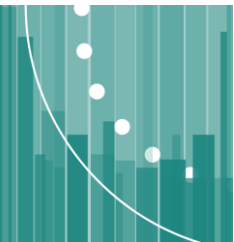


Введение в статистику

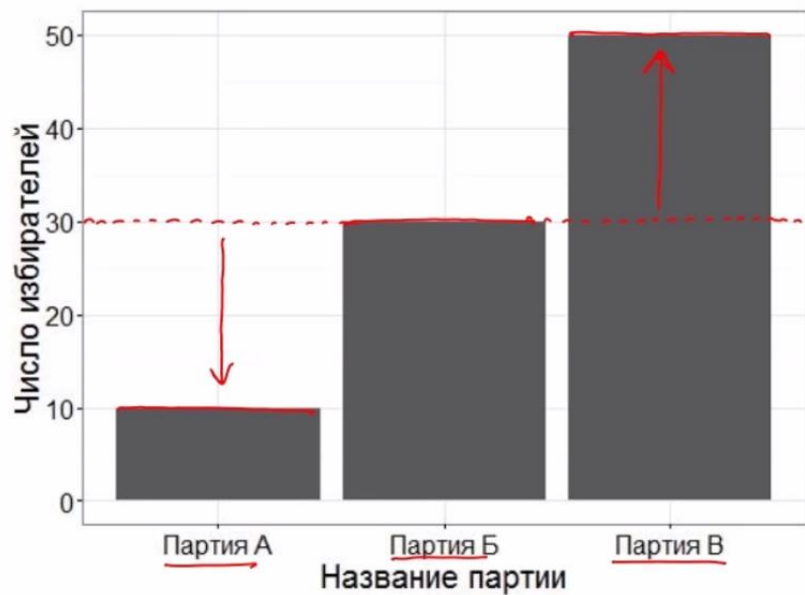
Часть 2

Неделя 1

Анализ номинативных данных



Проверка гипотезы о распределении номинативной переменной



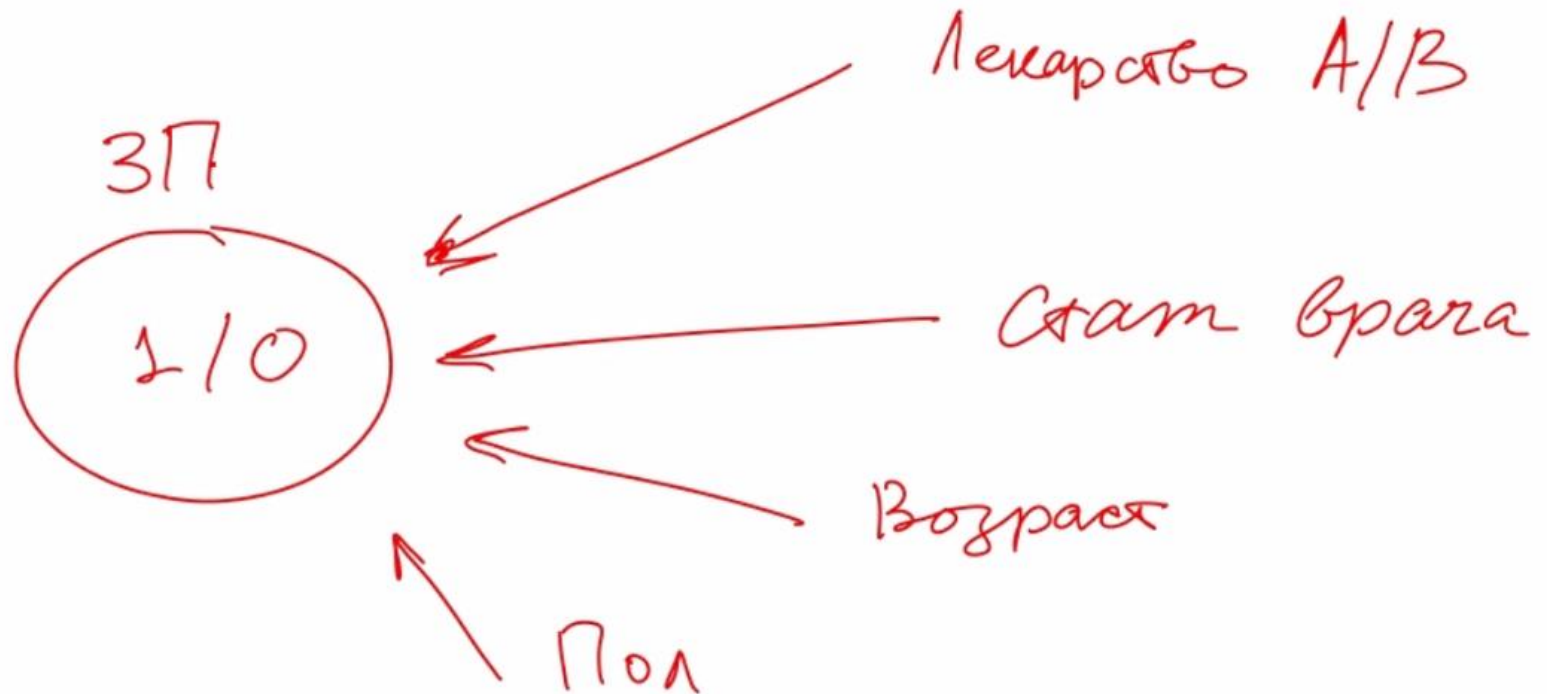
N	Партия	<u>A</u>	<u>B</u>	<u>B</u>
1	A			
2	B			
3	B			
...				
90	B			

