



# Статистика ФИВТ ПМИ

## Прикладной поток

Лекция 11



# 5. Проверка статистических гипотез

## 5.4. p-value



## Гипотезы (напоминание)

$X = (X_1, \dots, X_n)$  — выборка из неизвестного распределения  $P \in \mathcal{P}$ .

$H_0 : P \in \mathcal{P}_0$  — основная гипотеза;

$H_1 : P \in \mathcal{P}_1$  — альтернативная гипотеза.

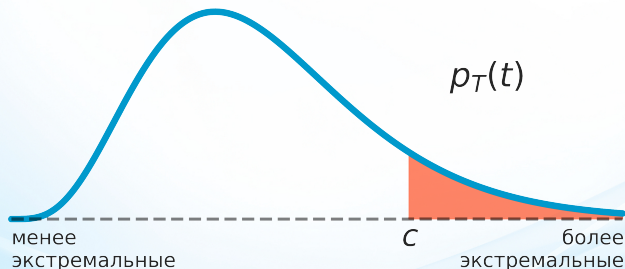
$S \subset \mathcal{X}$  — критерий уровня значимости  $\alpha$  для проверки  $H_0$  vs.  $H_1$ ,  
если  $P(X \in S) \leq \alpha, \forall P \in \mathcal{P}_0$ .

Варианты ответа:

1.  $X \in S \implies H_0$  отвергается  $\implies$  результат стат. значим;
2.  $X \notin S \implies H_0$  **не отвергается**  $\implies$  результат не стат. значим

## Гипотезы (напоминание)

Часто критерий имеет вид  $S = \{T(x) \geq c\}$ ,  
где  $T(X)$  — статистика критерия.



$H_0$  отвергается  $\iff T(X) \geq c_\alpha$

Значение  $t_1$  более экстремально, чем  $t_2$ , если  $t_1 > t_2$ .



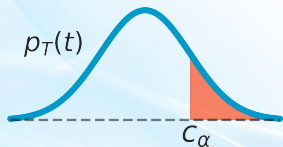
## Критерии (напоминание)

Часто критерий имеет вид  $S = \{T(x) \geq c_\alpha\}$ ,  
где  $T(X)$  — статистика критерия.

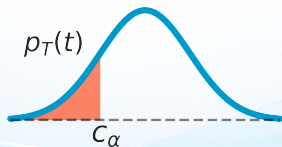
$\alpha$  выбирается **ДО** эксперимента,

$c_\alpha$  вычисляется из условия  $P_0(T(X) > c_\alpha) \leq \alpha$ .

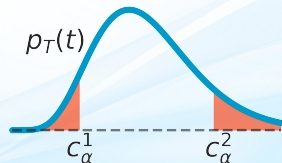
$$S = \{T(x) > c_\alpha\}$$



$$S = \{T(x) < c_\alpha\}$$



$$S = \{|T(x)| > c_\alpha\}$$



*Замечание.* Выбирать  $\alpha$  после эксперимента неправильно.

Так можно подогнать результат под желаемый.

*"Статистика может доказать что угодно, даже истину."*



# Пример (влияние нового препарата на выздоровление)

Испытуемые делятся случайно на две группы:

1. *Исследуемая группа* — принимает новый препарат;

$X = (X_1, \dots, X_n), X_i \sim \text{Bern}(p_1)$  — результаты лечения.

2. *Контрольная группа* — принимает плацебо;

$Y = (Y_1, \dots, Y_m), Y_i \sim \text{Bern}(p_2)$  — результаты лечения.

$H_0: p_1 = p_2$  — отсутствие эффекта

$H_1: p_1 > p_2$  — эффект присутствует



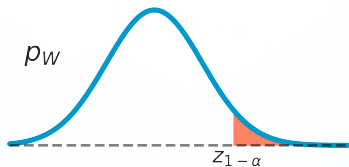
## Пример (влияние нового препарата на выздоровление)

$$\hat{p}_1 = \bar{X} \stackrel{d}{\approx} \mathcal{N}\left(p_1, \frac{p_1(1-p_1)}{n}\right) \text{ и } \hat{p}_2 = \bar{Y} \stackrel{d}{\approx} \mathcal{N}\left(p_2, \frac{p_2(1-p_2)}{m}\right) - \text{ОМП}$$

При справедливости  $H_0$  получаем

$$W(X, Y) = \frac{\hat{p}_1 - \hat{p}_2}{\hat{\sigma}} \stackrel{d}{\approx} \mathcal{N}(0, 1),$$

$$\text{где } \hat{\sigma} = \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n} + \frac{\hat{p}_2(1-\hat{p}_2)}{m}}.$$



// Сходимость  $W(X, Y) \stackrel{d}{\rightarrow} \mathcal{N}(0, 1)$  при  $n, m \rightarrow +\infty$  можно доказать строго.

