

Binomial tests

Chi square tests

P-values

Confidence Intervals

Null Hypothesis testing



Biostatistics Course 2024
Lecture 2
Tuesday, 9 July 2024
10:00am - 12:00pm

Example 1: Human Sex Ratio

Computing sex ratio of humans is one of the oldest applications of statistics

year	male	female
1629	5218	4683
1630	4858	4457
1631	4422	4102
1632	4994	4590
1633	5158	4839
1634	5035	4820
1635	5106	4928
1636	4917	4605
1637	4703	4457

:

Arbuthnott J (1711). An Argument for Divine Providence, taken from the Constant Regularity observed in the Births of both Sexes.

The screenshot shows the Sex Ratio application interface. The main window title is "sex_ratio.pzfx — Edited". The left sidebar contains a navigation tree with sections like "Data Tables", "Info", "Results", "Graphs", and "Layouts". A large arrow points from the "Graphs" section towards the central table area. The central area displays a table titled "Table format: Parts of whole" with columns labeled A through I. Column A is further divided into sub-headings: "year 1634", "total", and "Y". The data rows show the following values:

	A	B	C	D	E	F	G	H	I
1	boys	5035	484382						
2	girls	4820	453841						
3	Title								
4	Title								
5	Title								
6	Title								
7	Title								
8	Title								
9	Title								
10	Title								
11	Title								
12	Title								
13	Title								
14	Title								
15	Title								
16	Title								
17	Title								
18	Title								
19	Title								
20	Title								
21	Title								
22	Title								
23	Title								
24	Title								
25	Title								
26	Title								
27	Title								

The status bar at the bottom indicates "Row --, A: year 1634" and "Selected: Rows 1073741827, Columns 1".

Create New Analysis

Data to analyze

Table: Data 1

Type of analysis

Which analysis?

- ▼ Transform, Normalize...
 - Transform
 - Transform concentrations (X)
 - Normalize
 - Prune rows
 - Remove baseline and column math
 - Transpose X and Y
 - Fraction of Total
- XY analyses
- Column analyses
- Grouped analyses
- Contingency table analyses
- Survival analyses
- ▼ Parts of whole analyses
 - Fraction of Total
 - Compare observed distribution with ex...
- Multiple variable analyses
- Nested analyses
- Generate curve
- Simulate data
- Recently used

Analyze which data sets?

- A:year 1634
- B:total

Select All

Deselect All

?

Cancel

OK

Parameters: Compare observed distribution with expected

This analysis expects that each value in the data table is an actual number of events or items, and is not normalized in any way.

Data set to analyze

A: year 1634

Enter expected values as

Actual numbers of objects or events
 Percentages

With two rows, perform

Chi-square test (recommended)
 Chi-square test for goodness of fit

Expected distribution

Row	Outcome	Observed %	Expected %
1	boys	51.09	50
2	girls	48.91	50

Output

Method to calculate CI: Wilson/Brown (recommended)

Show this many significant digits (for everything except P values): 4

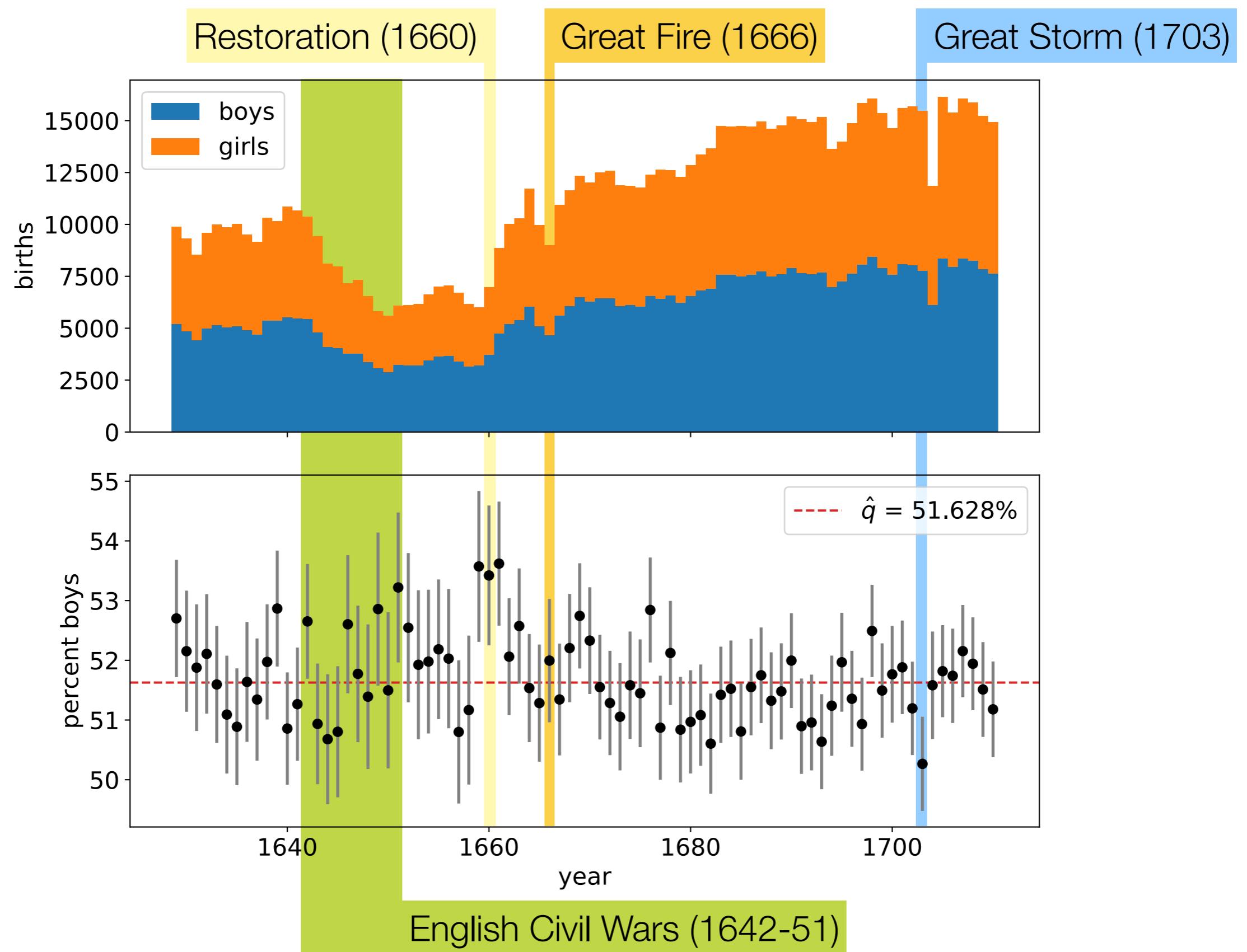
P value style: GP: 0.1234 (ns), 0.0332 (*), 0.0021 (**), 0.000... N= 6

?