Проверка гипотезы о взаимосвязи двух номинативных переменных

Лекарство	Результат
<u>A</u>	1
A	1
A	0
<u>B</u>	0
B	0
B	1

	Лекарство	
	A	B
1	2	1
0	1	2

Более сложные модели



Расчет Евклидова расстояния





	Решка	Орел
О	20	40
Е	30	30

$N=60$ H_0 $p_{орла} = 0,5$
 H_1 $p_{орла} \neq 0,5$

$$(20 - 30)^2 + (40 - 30)^2 = 200$$

$$\begin{array}{cc} 1020 & 1040 \\ 1030 & 1030 \end{array} \quad (1020 - 1030)^2 + (1040 - 1030)^2 = 200$$

Расчет расстояния Хи – квадрат

	Решка	Орел
О	<u>20</u> 	40 
Е	<u>30</u> 	30 

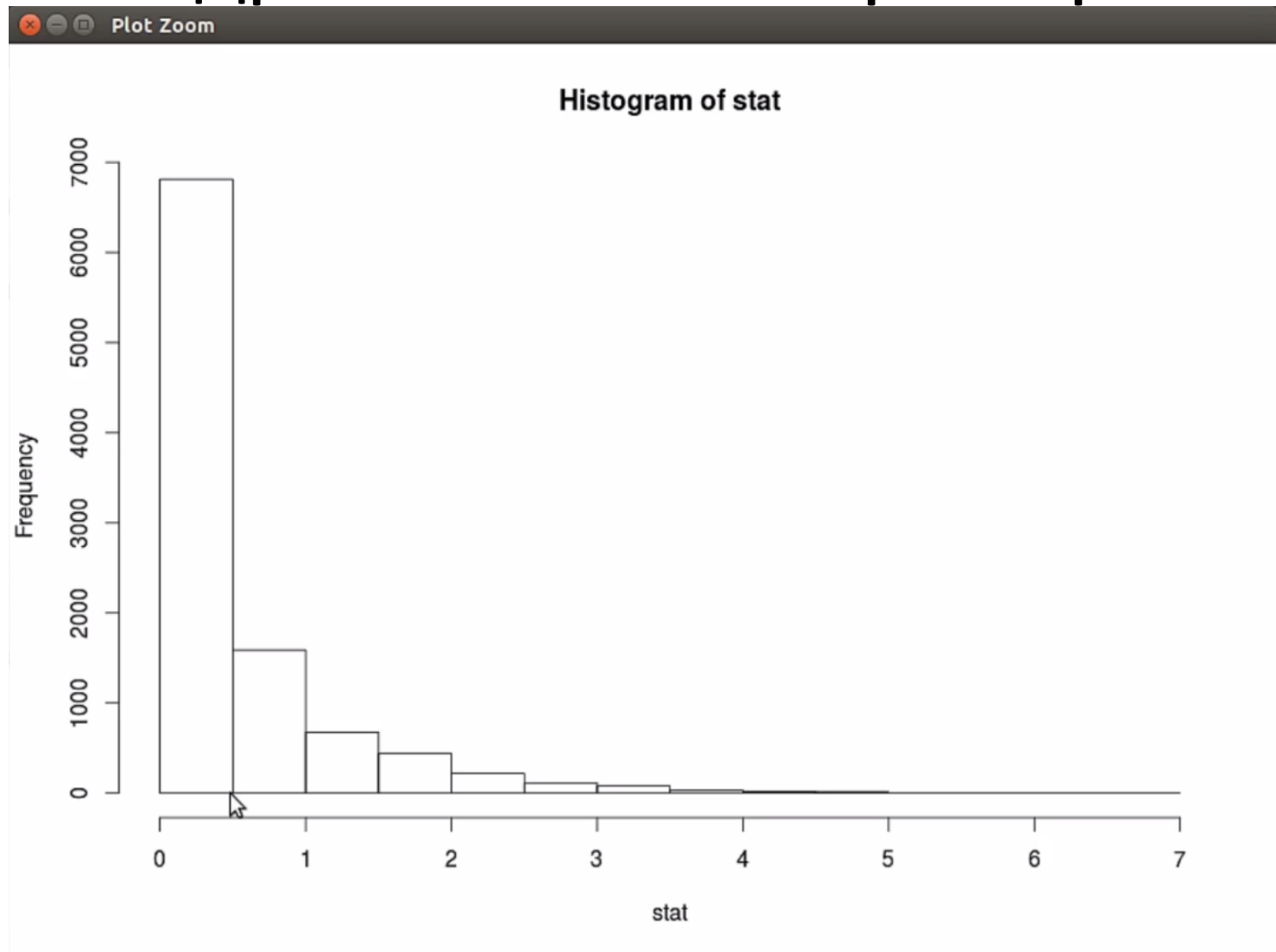
$N=60$ $\underline{H_0}$ $p_{орла} = 0,5$
 H_1 $p_{орла} \neq 0,5$