Критерий Вальда  $S = \{W(x, y) > z_{1-\alpha}\}$ .

$$\alpha = 0.05 \implies z_{1-\alpha} \approx 1.64, \quad S = \{W(x, y) > \mathbf{1.64}\}.$$

Дов. интервал для  $p_1 - p_2$  равен  $C = (\hat{p}_1 - \hat{p}_2 - z_{1-\alpha}\hat{\sigma}, 1)$ .

$H_0$  отвергается  $\iff 0 \notin C$



## Пример (влияние нового препарата на выздоровление)

- 1 группа:  $n = 30$  человек, 27 выздоровело  $\implies \hat{p}_1 = 0.9$   
2 группа:  $m = 30$  человек, 21 выздоровело  $\implies \hat{p}_2 = 0.7$   
 $W(x, y) \approx 2 \implies H_0$  отвергается, результат стат. значим  
дов. интервал  $(0.036, 1)$   $\leftarrow$  **слабая уверенность в результате**
- 1 группа:  $n = 30$  человек, 27 выздоровело  $\implies \hat{p}_1 = 0.9$   
2 группа:  $m = 30$  человек, 15 выздоровело  $\implies \hat{p}_2 = 0.5$   
 $W(x, y) \approx 3.76 \implies H_0$  отвергается, результат стат. значим  
дов. интервал  $(0.225, 1)$   $\leftarrow$  **хорошая уверенность в результате**
- 1 группа:  $n = 30$  человек, 27 выздоровело  $\implies \hat{p}_1 = 0.9$   
2 группа:  $m = 10$  человек, 7 выздоровело  $\implies \hat{p}_2 = 0.7$   
 $W(x, y) \approx 1.54 \implies H_0$  не отвергается, результат стат. незнач.  
дов. интервал  $(-0.017, 1)$   $\leftarrow$  **нет результата**





## Материал на доске





## p-value (достигаемый уровень значимости)

$$H_0 : P \in \mathcal{P}_0$$

$x_1, \dots, x_n$  — реализация выборки

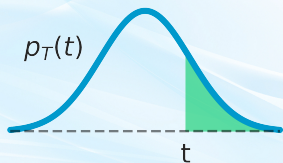
$T(x)$  — статистика критерия

$t = T(x_1, \dots, x_n)$  — реализация стат.

**p-value** — вероятность получить при справедливости  $H_0$  такое значение статистики  $t = T(x)$  или еще более экстремальное.

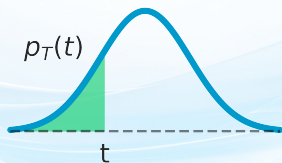
$$S = \{T(x) > c_\alpha\}$$

$$p(x) = P_0(T(X) \geq t),$$



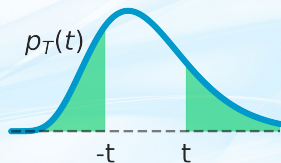
$$S = \{T(x) < c_\alpha\}$$

$$p(x) = P_0(T(X) \leq t),$$



$$S = \{|T(x)| > c_\alpha\}$$

$$p(x) = P_0(T(X) \geq |t|) + P_0(T(X) \leq -|t|),$$



Замечание. Если распр.  $T(X)$  при  $H_0$  не одинаково, то нужно добавить  $\sup$

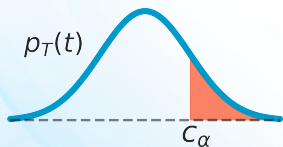
$P \in \mathcal{P}_0$



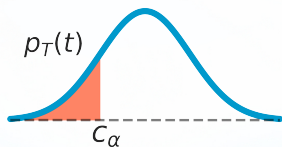
## В чем же разница? Графики одинаковые!!!

Еще раз посмотрим на них:

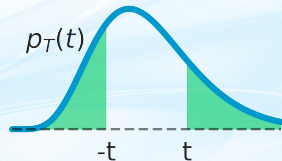
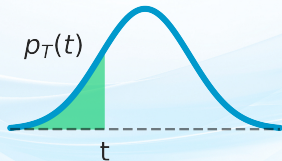
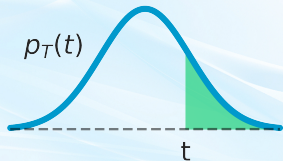
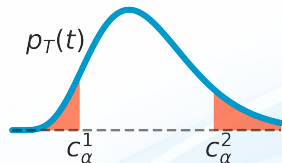
$$S = \{T(x) > c_\alpha\}$$



$$S = \{T(x) < c_\alpha\}$$



$$S = \{|T(x)| > c_\alpha\}$$

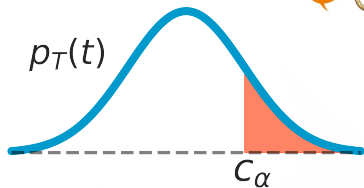




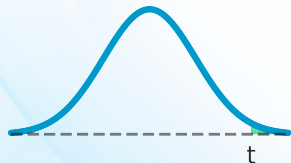
Рассмотрим случай  $S = \{T(x) > c_\alpha\}$

Критическое множество (слева) фиксировано.