Cancel

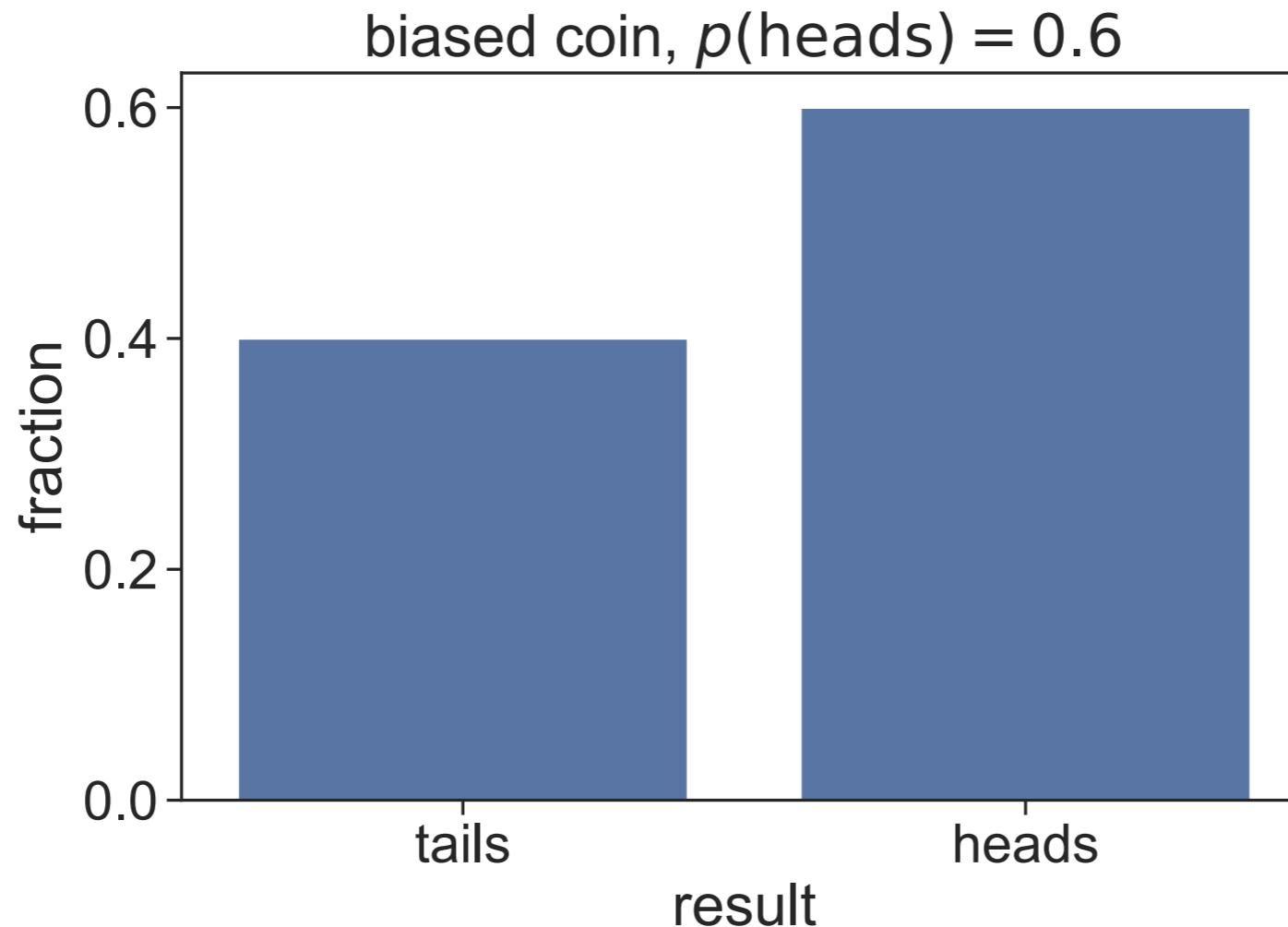
OK

Births in London, 1629-1710



Example 2: A biased coin

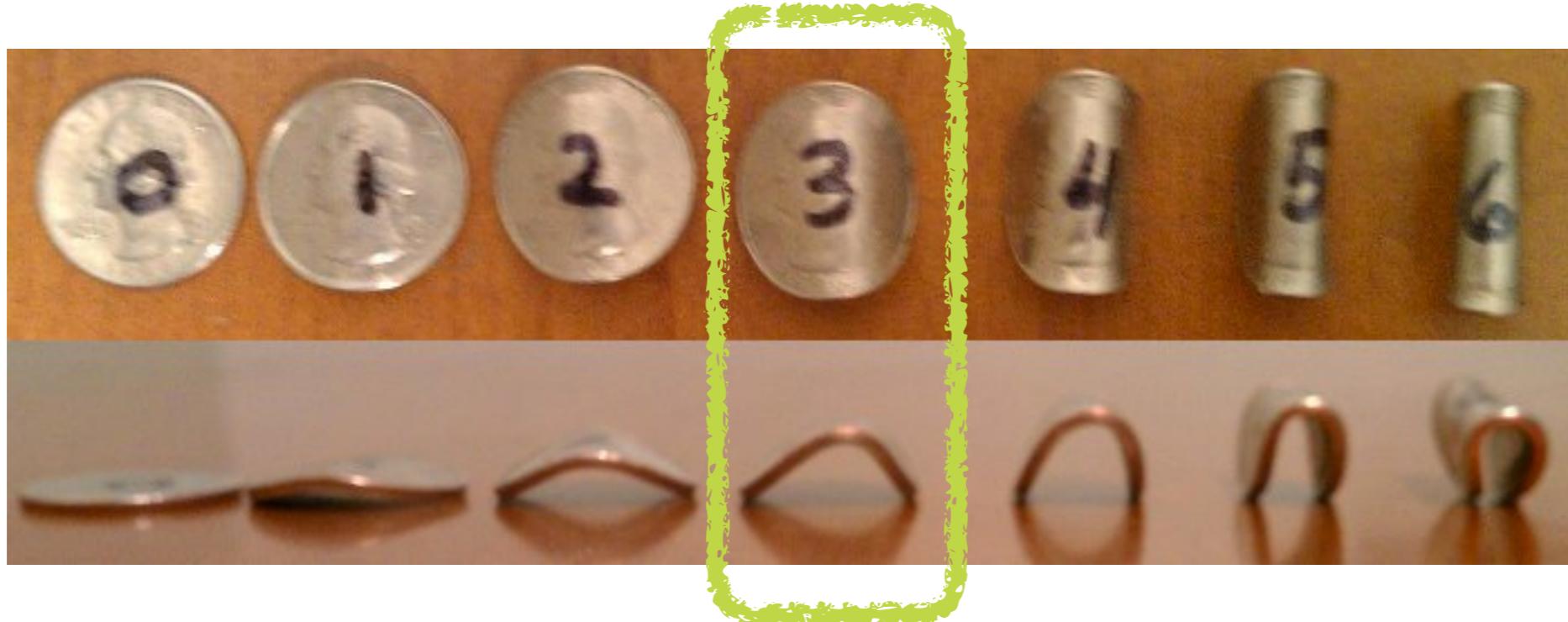
Biased coins are modeled using a Bernoulli distribution, which describes probabilities for a binary variable



Making a biased coin



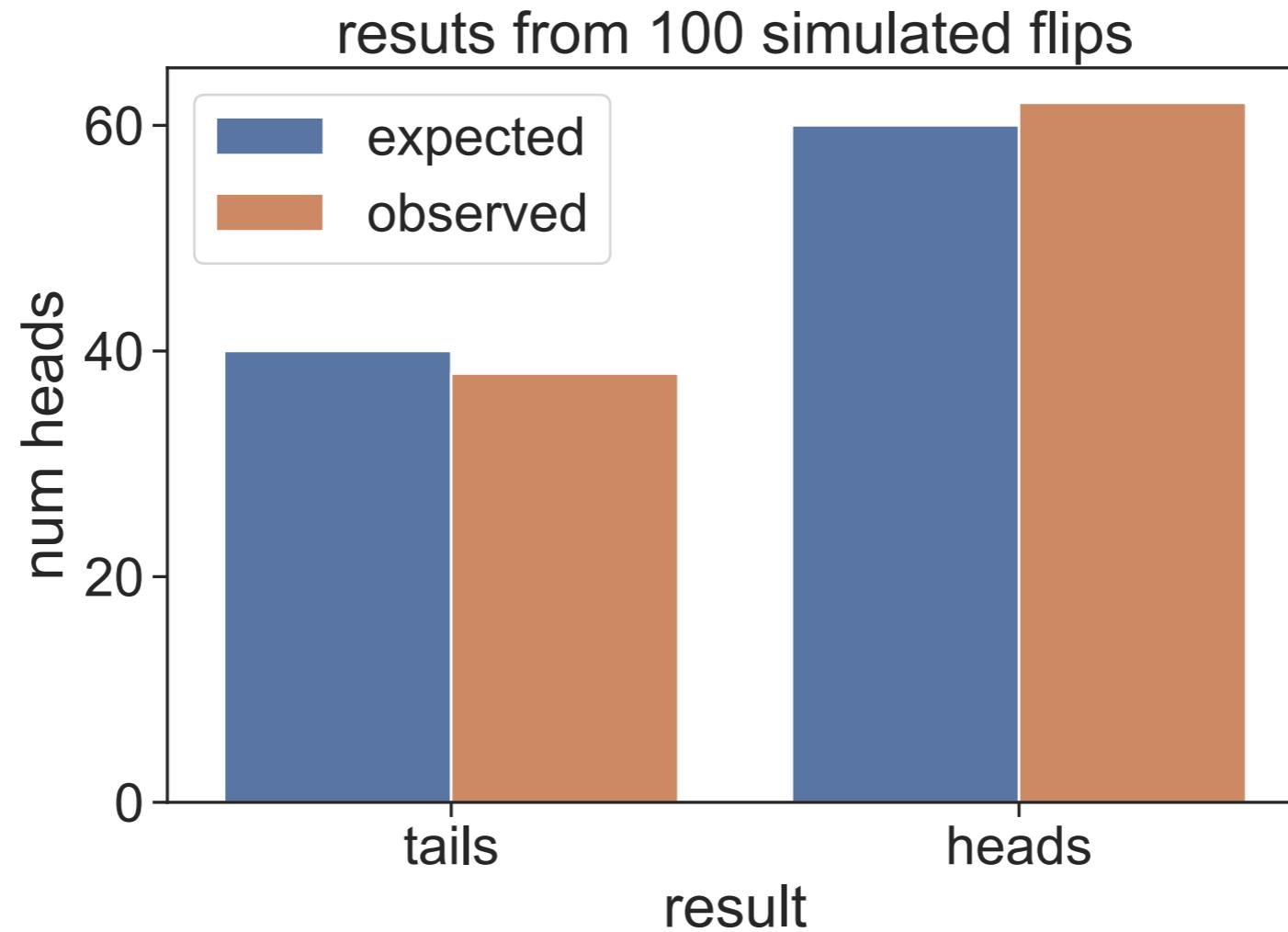
$p(\text{heads}) \approx 60\%$



Mike Izbicki (Claremont McKenna College)

<https://izbicki.me/blog/how-to-create-an-unfair-coin-and-prove-it-with-math.html>