$$\chi^2 = \left(\frac{20 - 30}{\sqrt{30}} \right)^2 + \left(\frac{40 - 30}{\sqrt{30}} \right)^2 = \frac{(20 - 30)^2}{30} + \frac{(40 - 30)^2}{30} =$$

$$= \frac{100}{30} + \frac{100}{30} \approx \textcircled{6,7}$$

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i}$$

Распределение расстояние Хи – квадрат из нашего примера



Распределение расстояние Хи – квадрат из нашего примера

Решка

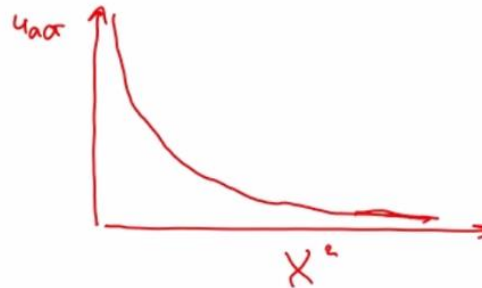
O_1

E_1

Орел

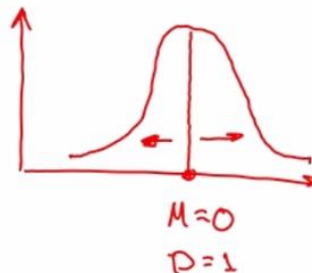
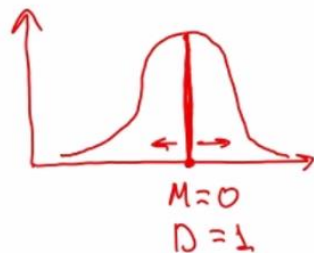
O_2

E_2



H_0

$$\chi^2 = \left(\frac{O_1 - E_1}{\sqrt{E_1}} \right)^2 + \left(\frac{O_2 - E_2}{\sqrt{E_2}} \right)^2 = (\text{wavy line})^2 + (\text{wavy line})^2$$



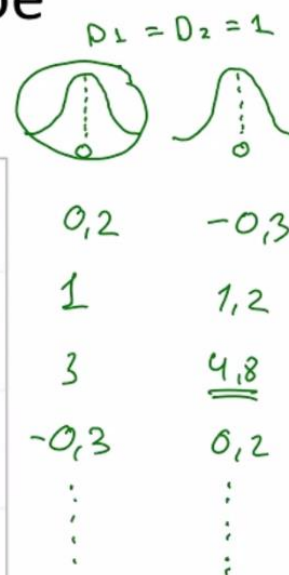
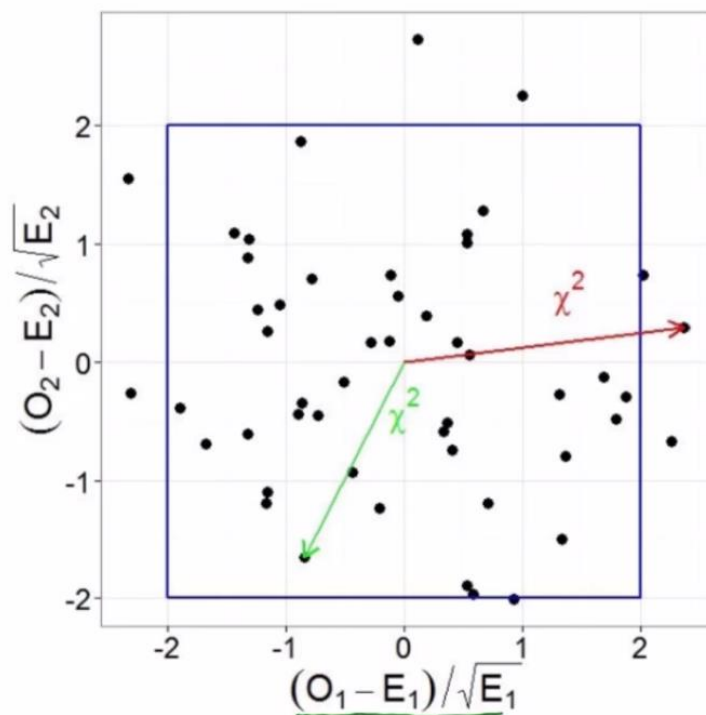
Определение