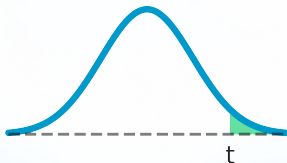
Ниже p-value для различных реализаций.



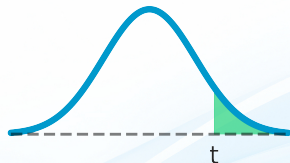
p-value(t) = 0.014



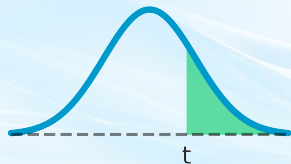
p-value(t) = 0.036



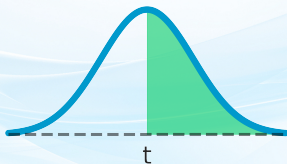
p-value(t) = 0.081



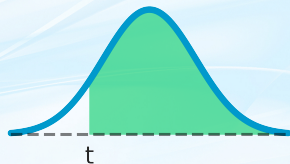
p-value(t) = 0.212



p-value(t) = 0.500



p-value(t) = 0.903





## Вывод:

$H_0$  отвергается  $\iff$  p-value  $\leq \alpha$

В этом случае p-value —  
степень уверенности в отвержении  $H_0$ .  
(чем p-value меньше, тем увереннее)



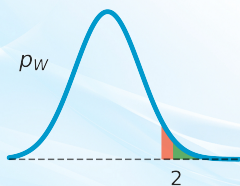
# Пример (влияние нового препарата на выздоровление)

Критерий  $S = \{W(x, y) > z_{1-\alpha}\}$ , где  $W(X, Y) \xrightarrow{d} \mathcal{N}(0, 1)$ .

p-value:  $p(w) = P(W(X, Y) > w) = \text{scipy.stats.norm.sf}(w)$ .

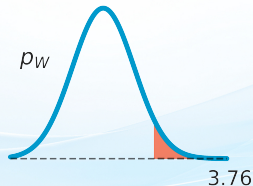
$$w = W(x) = 2$$

$$p(w) = 0.0228$$



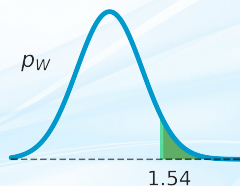
$$w = W(x) = 3.76$$

$$p(w) = 0.00008$$



$$w = W(x) = 1.54$$

$$p(w) = 0.0618$$





## Свойство p-value

### **Утверждение.**

Если при справедливости  $H_0$   
распр. статистики  $T(X)$  одинаково и непрерывно,  
то  $p(T(X)) \sim U[0, 1]$  при  $H_0$ .

### *Замечание.*

Часто на практике это верно, т.к. многие критерии так и строятся.

### *Следствие.*

Значения 0.2 и 0.9 одинаковы с точки зрения справедливости  $H_0$ ,  
т.е. p-value не есть степень уверенности в справедливости  $H_0$ .

Что возможно, если p-value большой?

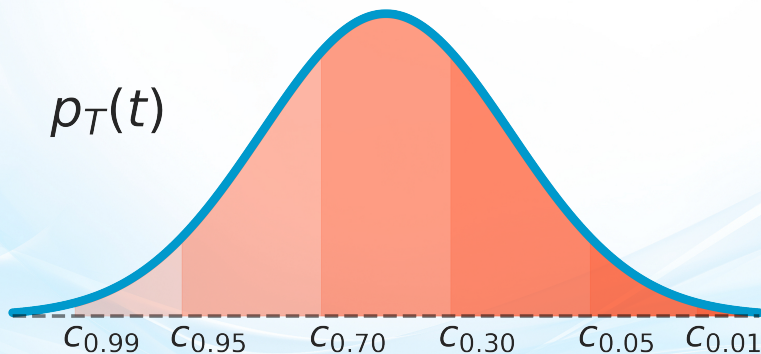
1.  $H_0$  верна;
2. Критерий недостаточно мощный.



## Общий случай p-value

$\{S_\alpha \mid \alpha \in [0, 1]\}$  — семейство критериев для разных уровней значимости.