The number of heads after 100 flips of the biased coin will resemble the underlying probabilities, but will not match exactly



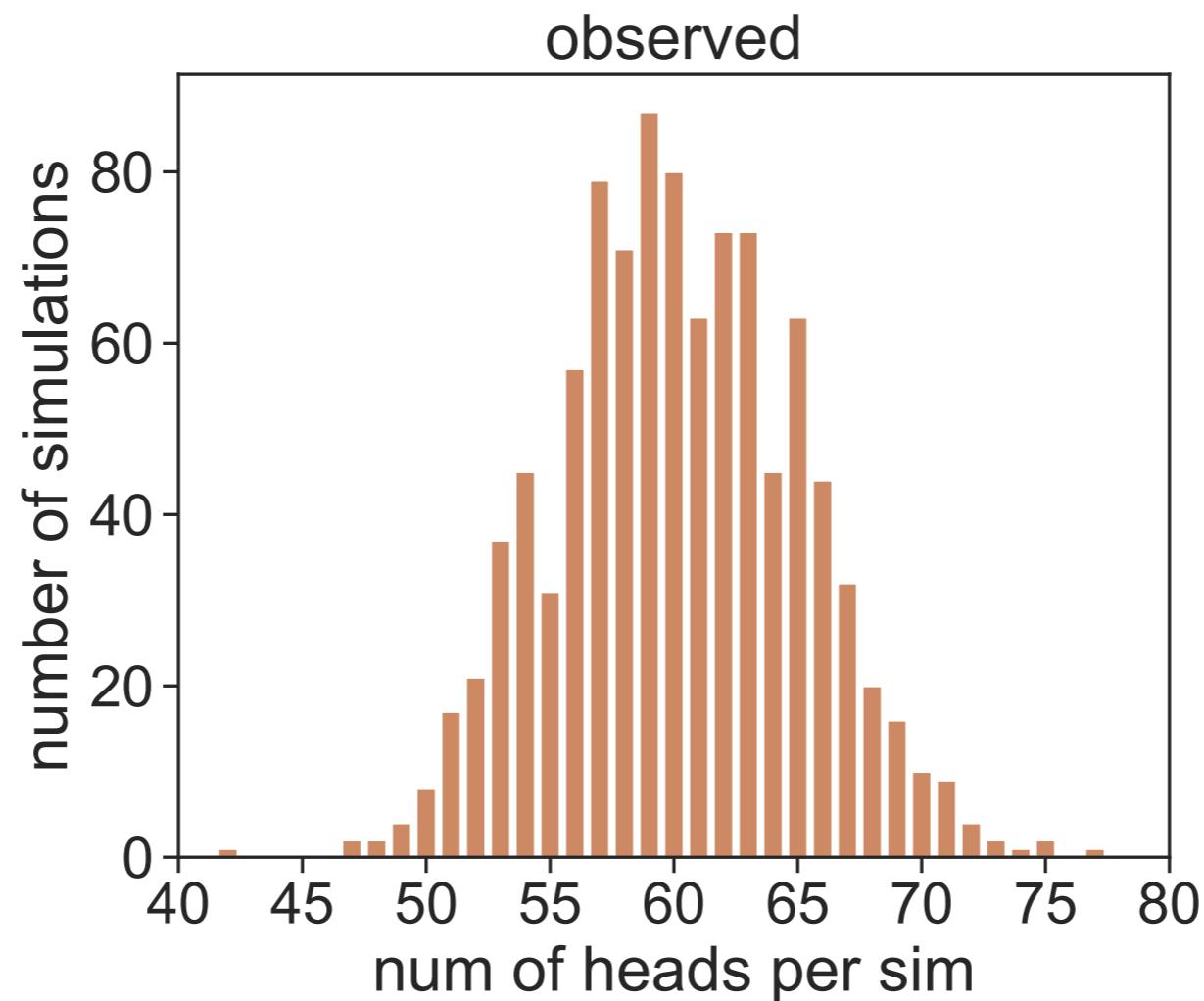
expected: 60 heads, 40 tails

observed: 62 heads, 38 tails

How much deviation from the expected values do we expect?

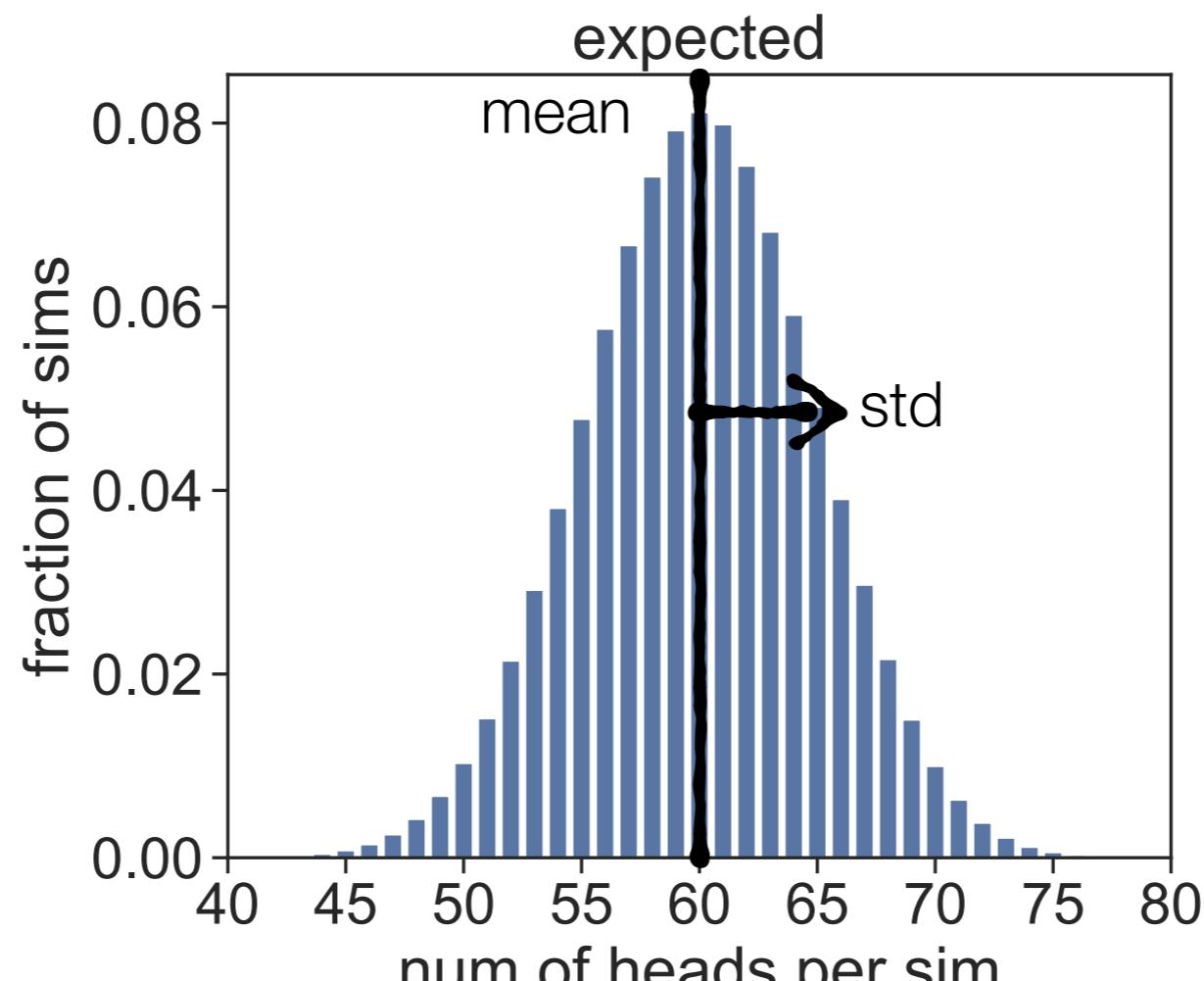
There is substantial variation across replicates. This is to be expected.

Results from 1000 simulations, 100 flips per simulation



The variation in the number of heads from replicate to replicate is described by a binomial distribution

Results from 1000 simulations, 100 flips per simulation



mean = 60

standard deviation (std) = 4.9

Can we determine whether or not a coin is biased by flipping it 100 times?

Suppose we flip a coin 100 times and observe **62 heads** (and 38 tails).

Null hypothesis: heads and tails are equally likely, i.e.

$$p(\text{heads}) = 50\%$$

Alternative hypothesis: heads and tails are not equally likely, i.e.

$$p(\text{heads}) \neq 50\%$$

Our observation (62 heads) may or may not allow us to reject the null hypothesis and thus accept the alternative hypothesis.

No amount of data, however, can cause us to accept the null hypothesis.

flips.pzfx — Edited

	A	B	C	D	E	F	G	H	I
	flips	Title							
1	heads	62							
2	tails	38							
3	Title								
4	Title								
5	Title								
6	Title								
7	Title								
8	Title								
9	Title								
10	Title								
11	Title								
12	Title								
13	Title								
14	Title								
15	Title								
16	Title								
17	Title								
18	Title								
19	Title								
20	Title								
21	Title								
22	Title								
23	Title								
24	Title								
25	Title								

Family

Data 1
Data 1

Data 1

Row 6, Column C

Create New Analysis

Data to analyze

Table: Data 1

Type of analysis

Which analysis?