Распределение хи-квадрат с k степенями свободы — это распределение суммы квадратов k независимых стандартных нормальных случайных величин.

$$\chi^2 = \left(\underbrace{\text{Normal}(0,1)}_{\substack{M=0 \\ D=1}} \right)^2 + \left(\underbrace{\text{Normal}(0,1)}_{\substack{M=0 \\ D=1}} \right)^2$$


График распределения двух независимых случайных стандартных нормальных величин

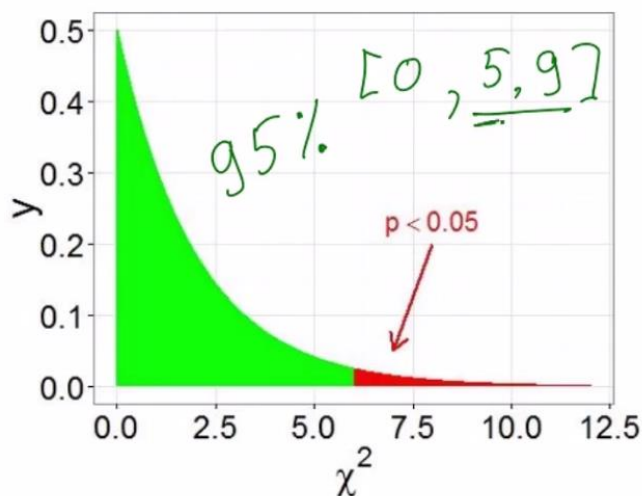
Двумерное нормальное
распределение



Критическое значение

Итак, как ведет себя χ^2

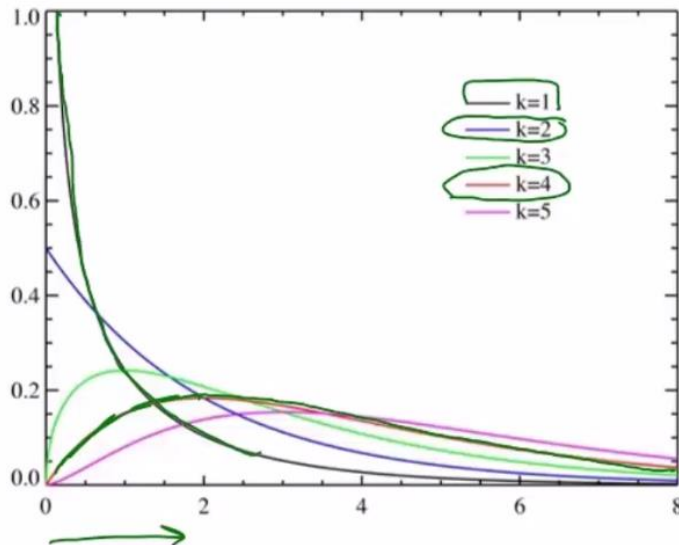
Распределение χ^2 с **двумя степенями (df = 2)** свободы имеет следующий вид:



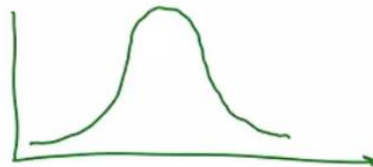
Критическое значение χ^2 для $p < 0.05$ равняется 5.9

Степени свободы

Распределение хи-квадрат с k степенями свободы — это распределение суммы квадратов k независимых стандартных нормальных случайных величин.



$$\chi^2 = \underbrace{(\overset{\approx 0}{\sim})^2 + (\overset{\approx 0}{\sim})^2 + (\overset{\approx 0}{\sim})^2 + (\overset{\approx 0}{\sim})^2}_{df=4}$$
$$\chi^2 = \left(\underset{0}{\sim} \right)^2$$