$S_\alpha = \{T(x) > c_\alpha\}$  — критерий

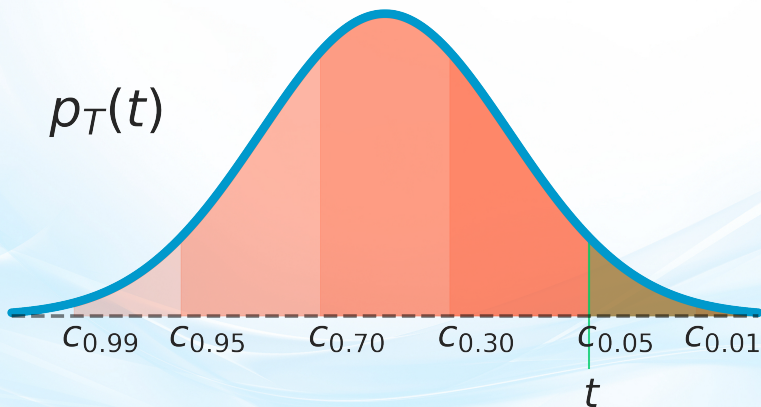




## Общий случай p-value

$\{S_\alpha \mid \alpha \in [0, 1]\}$  — семейство критериев для разных уровней значимости.

$S_\alpha = \{T(x) > c_\alpha\}$  — критерий

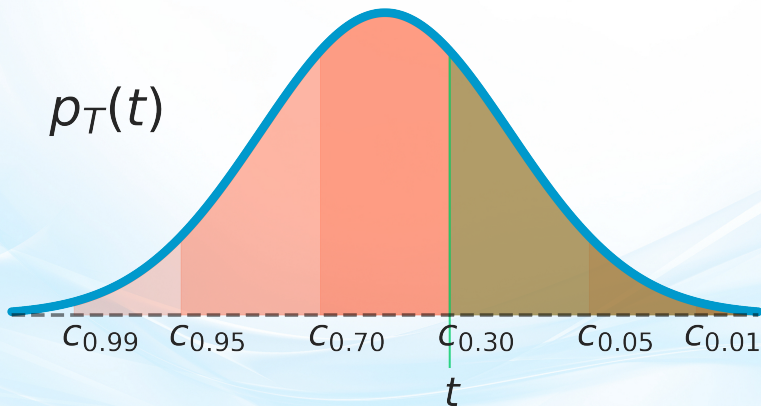


$p\text{-value}(t) = 0.05$

## Общий случай p-value

$\{S_\alpha \mid \alpha \in [0, 1]\}$  — семейство критериев для разных уровней значимости.

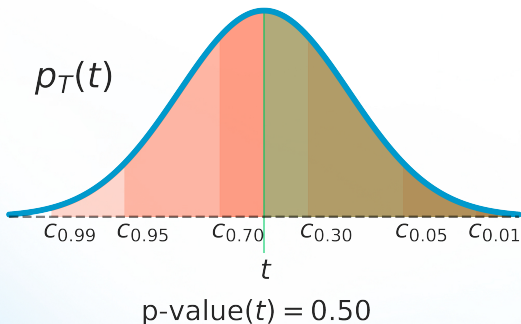
$S_\alpha = \{T(x) > c_\alpha\}$  — критерий



$$p\text{-value}(t) = 0.30$$



## Общий случай p-value



$t = c_{0.5} \implies$  при  $\alpha \geq 0.5$  гипотеза  $H_0$  отвергается.  
при  $\alpha < 0.5$  гипотеза  $H_0$  не отвергается.

**Ключевое наблюдение:** Если отвергнуть  $H_0$  можно только совершив большую ошибку, то скорее ее не стоит отклонять.



## Общий случай p-value

### Вывод:

p-value — наименьший уровень значимости, при котором  $H_0$  можно отвергнуть для данной реализации выборки  $x$ .

$$p(x) = \inf \{ \alpha \in [0, 1] \mid x \in S_\alpha \}$$



# Что не есть p-value

Величина p-value не является

- ▶ уровнем значимости, реальным уровнем значимости, вероятностью ошибки первого рода;  
**не зависят от выборки**
- ▶ вероятностью  $H_0$ , вероятностью  $H_0$  при условии выборки;  
**она либо верна, либо нет**
- ▶ много чем еще.

## Что не есть p-value (Пример)

На ЧМ по футболу в 2010 г. осьминог Пауль предсказывает результаты матчей с участием сборной Германии, выбирая кормушку с флагом страны-победителя.



$X_1, \dots, X_n \sim \text{Bern}(p)$  — результаты предсказания (правильно/нет).

