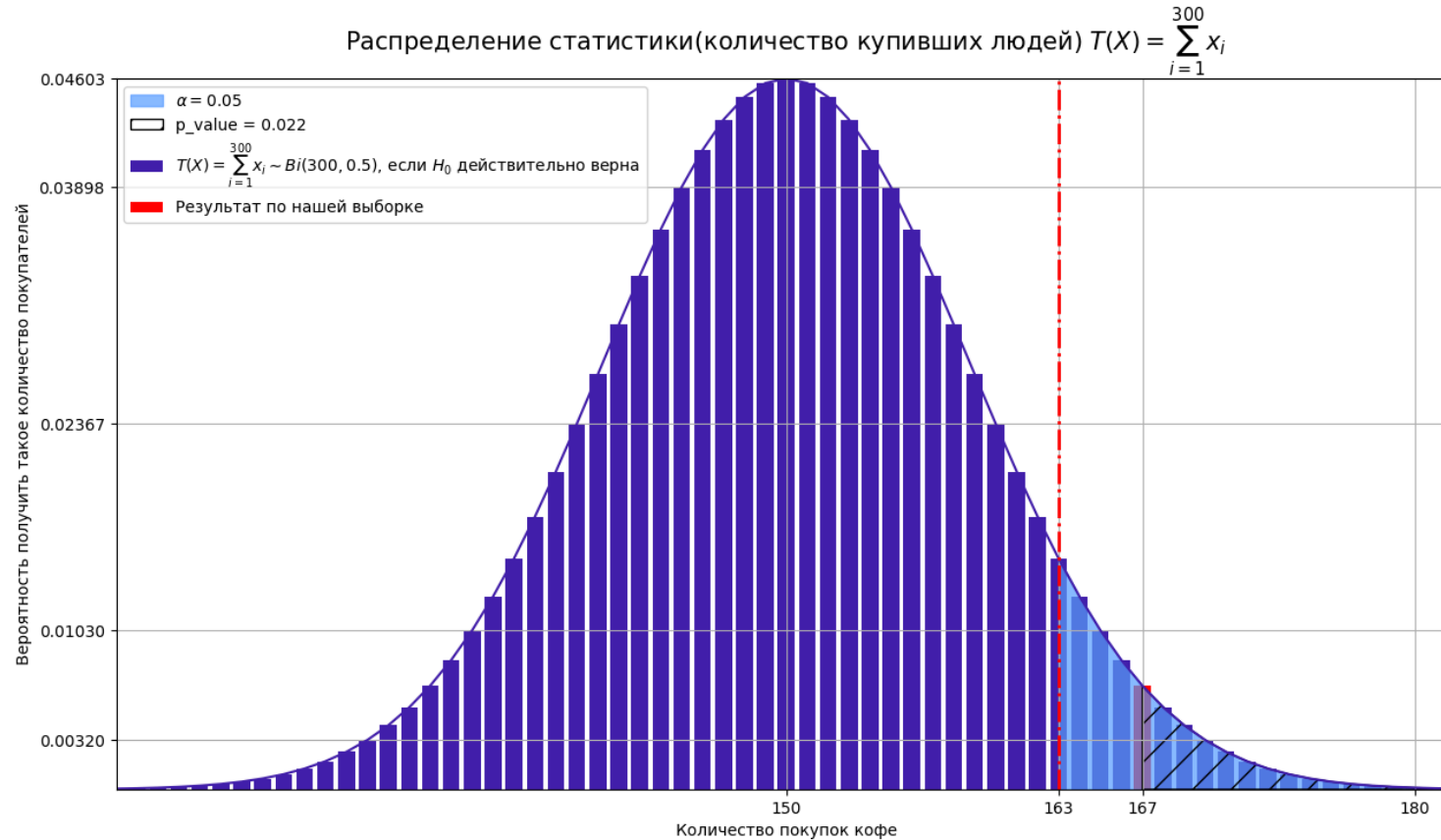
Критическая область

Критической областью называется область значений статистики критерия, при которых отвергается H_0 . А критические значения – это граница критической области.



p-value

- **P-value** можно интерпретировать как вероятность ошибиться, если мы выбираем гипотзу H_1 .



- Статистические тесты позволяют ответить есть ли **статистически значимый результат**
- Ошибки 1-го рода(α) и 2-го рода(β) не хороши для нас, однако в большинстве случаев их не избежать. Катастрофичность каждой из них зависит от конкретной задачи
- **p-value** позволяет оценить вероятность ошибки и принять решение