Число степеней свободы в нашем примере

Stepic.org

Решка	Орел
20	40
30	30

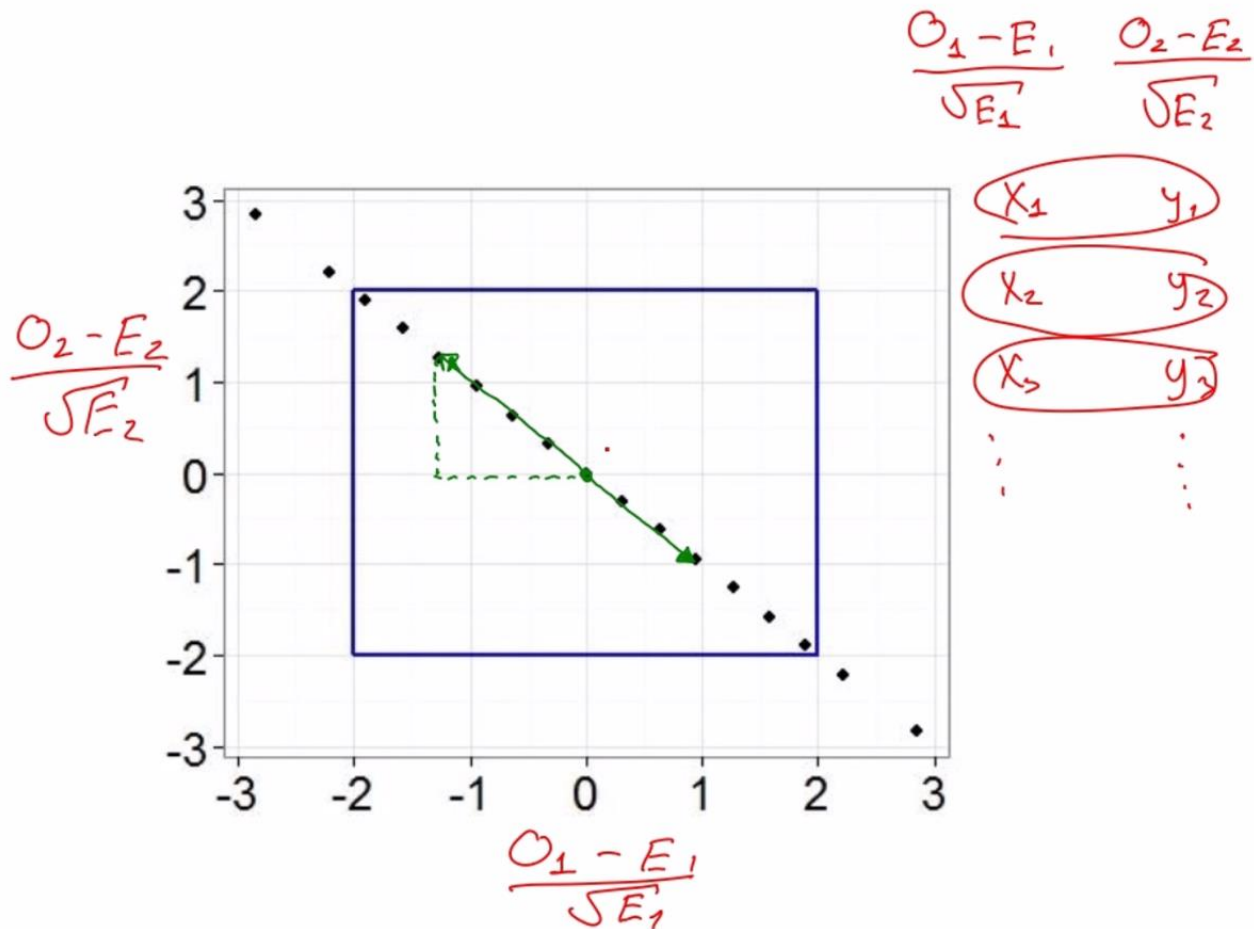
$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

df = 1

10 50
40 20

$$\chi^2 = (\underline{10})^2 + (\underline{10})^2$$

Решка	Орел
$\frac{O_1 - E_1}{\sqrt{E_1}}$	$\frac{O_2 - E_2}{\sqrt{E_2}}$
$\frac{-20}{\sqrt{30}}$	$\frac{20}{\sqrt{30}}$
$\frac{10}{\sqrt{30}}$	$\frac{-10}{\sqrt{30}}$