$H_0: p = 1/2$  vs.  $H_1: p > 1/2$  (наугад vs. не наугад)

Критерий  $S = \{T(x) > c_\alpha\}$ , где  $T(X) = \sum X_i \sim \text{Bin}(n, p)$ .

$$p\text{-value: } p(t) = \frac{1}{2^n} \sum_{j=t}^n C_n^j$$

13 матчей: Пауль верно угадывает исход матча в 11 случаях.

$$p(11) = 2^{-13} (C_{13}^{11} + C_{13}^{12} + C_{13}^{13}) \approx 0.0112 < 0.05 \implies H_0 \text{ отвергается;}$$

0.0112 не является вероятностью того, что Пауль выбирает кормушку наугад







# 5. Проверка статистических гипотез

## 5.5. Практическая значимость результата

# Большие выборки

$X_1, \dots, X_n \sim \text{Bern}(\theta)$  — результаты испытания Пауля.

$H_0 : \theta = 1/2$  vs.  $H_1 : \theta > 1/2$

Критерий  $S = \{T(x) \geq c_\alpha\}$ , где  $T(X) = \sum X_i \sim \text{Bin}(n, p)$ .

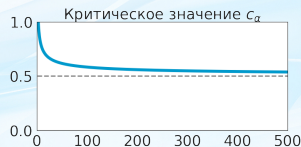
Как критическое значение  $c_\alpha$  зависит от  $n$ ?

Рассмотрим асимптотически эквивалентный критерий Вальда

$$W(X) = \sqrt{n} \frac{\bar{X} - 1/2}{\sqrt{1/4}} \xrightarrow{d_{1/2}} \mathcal{N}(0, 1).$$

Тогда критерий

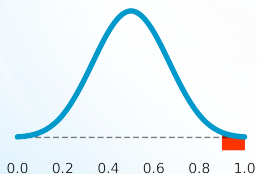
$$S_W = \{W(x) > z_{1-\alpha}\} = \left\{ \bar{x} > \frac{1}{2} + \frac{z_{1-\alpha}}{2\sqrt{n}} \right\}.$$



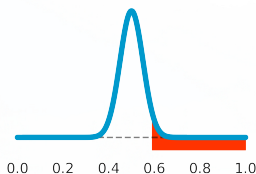


# Большие выборки

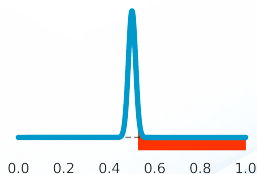
$n = 10, c_\alpha = 0.9$



$n = 100, c_\alpha = 0.59$



$n = 1000, c_\alpha = 0.527$



**Вывод:** при  $n \rightarrow +\infty$  мощность критерия сходится к 1.

**Теория:** это замечательно!

**Практика:** не совсем...

**Теория:** даже при небольшом отличии истины от  $H_0$  мы ее отклоним!

**Практика:** Хахаха, какой смысл в осьминоге, который угадывает с вероятностью 0.51? На сковородку его!



# Вывод с точки зрения практики

Как правило, на практике:

**1. При малом размере выборки:**

Почти ничего не отклоняется, т.к. мощность небольшая.

*Недообучение*

**2. При большом размере выборки:**

Отклоняется почти все, т.к. небольшие отличия от  $H_0$  есть почти всегда.

*Переобучение*



# Практическая значимость

**Размер эффекта** — величина, оценивающая по данным, насколько основная гипотеза отличается от истины.

## Пример 1.

В течении трех лет женщины выполняли физические упражнения в двух группах: группа 1 — не менее часа в день, группа 2 — не более 20 минут в день.

$H_0$ : изменение веса в обеих группах одинаково

$H_1$ : в 1 группе уменьшение веса больше, чем во 2-й

$p\text{-value} < 0.001 \implies$  результат статистически значим

Разница в весе **150 грамм**  $\implies$  результат практически не значим



## Практическая значимость

**Размер эффекта** — величина, оценивающая по данным, насколько основная гипотеза отличается от истины.

### Пример 2.

В 2002 году проводились клинические испытания гормонального препарата Премарин, которые были досрочно прерваны.

На **0.08%** увеличивается риск развития рака груди;

На **0.08%** увеличивается риск инсульта;

На **0.07%** увеличивается риск инфаркта.

С учетом численности населения есть практическая значимость.



## Практическая значимость

	Есть практ. значимость	Нет практ. значимости
Есть стат. значимость	$H_0$ отвергается: Эффект присутствует и доказан статистически	Скорее всего в полученном результате смысла нет
Нет стат. значимости	Эффект присутствует, но не доказан статистически; нужно продолжать эксперимент	$H_0$ не отвергается: результата нет

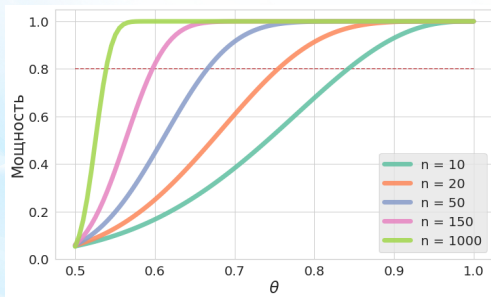
# План эксперимента

Как определить размер выборки **до** эксперимента?

$$X_1, \dots, X_n \sim \text{Bern}(\theta)$$

$$H_0 : \theta = 1/2 \text{ vs. } H_1 : \theta > 1/2$$

Графики мощности для критерия  $S = \{\sum X_i > c_\alpha\}$ :



$\alpha = 0.05$  — ур. знач.