- ▼ Transform, Normalize...
 - Transform
 - Transform concentrations (X)
 - Normalize
 - Prune rows
 - Remove baseline and column math
 - Transpose X and Y
 - Fraction of Total
- XY analyses
- Column analyses
- Grouped analyses
- Contingency table analyses
- Survival analyses
- ▼ Parts of whole analyses
 - Fraction of Total
 - Compare observed distribution with ex... 
- Multiple variable analyses
- Nested analyses
- Generate curve
- Simulate data
- Recently used

Analyze which data sets?

A:flips

When you analyze tables or graphs with more than one data set, use this space to select which data set(s) to analyze.

Select All

Deselect All

?

Cancel

OK 

Parameters: Compare observed distribution with expected

This analysis expects that each value in the data table is an actual number of events or items, and is not normalized in any way.

Data set to analyze

A: flips

Enter expected values as

- Percentages

With two rows, perform

- Chi-square test for goodness of fit

Expected distribution

Row	Outcome	Observed %	Expected %
1	heads	62	50
2	tails	38	50

Output

Method to calculate CI: Wilson/Brown (recommended)

Show this many significant digits (for everything except P values): ^
v

P value style: GP: 0.1234 (ns), 0.0332 (*), 0.0021 (**), 0.000... N= 6



Cancel



flips.pzfx — Edited

Search

▼ Data Tables >>

- Data 1
- + New Data Table...

▼ Info >>

- Project info 1
- + New Info...

▼ Results >>

- O vs. E of Data 1
- + New Analysis...

▼ Graphs >>

- Data 1
- + New Graph...

▼ Layouts >>

- + New Layout...

Family >>

Data 1

O vs. E

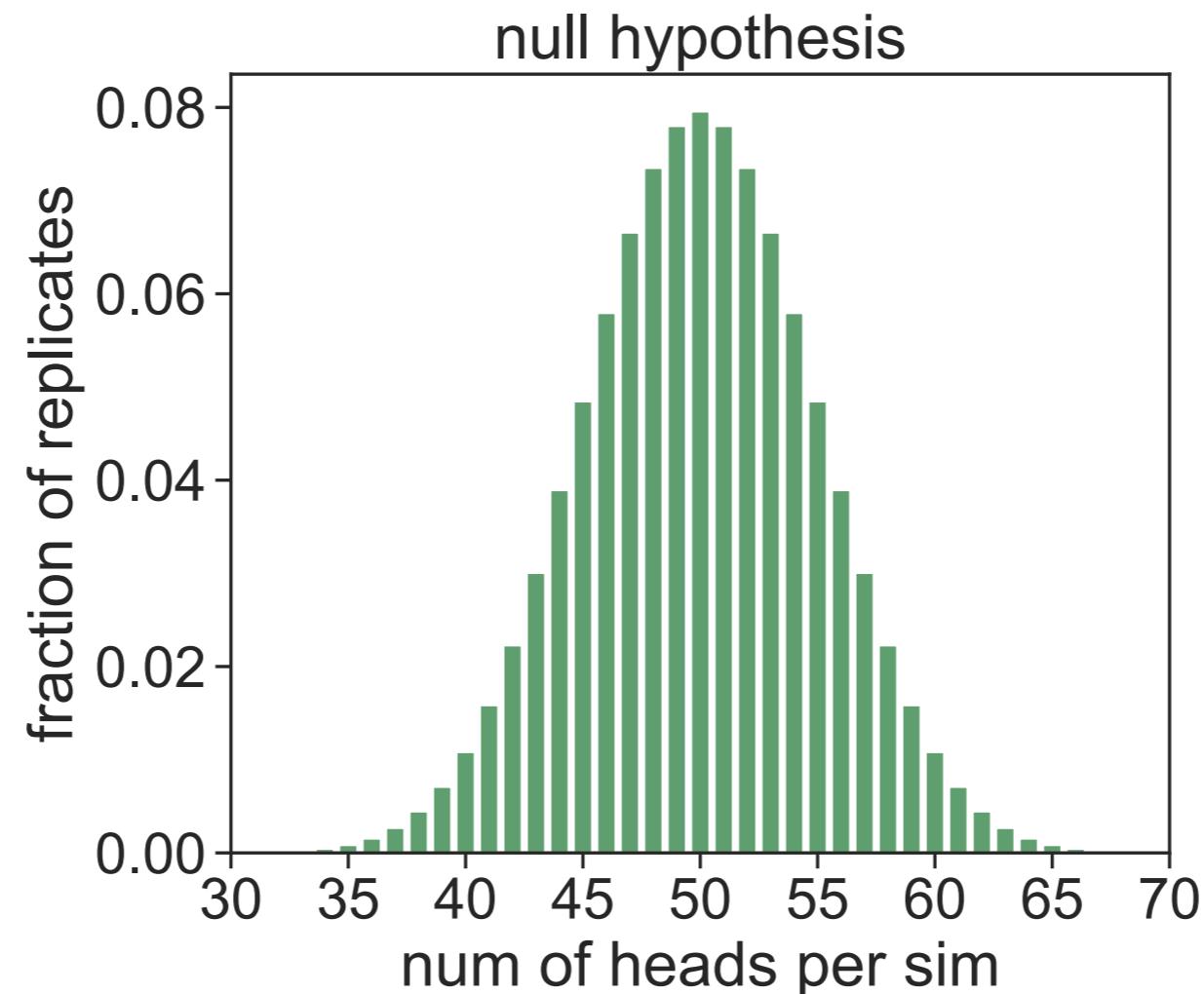
O vs. E

1	Table analyzed	Data 1					
2	Column analyzed	Column A					
3							
4	Binomial test						
5	P (one-tailed)	0.0105					
6	P (two-tailed)	0.0210					
7	P value summary	*					
8	Is discrepancy significant ($P < 0.05$)?	Yes					
9							
10	Outcome	Expected #	Observed #	Expected %	Observed %	95% CI of Observed %	
11	heads	50.00	62	50.00	62.00	52.21 to 70.90	
12	tails	50.00	38	50.00	38.00	29.10 to 47.79	
13	TOTAL	100.0	100.0	100.0	100.00		
14							
15							
16							
17							
18							
19							
20							
21							
22							
23							
24							
25							

O vs. E of Data 1

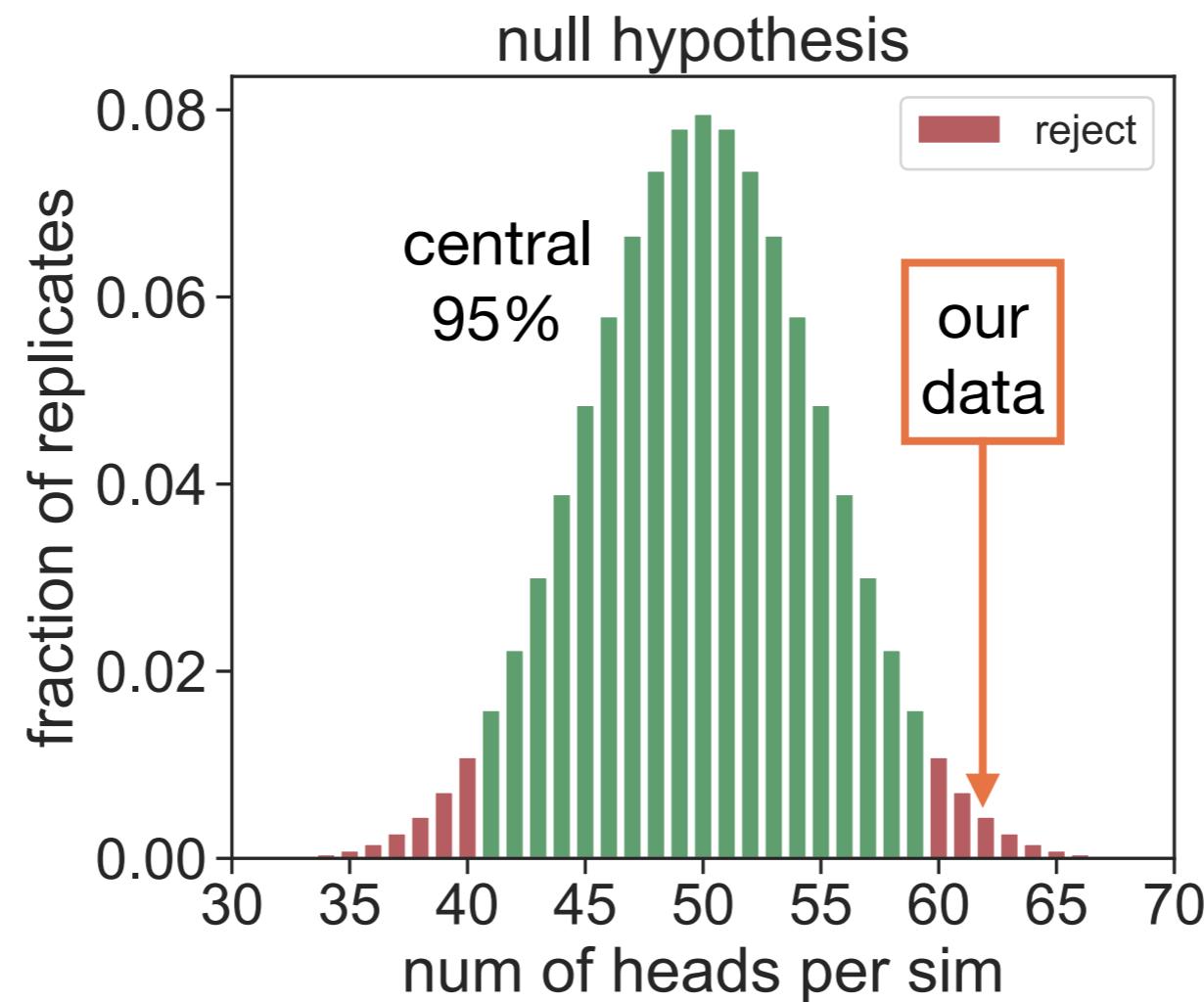
Row 1, Column A

The null hypothesis is assessed by where the data fall within the null distribution