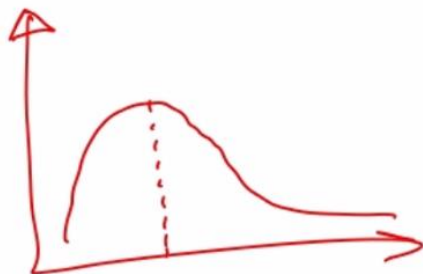


Число степеней свободы для игральной кости

	1	2	3	4	5	6
O	10	10	10	5	10	15
E	10	10	10	10	10	10

$$\underline{\underline{N=60}}$$

$$\sum \frac{(O_i - E_i)^2}{E_i} \sim \chi^2 \quad df=5$$



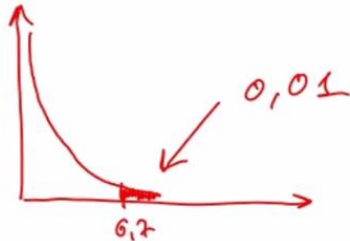
Расчет p - value

Решка	Орел
20	40
<u>30</u>	30

$$\underline{H_0 : p_{орла} = 0,5}$$

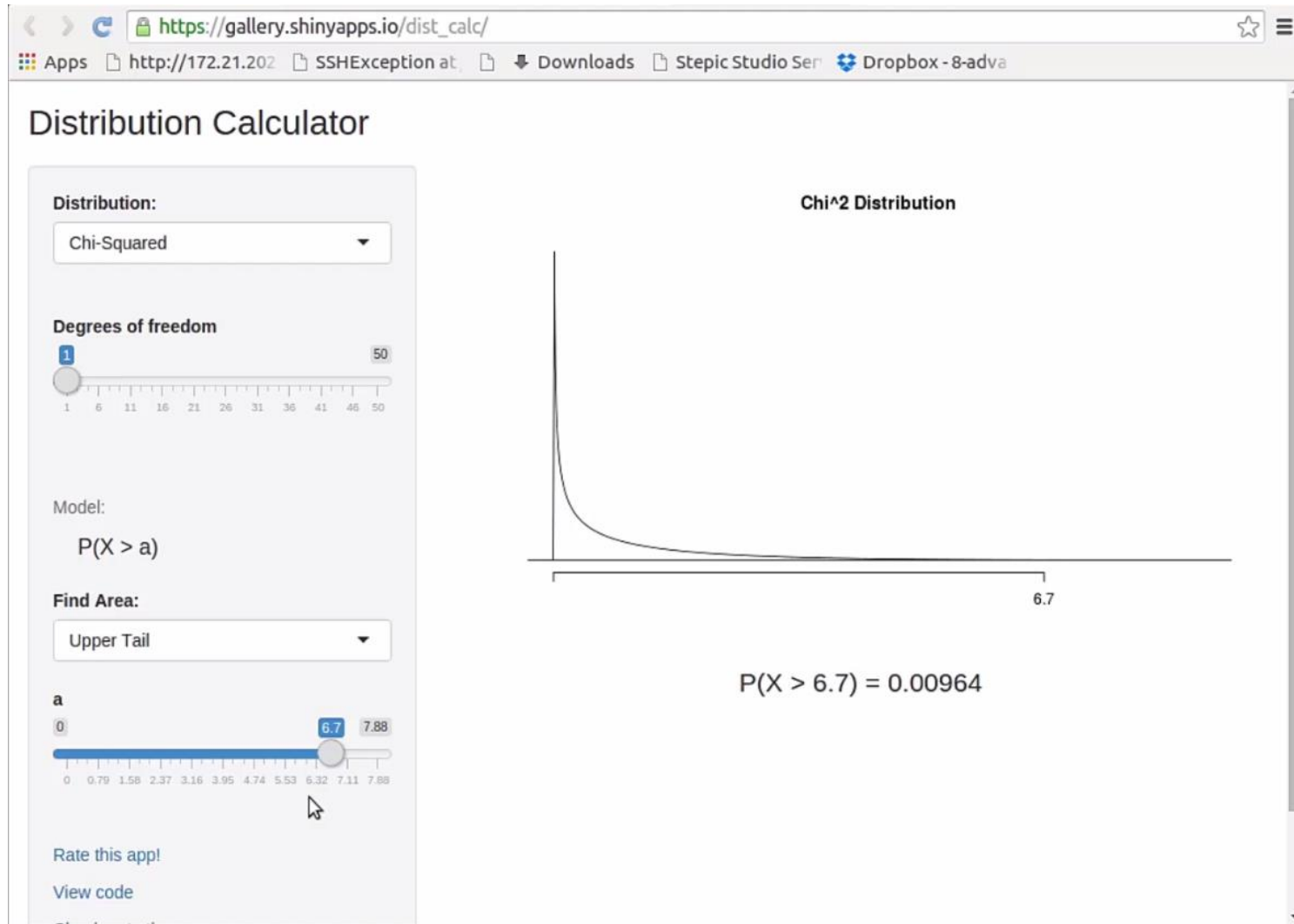
$$\underline{\chi^2} = \frac{(20 - 30)^2}{30} + \frac{(40 - 30)^2}{30} = \frac{100}{30} + \frac{100}{30}$$

$$\underline{\underline{\chi = 6,7}}$$



$$\underline{\underline{p < 0,05}} \quad \text{No.}$$

Расчет p - value



Анализ таблиц сопряженности

	Юноши	Девушки
Биологи	15	9
Информатики	11	6

Две номинативные переменные пол и профессия, обе с двумя градациями.

Нулевая гипотеза - распределение не отличается от ожидаемого.

Альтернативная гипотеза – распределение отличается или иными словами: две переменные взаимосвязаны между собой

Расчет ожидаемых значений

	Юноши	Девушки	Всего
Биологи	15	9	24
Информатики	11	6	17
Всего	26	15	41

Юношей 26 человек из 41 – 63.4 %

Девушек 15 человек из 41 – 36.6 %

Примем нулевую гипотезу, что профессия никак не связана с полом, тогда юноши должны с равной частотой наблюдаться у информатиков и биологов.

63.4 % от 24 и 17 составляют 15.2 и 10.8 соответственно.