*Желаемые значения:*

$\beta = 0.8$  — мощность;

$\theta \geq 0.6$  — значимый эффект.

Выбираем  $n$ , для которого кривая мощности проходит через точку  $(0.6, 0.8)$ .





# 5. Проверка статистических гипотез

## 5.6. Множественная проверка гипотез



# Поиск экстрасенсов

Этап 1: Угадайте цвета (**синий** и **оранжевый**) с учетом порядка.





# Поиск экстрасенсов

**Этап 1:** ответы.





## Поиск экстрасенсов

Этап 2: Угадайте цвета (**синий** и **оранжевый**) с учетом порядка.





# Поиск экстрасенсов

**Этап 2:** ответы.





## Поиск экстрасенсов

В 1950 г. проводились испытания  
возможности экстрасенсорного восприятия.

Этап 1: поиск экстрасенсов — испытуемому нужно угадать цвет 10 карт.

$X_1, \dots, X_{10} \sim \text{Bern}(\theta)$  — результаты (правильно / нет).

$H_0 : \theta = 1/2$  vs.  $H_1 : \theta > 1/2$  (наугад vs. не наугад)

Критерий  $S = \{T(x) \geq c_\alpha\}$ , где  $T(X) = \sum X_i \sim \text{Bin}(n, p)$ .

$c$	7	8	9	10
$P_{1/2}(T(X) \geq c)$	0.172	0.055	0.010	0.001

Берем  $c_\alpha = 9$ , т.е.  $H_0$  отклоняется если  $\sum X_i \geq 9$ .



## Поиск экстрасенсов

Вывод: если человек верно отгадывает хотя бы 9 карт из 10, то он становится предполагаемым экстрасенсом.

В эксперименте приняли участие 1000 человек, при этом

- ▶ 9 карт верно отгадали 9 человек;
- ▶ 10 карт верно отгадали 2 человека.

В дальнейшем ни один из них не подтвердил свои способности...

$$P_{1/2}(\text{хотя бы один из 1000 угадает 9 или 10 карт верно}) = \\ = 1 - (1 - C_{10}^9/2^{10} - C_{10}^{10}/2^{10})^{1000} = 1 - (1 - 11/2^{10})^{1000} \approx 0.99997$$



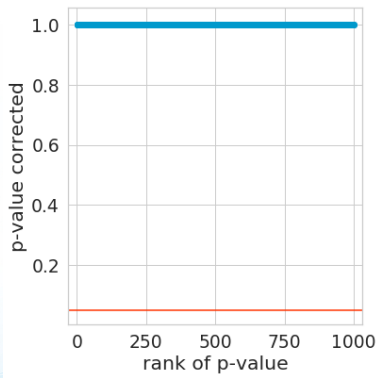
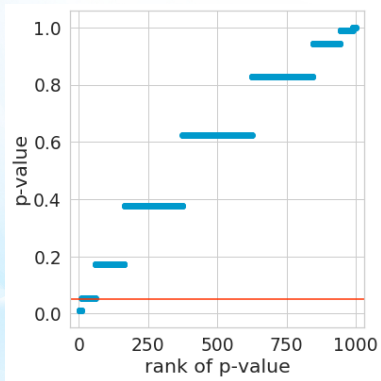
## Материал на доске





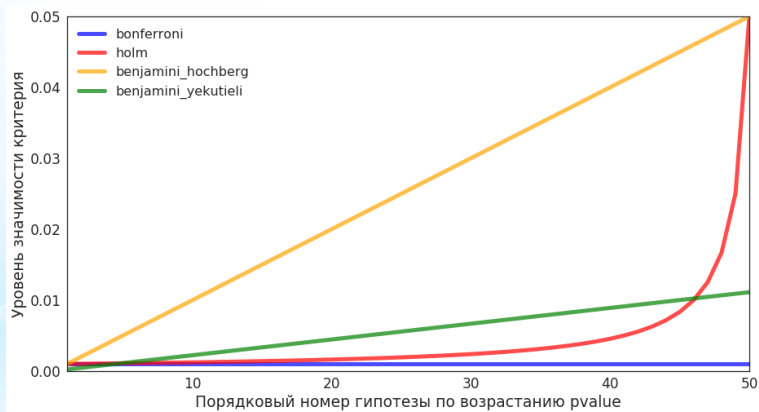


# Поиск экстрасенсов





# Сравнение методов МПГ





# Реализация МПГ

```
statsmodels.stats.multitest.multipletests  
(pvals, alpha=0.05, method='hs',  
is_sorted=False, returnsorted=False)
```