We reject the null hypothesis if the data fall too far away from the bulk of the distribution

If the null hypothesis is true, data should fall within the green region 95% of the time, and within the red “reject” region 5% of the time.

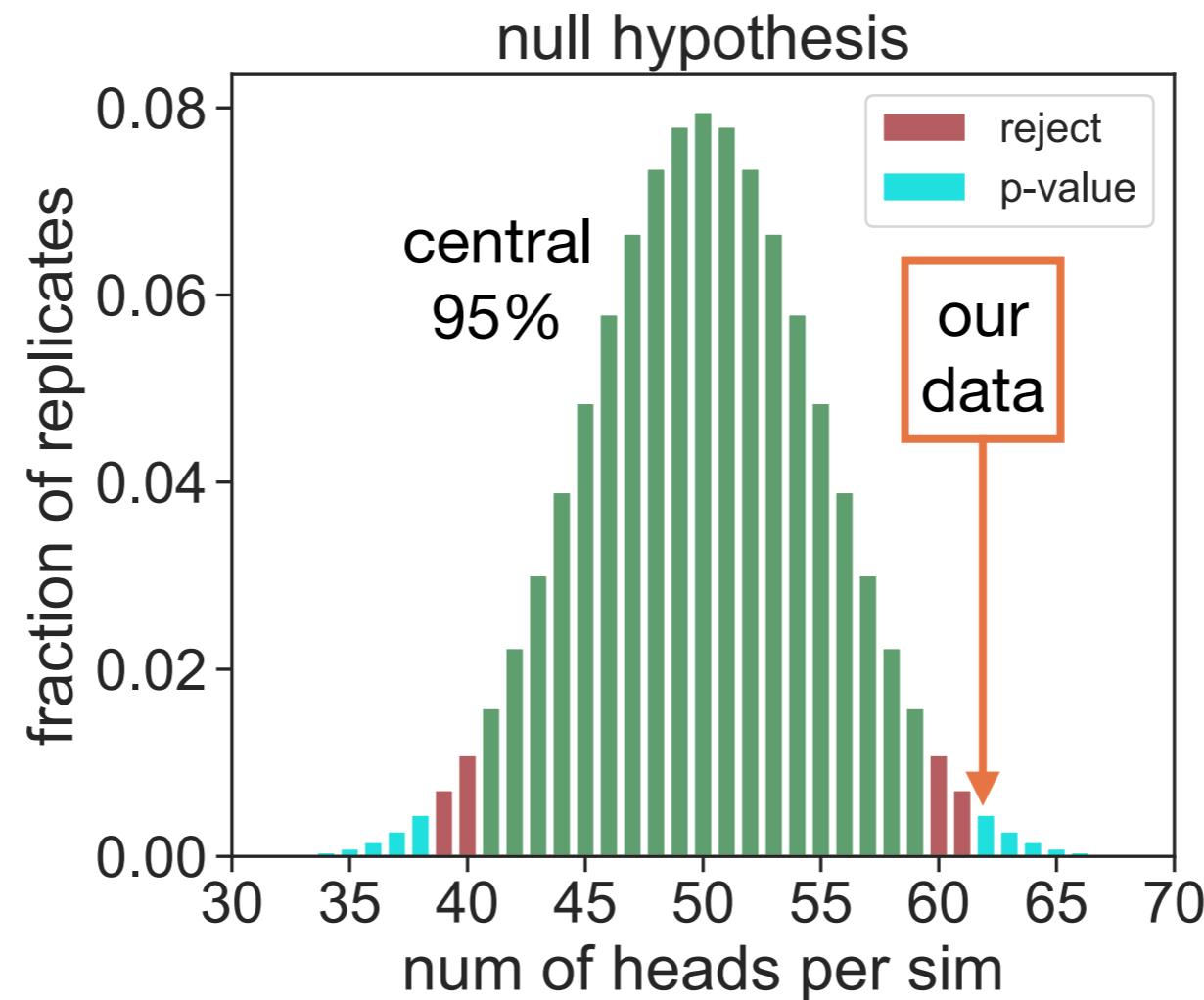


Our assumed dataset (62 heads) lies outside the central 95%.

We can therefore reject the null hypothesis with 95% confidence.

P-values quantify the probability of data being as or more extreme than the data in hand were the null hypothesis true.

The P-value threshold of 0.05 comes from adopting a confidence threshold of 95%.

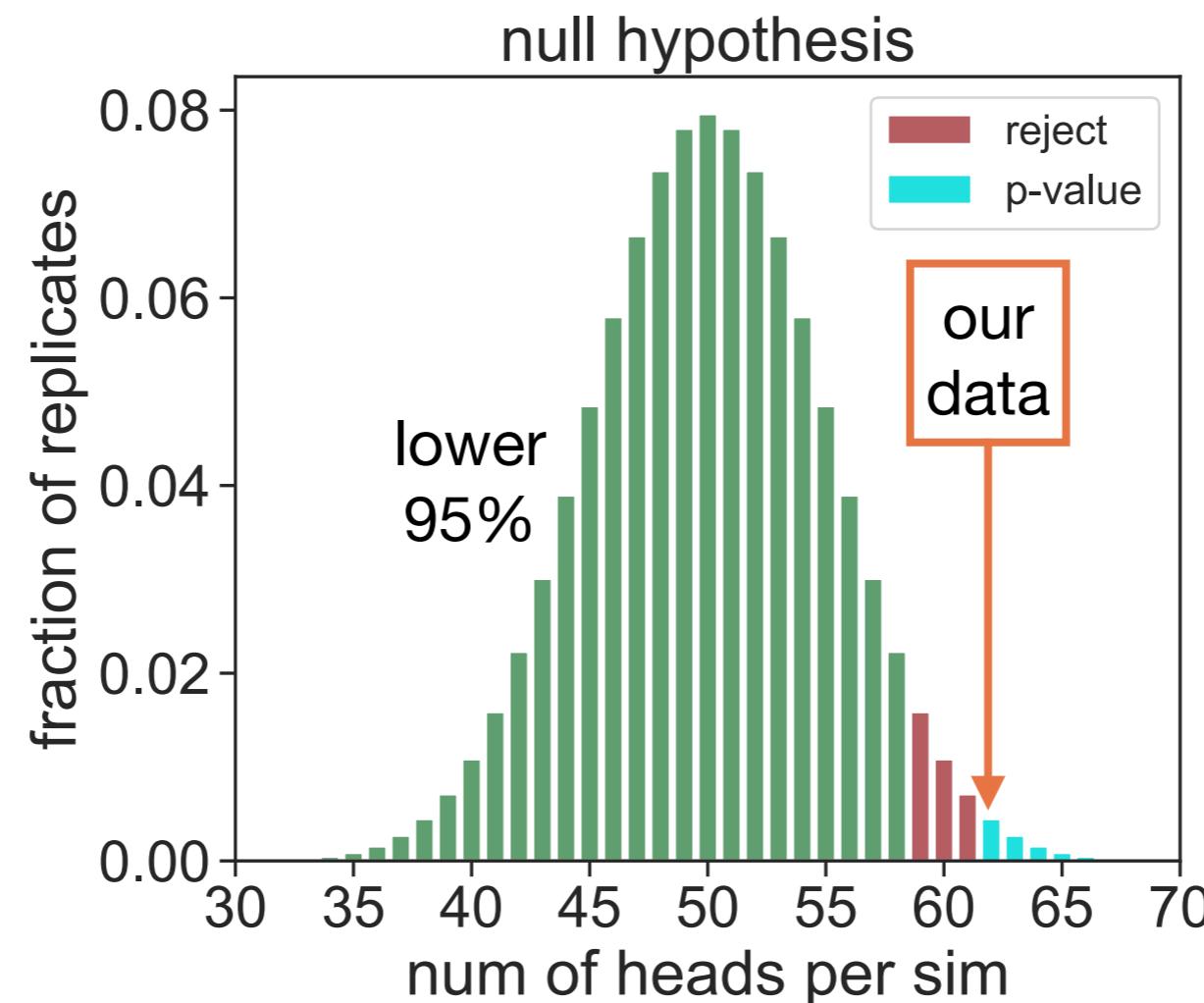


We find that **p=0.0210** for the two-sided test.