36.6 % от 24 и 17 составляет 8.8 и 6.2 соответственно.

Расчет ожидаемых значений

	Юноши	Девушки	Всего
Биологи	15	9	24
Информатики	11	6	17
Всего	26	15	41

VS

	Юноши	Девушки	Всего
Биологи	15.2	8.8	24
Информатики	10.8	6.2	17
Всего	26	15	41

$$f_{ij} = \frac{f_i \cdot f_j}{N}$$

Критерий χ^2 - Пирсона

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

Проверяет гипотезу о том, что наблюдаемое распределение номинативной переменной отличается от ожидаемого.

Для одной номинативной переменной **df = n - 1**, где n - количество столбцов таблицы.

Для таблиц сопряженности **df = (n - 1) * (m - 1)**, где n - количество столбцов таблицы, m - число строк таблицы.

✓ *Все наблюдения независимы*

✓ *Минимальное количество наблюдений в каждой из ячеек должно быть больше 5.*

Поправка Йетса

	Ю	Д
Б	15	9
И	11	6

	Ю	Д
Б	15.2	10.8
И	8.8	6.2

$$\chi^2_{Yates} = \sum_k \frac{(|f_0 - f_e| - 0.5)^2}{f_e}$$

В теории распределение χ^2 непрерывно, тогда как вычисляемые значения всегда дискретны, в результате H_0 может отвергаться слишком часто. Чтобы скорректировать значения p - уровня значимости применяется поправка Йетса на непрерывность. Обычно применяется, когда некоторые ожидаемые частоты меньше 10.

Расчет критерия

	Ю	Д
Б	15	9
И	11	6

	Ю	Д
Б	15.2	10.8
И	8.8	6.2

$$\chi_{Yates}^2 = \sum_k \frac{(|f_0 - f_e| - 0.5)^2}{f_e}$$

$$\chi^2 = (|15 - 15.2| - 0.5)^2 / 15.2 + \dots + (|6 - 6.2| - 0.5)^2 / 6.2$$

$$df = (n - 1)(m - 1) = 1$$

```
> students <- rbind(c(15,9), c(11, 6))  
> chisq.test(students)
```

Pearson's Chi-squared test with Yates'
continuity correction

```
data:  students  
X-squared = 1.2684e-31, df = 1, p-value = 1
```

Нельзя ли снизить риск тромбоза назначением небольших доз аспирина (160 мг/сут)? *

	Есть тромбоз	Нет тромбоза
Плацебо	18	7
Аспирин	6	13

*H. R. Harter, J. W. Burch, P. W. Majerus, N. Stanford, J. A. Delmez, C. B. Anderson, C. A. Weerts. Prevention of thrombosis in patients in hemodialysis by low-dose aspirin. *N. Engl. J. Med.*, 301:577—579, 1979.

Расчет критерия

	Есть тромбоз	Нет тромбоза
Плацебо	18	7
Аспирин	6	13

```
> patients <-rbind(c(18,7), c(6, 13))  
> chisq.test(patients)
```

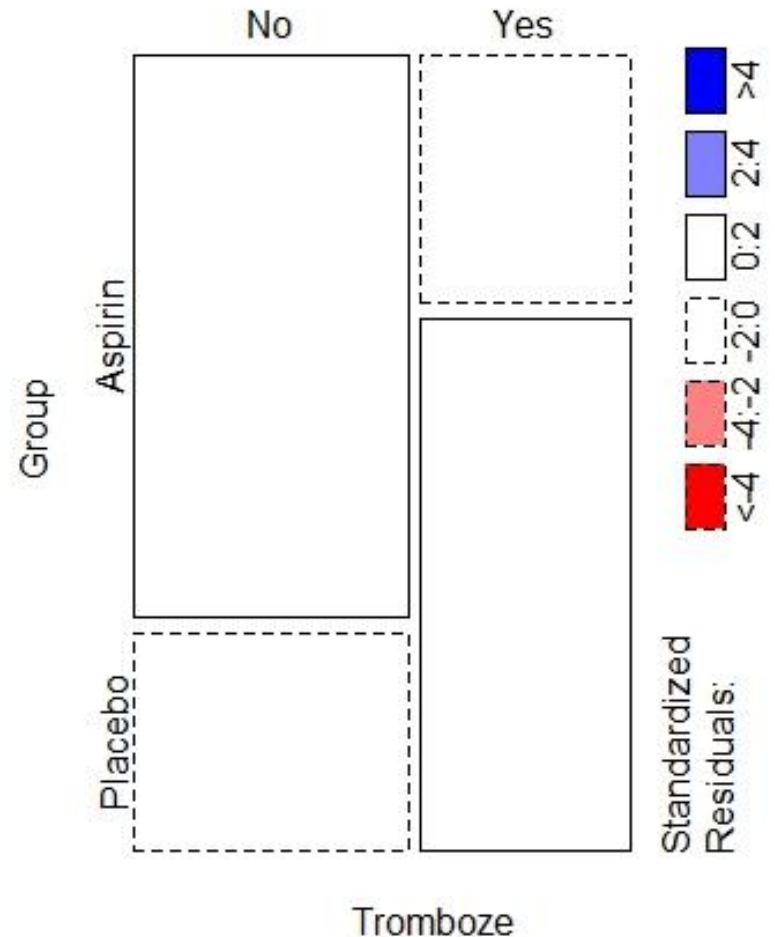
```
Pearson's Chi-squared test with Yates'  
continuity correction
```

```
data: patients  
X-squared = 5.5772, df = 1, p-value = 0.0182
```