- ▶ `pvals` — значения `p-value` по всем критериям
- ▶ `alpha` — желаемый уровень значимости
- ▶ `method`:
  - ▶ `bonferroni`
  - ▶ `sidak`
  - ▶ `fdr_bh`
  - ▶ `holm`
  - ▶ `holm-sidak`
  - ▶ `fdr_by`

Возвращает:

- ▶ `reject` — для отвергаемых гипотез `True`
- ▶ `pvals_corrected` — скорректированные `p-value`

## Простой пример

Знакомая задача:

$$X_1, \dots, X_n \sim \mathcal{N}(\theta, 1)$$

$$H_0: \theta \geq 0 \text{ vs } H_1: \theta < 0$$

$$\text{PHMK: } S = \{x \in \mathbb{R} \mid \bar{x} \leq c_\alpha\}$$

Пусть теперь две одинаковые задачи с независимыми выборками:

$$X_1, \dots, X_n \sim \mathcal{N}(\theta_1, 1)$$

$$Y_1, \dots, Y_n \sim \mathcal{N}(\theta_2, 1)$$

$$H_1: \theta_1 \geq 0 \text{ vs } H'_1: \theta_1 < 0$$

$$H_2: \theta_2 \geq 0 \text{ vs } H'_2: \theta_2 < 0$$

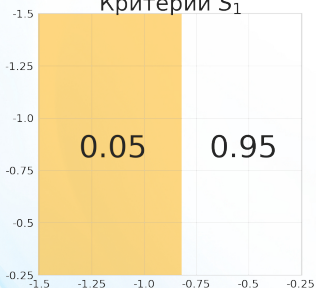
$$\text{Критерии: } S_1 = \{(x, y) \in \mathbb{R}^2 \mid \bar{x} \leq c_\alpha\}$$

$$S_2 = \{(x, y) \in \mathbb{R}^2 \mid \bar{y} \leq c_\alpha\}$$

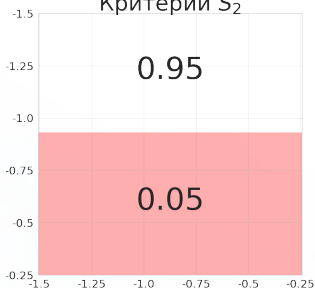
**Частая ошибка:** Выборки независимы  $\implies$  МПГ не нужна.



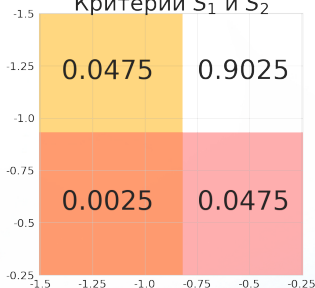
Критерий  $S_1$



Критерий  $S_2$



Критерии  $S_1$  и  $S_2$

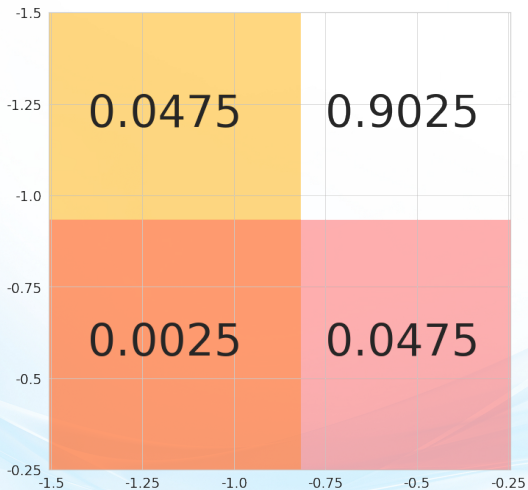


## Вывод:

вероятность допустить хотя бы одну ошибку равна 0.0975,  
если обе основные гипотезы верны