We therefore say that our result is “statistically significant”

P-values quantify the probability of data being as or more extreme than the data in hand were the null hypothesis true.

A one-sided hypothesis test only considers one side of the distribution.

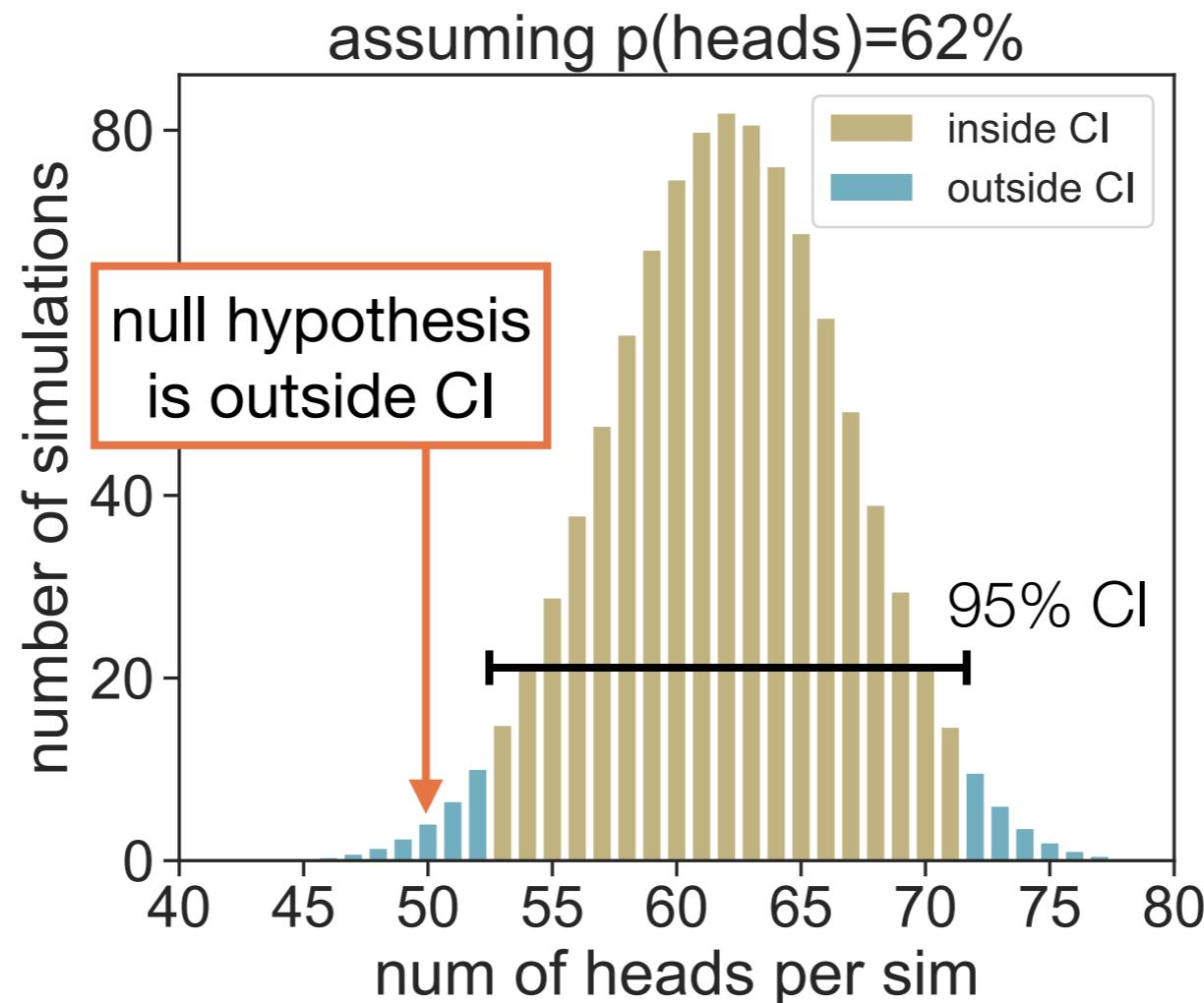


We find that **p=0.0105** for the one-sided test.

In general, two-sided tests are more conservative than one-sided tests.

Unless you have good reason to do otherwise, use two-sided tests.

Confidence intervals (CIs) are more informative than P-values



We conclude that $p(\text{heads})$ lies within $[52.5\%, 71.5\%]$ with 95% confidence.

We can reject the null hypothesis because it lies outside this confidence interval.

P-values have multiple pitfalls

- “Statistically significant” does not actually mean “significant” in the normal sense.
At best, it means “detectable”.
- P-values do not say how big an observed effect is.
- P-values do not say how important that observed effect is.
- P-values calculations rely on assumptions, and violation of any of those assumptions can render P-values meaningless.
- Perhaps most severe is the fact that P-values do not actually quantify you how likely or unlikely your null hypothesis is!

Why are Confidence Intervals better than P-values?

- Like a P-value, a CI communicates statistical significance (i.e. detectability).
- A CI also communicates effect size, as well as the uncertainty in that effect size.
- A 95% CI does not actually mean that the true value of a parameter lies within that interval with 95% probability. Still, this (extremely common) misinterpretation is largely benign compared to the misinterpretation of P-values.
- However, P-values are more commonly reported than confidence intervals.

The perils of null hypothesis testing

Summary of null hypothesis testing

Step 1: Specify a null hypothesis.

Step 2: Specify a confidence level (usually 95%)

Step 3: Identify the appropriate statistical test

Then:
evaluate on data



Result: P-value summarizing how unlikely the data is compared to null hypothesis expectations.

Perhaps most problematic is how easily P-values are misinterpreted.

Roughly speaking, P-values quantify how likely our data would be if the null hypothesis were true.

$$p(\text{data} \mid \text{null hypothesis})$$

P-values do not quantify the probability of the null hypothesis given our data. Unfortunately, this is the quantity that we actually care about.

$$p(\text{null hypothesis} \mid \text{data})$$

My opinion: the use of P-values to reject hypotheses is predicated on the base rate fallacy

By convention $P < 0.05$, then one rejects null hypothesis, supposedly because $p(\text{null hypothesis} \mid \text{data})$ is small.

For this to make sense, one has to accept the base rate fallacy, i.e.,

$$p(\text{data} \mid \text{null hypothesis}) \approx p(\text{null hypothesis} \mid \text{data})$$

Whether or not this is true in a specific case depends on the prior odds,

$$p(\text{null hypothesis}),$$

which Frequentist statistics refuses to consider.

The misinterpretation of P-values reflects the Frequentist / Bayesian divide

Frequentist statistics (a.k.a. classical statistics) focuses on likelihood:

$$p(\text{data} \mid \text{hypothesis}).$$

Iron Law of Frequentist Statistics:

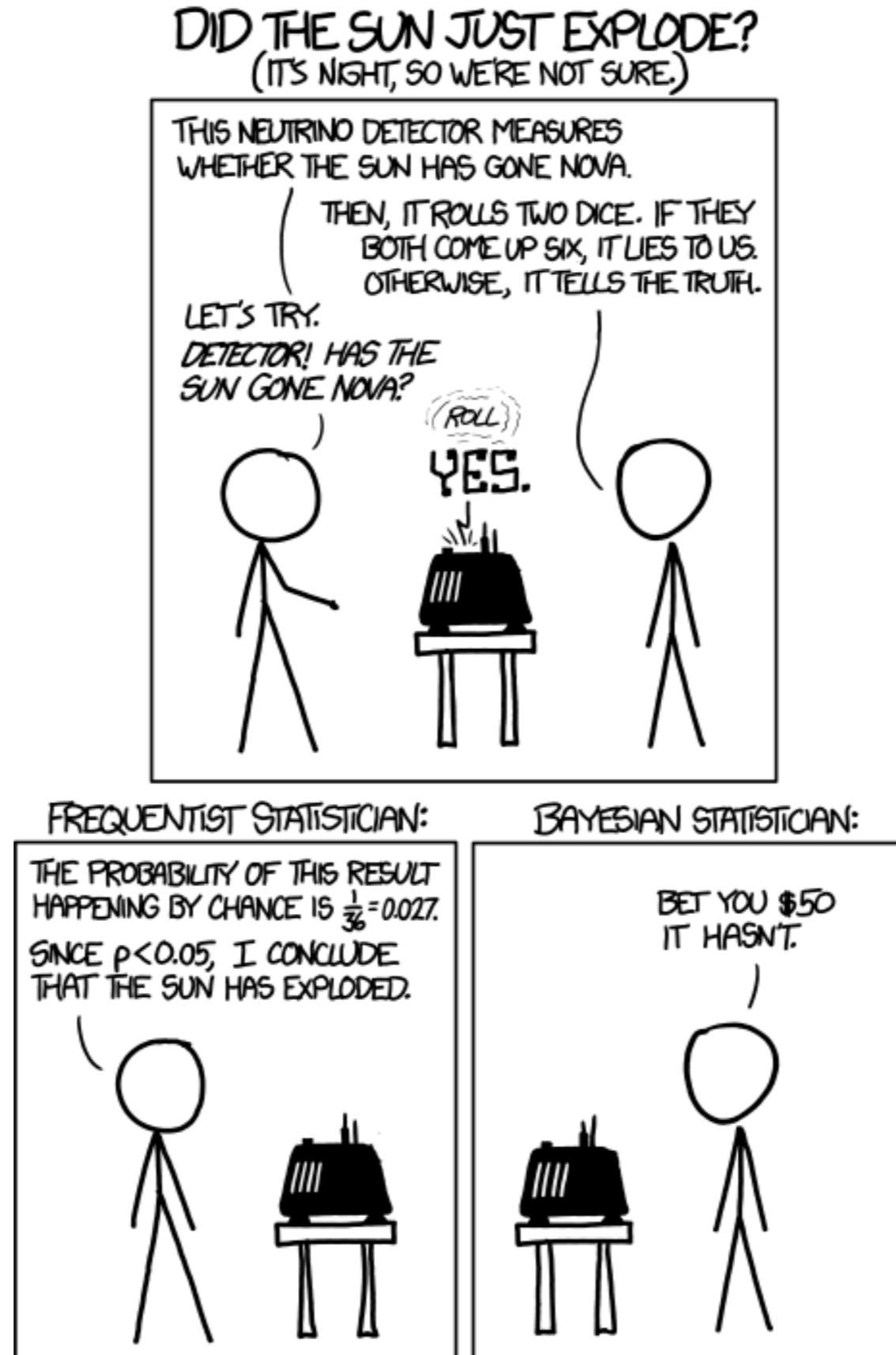
Never compute the probability of a hypothesis.

Bayesian statistics focuses on computing posterior probabilities:

$$p(\text{hypothesis} \mid \text{data}).$$

Example 3: Supernova detection machine

Exercise: supernova

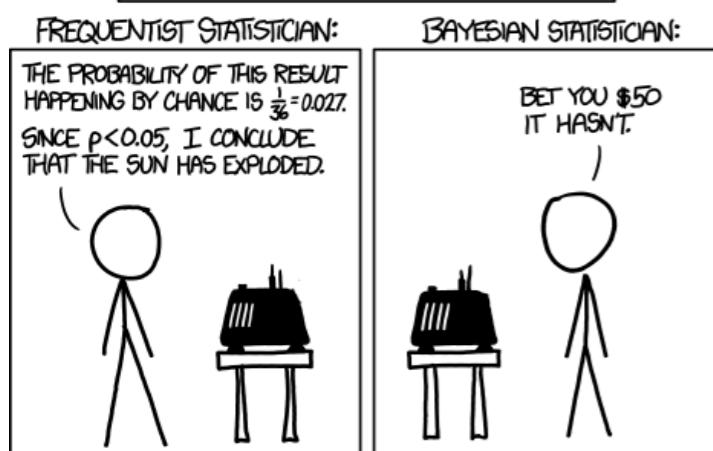
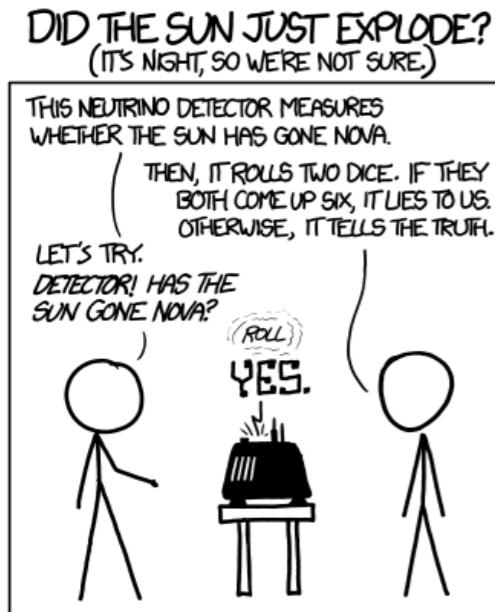


Exercise: supernova

Bayes's theorem (from yesterday):

$$\frac{p(\text{nova}^+ | \text{detector}^+)}{p(\text{nova}^- | \text{detector}^+)} = \frac{p(\text{detector}^+ | \text{nova}^+)}{p(\text{detector}^+ | \text{nova}^-)} \times \frac{p(\text{nova}^+)}{p(\text{nova}^-)}$$

$$\left[\frac{35/36}{1/36} = 35 \right]$$



If our prior belief is that a supernova is very unlikely, i.e.

$$\frac{p(\text{nova}^+)}{p(\text{nova}^-)} \ll \frac{1}{35},$$

then we still shouldn't believe the sun has gone nova.

<https://xkcd.com/1132/>

Even though, with a null hypothesis of nova^- ,

$$\text{P value} = p(\text{detector}^+ | \text{nova}^-) = \frac{1}{36} = 0.028 < 0.05$$

Example 4: Mendel's Peas

	Flower Colour	Plant Height	Seed Color	Seed Shape	Pod Colour	Pod Shape	Flower Position
Dominant Trait	Purple	Tall	Yellow	Round	Green	Inflated (full)	Axial
Recessive Trait	White	Short	Green	Wrinkled	Yellow	Constricted (flat)	Terminal
			3/4 Yellow	3/4 Round			
			1/4 Green	1/4 Wrinkled			

Chi square test (known proportions)

Example: Mendel's peas

	observed	expected proportion	expected counts
Round & yellow	315	9/16	312.75
Round & green	108	3/16	104.25
Angular & yellow	101	3/16	104.25
Angular & green	32	1/16	34.75
Total	556	16/16	556.00

Null Hypothesis:

observations in $K = 4$ different categories occur in the expected proportions

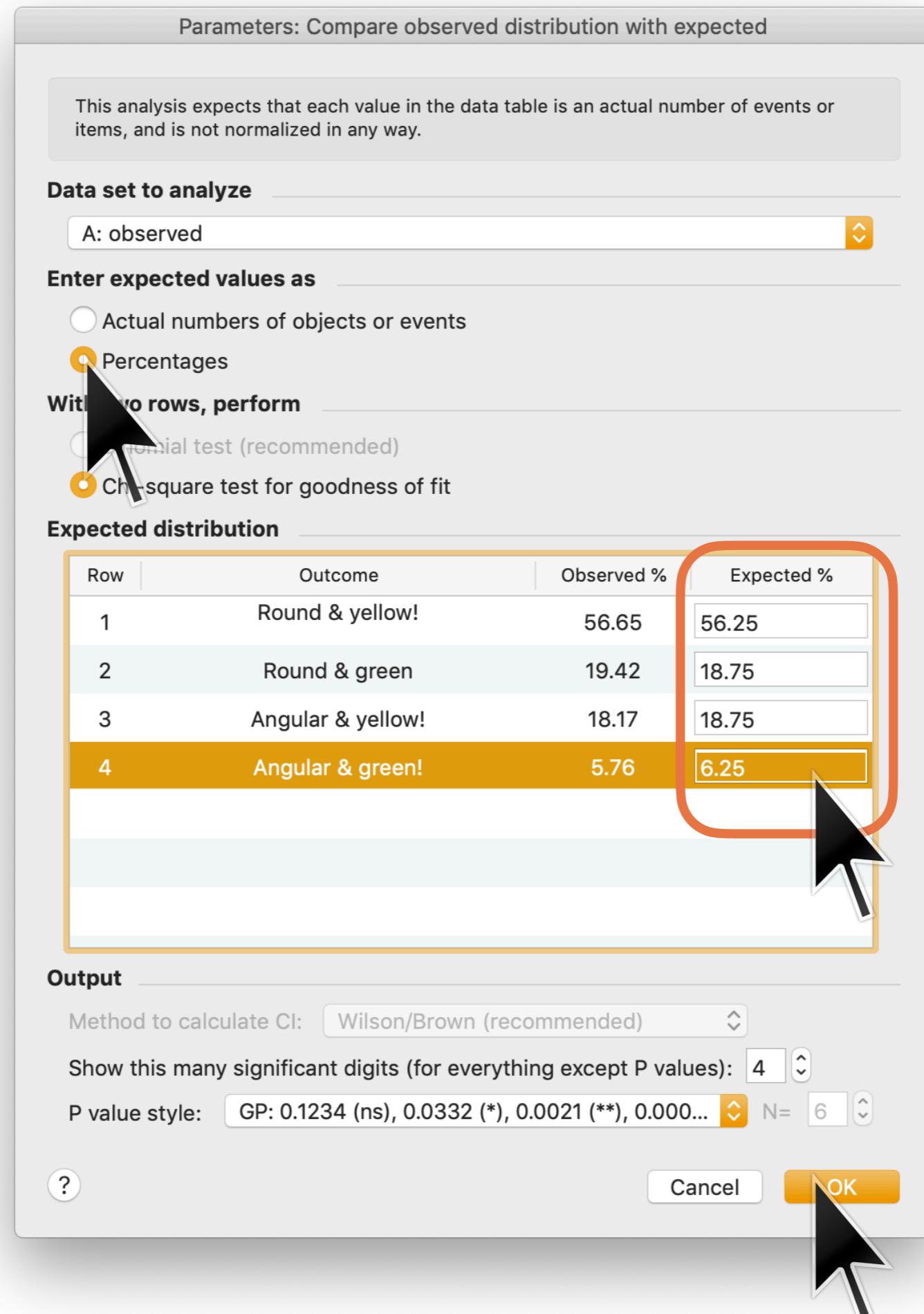
Data: number of observations in each category

$$\text{Statistic: } \chi^2 = \sum \frac{(\text{observed} - \text{expected})^2}{\text{expected}}$$