Уточним результаты

Анализ остатков позволяет выявить какие именно частоты значимо отклоняются от ожидаемых значений.

```
> test$observed
      Aspirin Placebo
No         18         7
Yes         6        13
> test$expected
      Aspirin  Placebo
No  13.63636 11.363636
Yes 10.36364  8.636364
```



Точный критерий Фишера

	Поправился <u>+</u>	Не поправился <u>-</u>	Сумма
Лекарство №1	<u>a</u>	<u>b</u>	<u>a+b</u>
Лекарство №2	c	d	<u>c+d</u>
Сумма	<u>a+c</u>	<u>b+d</u>	<u>n</u>

$$H_0: P_1 = P_2 = P$$

X - число "+" N1

Y - число "+" N2

$$X \sim \text{binomial}(a+b, p)$$

$$Y \sim \text{binomial}(c+d, p)$$

$$X+Y \sim \text{binomial}(n, p)$$

$$P(X=a) = C_{a+b}^a \cdot p^a \cdot (1-p)^b$$

$$P(Y=c) = C_{c+d}^c \cdot p^c \cdot (1-p)^d$$

$$P(X+Y=a+c) = C_n^{a+c} \cdot p^{a+c} \cdot (1-p)^{b+d}$$

Точный критерий Фишера

$$H_0: P_1 = P_2 = P$$

X - кол-во "+" N1

Y - кол-во "+" N2

	Поправился +	Не поправился -	Сумма
<u>Лекарство №1</u>	<u>a</u>	<u>b</u>	<u>a+b</u>
<u>Лекарство №2</u>	c	d	<u>c+d</u>
<u>Сумма</u>	<u>a+c</u>	<u>b+d</u>	<u>n</u>

$$P(X=a) = C_{a+b}^a \cdot p^a \cdot (1-p)^b$$

$$P(Y=c) = C_{c+d}^c \cdot p^c \cdot (1-p)^d$$

$$P(X+Y=a+c) = C_n^{a+c} \cdot p^{a+c} \cdot (1-p)^{b+d}$$

$$\begin{aligned} P(X=a | X+Y=a+c) &= \frac{P(X=a, X+Y=a+c)}{P(X+Y=a+c)} = \\ &= \frac{P(X=a) \cdot P(Y=c)}{P(X+Y=a+c)} \end{aligned}$$

$$\frac{C_{a+b}^a \cdot \cancel{p^a} \cdot \cancel{(1-p)^b} \cdot C_{c+d}^c \cdot \cancel{p^c} \cdot \cancel{(1-p)^d}}{C_n^{a+c} \cdot \cancel{p^{a+c}} \cdot \cancel{(1-p)^{b+d}}} =$$

$$\frac{C_{a+b}^a \cdot C_{c+d}^c}{C_n^{a+c}}$$

Точный критерий Фишера



	Поправился	Не поправился	Сумма
Лекарство №1	a=4 0	b=0 4	a+b=4
Лекарство №2	c=0 4	d=4 0	c+d=4
Сумма	a+c=4	b+d=4	n=8

$$\frac{C_a^a \cdot C_c^c}{C_n^{a+c}}$$

$$\frac{C_4^3 \cdot C_4^1}{C_8^4} \approx 0,229$$

3	1	4	0
1	3	0	4
1	3	0	4
3	1	4	0

$$(0,229 + 0,014) * 2 \approx \underline{\underline{0,49}}$$

$$\frac{C_4^4 \cdot C_4^0}{C_8^4} \approx 0,014$$

RStudio

File Edit Code View Plots Session Build Debug Tools Help

Go to file/function

Project: (None)

Environment History

Global Environment

Data

avian 1070 obs. of 1...

Values

a	331776
b	23224320
cover	chs [1.71 "DNR"

Files Plots Packages Help Viewer

Zoom Export

```
1 fisher.test(cbind(c(1,3),c(3,1)))
2
```

2:1 (Top Level) R Script

Console ~/R/RCourse/

Fisher's Exact Test for Count Data

```
data: cbind(c(1, 3), c(3, 1))
p-value = 0.4857
alternative hypothesis: true odds ratio is not equal to 1
95 percent confidence interval:
 0.001607888 4.722931239
sample estimates:
```