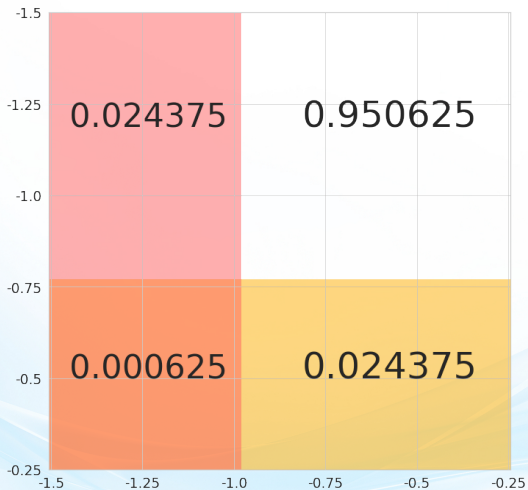
## Сравнение: без корректировки



Вероятности указаны при справедливости  $H_1$  и  $H_2$



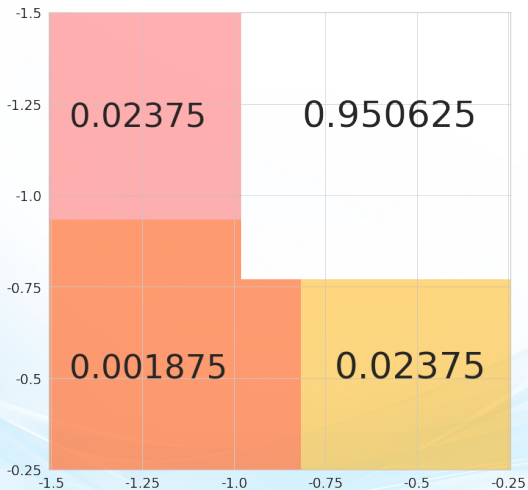
## Сравнение: метод Бонферрони



Вероятности указаны при справедливости  $H_1$  и  $H_2$



## Сравнение: метод Холма

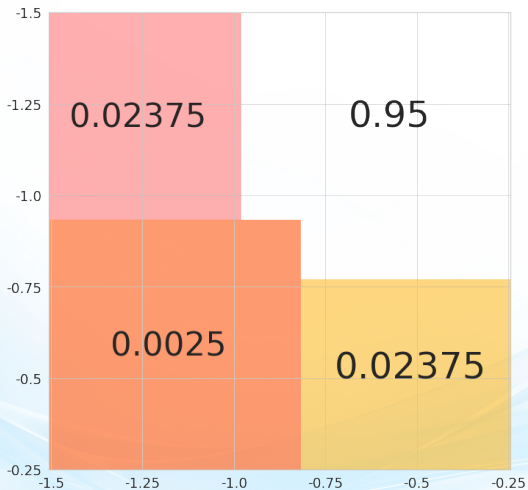


Вероятности указаны при справедливости  $H_1$  и  $H_2$





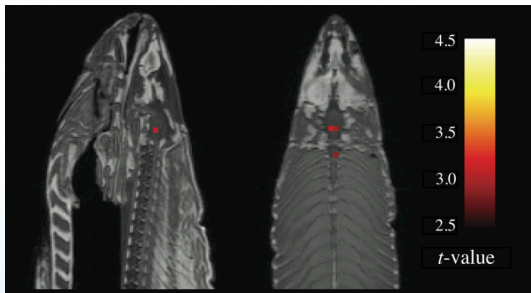
## Сравнение: метод Бенджамини-Хохберга (не контр. FWER)



Вероятности указаны при справедливости  $H_1$  и  $H_2$

## Удивительные открытия

2009 год. МРТ мозга мертвого самца лосося:



МРТ дает 3D-изображение на 130 000 вокселей.

**Эксперимент:** Лососю показывали фото и просили его пояснить, какие эмоции испытывают люди с картинки.

**Обработка:** Для каждого вокселя тестируется гипотеза о наличии активации этого участка мозга.

## Удивительные открытия

**Результат:** Для каждой картинке для нескольких вокселей мозга p-value оказывалось меньше 0.001.

**Вывод:** мертвый лосось реагирует на все!!!

Авторы удостоились Шнобелевской премии (2012 год) за открытие в области неврологии.

При применении МПГ лосось переставал на что-либо реагировать...

<http://prefrontal.org/files/posters/Bennett-Salmon-2009.pdf>



### Neural correlates of interspecies perspective taking in the post-mortem Atlantic Salmon: An argument for multiple comparisons correction

Craig M. Bennett<sup>1</sup>, Abigail A. Baird<sup>2</sup>, Michael B. Miller<sup>1</sup>, and George L. Wolford<sup>3</sup>

<sup>1</sup> Psychology Department, University of California Santa Barbara, Santa Barbara, CA; <sup>2</sup> Department of Psychology, Vassar College, Poughkeepsie, NY;

<sup>3</sup> Department of Psychological & Brain Sciences, Dartmouth College, Hanover, NH

#### INTRODUCTION

With the extreme dimensionality of functional neuroimaging data comes extreme risk for false positives. Across the 130,000 voxels in a typical fMRI volume the probability of a false positive is almost certain. Correction for multiple comparisons should be completed with these datasets, but is often

#### GLM RESULTS





**ВСЁ!**