Null distribution: Chi square distribution with $K - 1 = 3$ degrees of freedom (DOF)

peas.pzfx

Table format: Parts of whole		A	B
		observed	Title
1	Round & yellow	315	
2	Round & green	108	
3	Angular & yellow	101	
4	Angular & green	32	
5	Title		
6	Title		
7	Title		
8	Title		
9	Title		
10	Title		
11	Title		
12	Title		
13	Title		
14	Title		



enter
manually

peas.pzfx — Edited

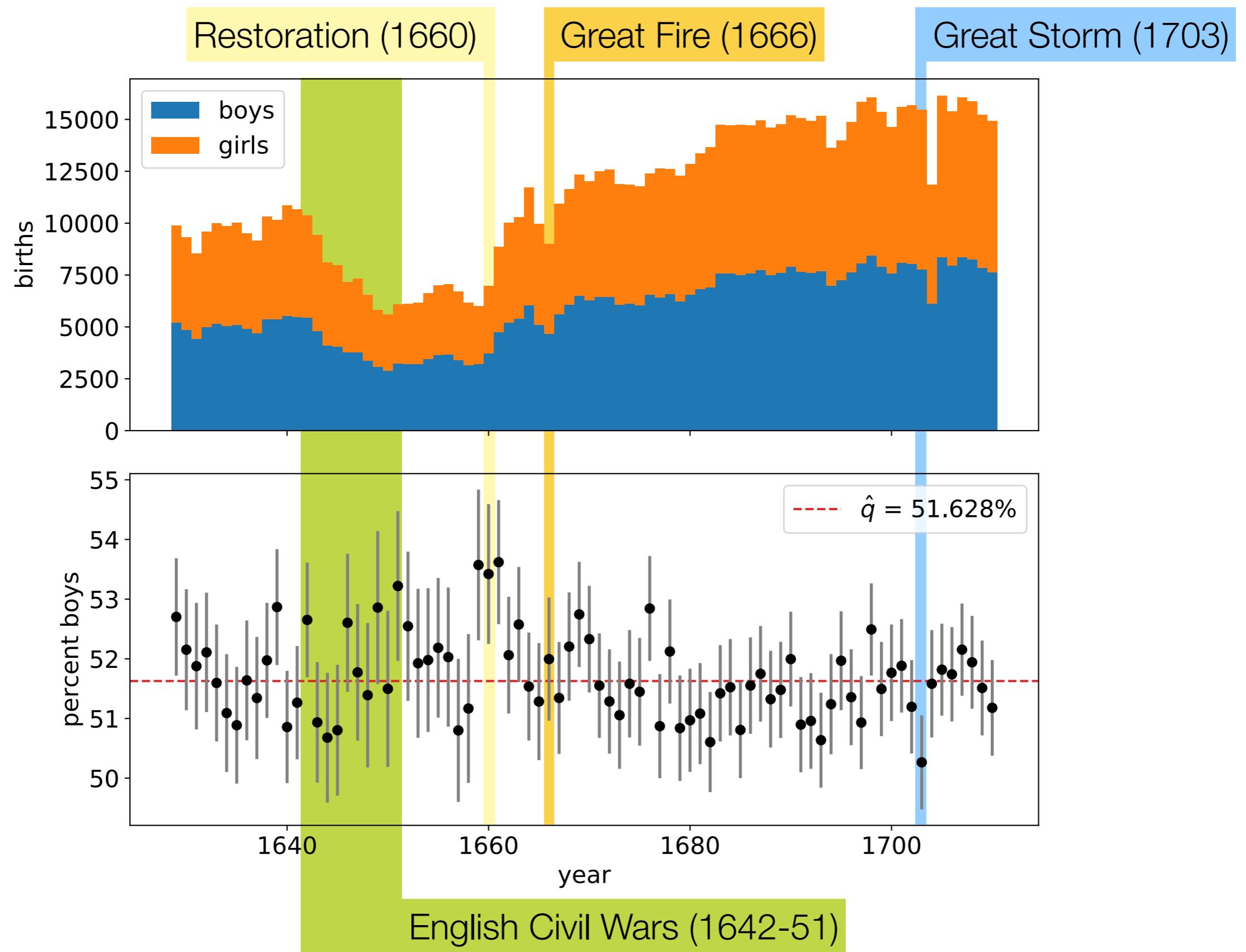
O vs. E

	1 Table analyzed	Data 1
	2 Column analyzed	Column A
4 Chi-square test		
5 Chi-square	0.4700	
6 DF	3	
7 P value (two-tailed)	0.9254	
8 P value summary	ns	
9 Is discrepancy significant ($P < 0.05$)?	No	
10		
11 Outcome	Expected #	Observed #
12 Round & yellow	312.8	315
13 Round & green	104.3	108
14 Angular & yellow	104.3	101
15 Angular & green	34.75	32
16 TOTAL	556.0	556.0
17		

data fits expectations

Example 4: Human sex ratio in London over time

Is it possible that the boy/girl ratio changes from year to year?



Chi square test (unknown proportions)

	sex	
	male	female
year		
1629	5218	4683
1630	4858	4457
1631	4422	4102
1632	4994	4590
1633	5158	4839
1634	5035	4820
1635	5106	4928
1636	4917	4605
1637	4703	4457
1638	5359	4952

Null Hypothesis:

Two multi-category variables A and B are independent, i.e.,
 $p(A, B) = p(A) \cdot p(B)$

Statistic:

$$\chi^2 = \sum \frac{(\text{observed} - \text{expected})^2}{\text{expected}}$$

Null distribution:

Chi square distribution with
DOF = $nm - m - n + 1$
where

m = number of possible values for A
 n = number of possible values or B

arbuthnot.pzfx

Table format: **Contingency**

		Outcome A	Outcome B	Outcome C	Outcome D	Outcome E	Outcome F	Outcome G	Outcome H
		boys	girls	Title	Title	Title	Title	Title	Title
		Y	Y	Y	Y	Y	Y	Y	Y
1	1629		5218	4683					
2	1630		4858	4457					
3	1631		4422	4102					
4	1632		4994	4590					
5	1633		5158	4839					
6	1634		5035	4820					
7	1635		5106	4928					
8	1636		4917	4605					
9	1637		4703	4457					
10	1638		5359	4952					
11	1639		5366	4784					
12	1640		5518	5332					
13	1641		5470	5200					
14	1642		5460	4910					
15	1643		4793	4617					
16	1644		4107	3997					
17	1645		4047	3919					
18	1646		3768	3395					
19	1647		3796	3536					
20	1648		3363	3181					
21	1649		3079	2746					
22	1650		2890	2722					
23	1651		3231	2840					
24	1652		3220	2908					
25	1653		3196	2959					
26	1654		3441	3179					
27	1655		3655	3349					
28	1656		3668	3382					
29	1657		3396	3289					

Create New Analysis

Data to analyze

Table: arbuthnot

Type of analysis

Which analysis?

- ▼ Transform, Normalize...
 - Transform
 - Transform concentrations (X)
 - Normalize
 - Prune rows
 - Remove baseline and column math
 - Transpose X and Y
 - Fraction of Total
- XY analyses
- Column analyses
- Grouped analyses
- ▼ Contingency table analyses
 - Chi-square (and Fisher's exact) test**
 - Row means with SD or SEM
 - Fraction of Total
- Survival analyses
- Parts of whole analyses
- Multiple variable analyses
- Nested analyses
- Generate curve
- Simulate data
- Recently used

Analyze which data sets?

- A:boys
- B:girls

Select All

Deselect All

?

Cancel

OK

Parameters: Chi-square (and Fisher's exact) test

Main Calculations Options

Effect sizes to report

Relative Risk

Used for prospective and experimental studies

Difference between proportions (attributable risk) and NNT

Used for prospective and experimental studies

Odds ratio

Used for retrospective case-control studies

Sensitivity, specificity and predictive values

Used for diagnostic tests

Method to compute the P value

Fisher's exact test

Yates' continuity corrected chi-square test

Chi-square test

Chi-square test for trend

Looking for the z test to compare proportions? Choose the chi-square test (with or without the Yates' correction). The chi-square and z tests are equivalent.



Cancel

OK

