



# HOW TO TEST A HYPOTHESIS?

Sandeep Soni

09/17/2024

# QUESTION FOR THE DAY

“How to statistically test claims using text?”

# HYPOTHESES

# HYPOTHESES

- Many claims can be framed as hypotheses

# HYPOTHESES

- Many claims can be framed as hypotheses

## Examples

Gender bias is decreasing in books

Reframing a tweet can increase its retweet rate

Institutional mistrust is predictive of hate speech

Classifier A is better than classifier B

# NULL HYPOTHESIS

# NULL HYPOTHESIS

- Null hypothesis is a claim that is assumed to be true

# NULL HYPOTHESIS

- Null hypothesis is a claim that is assumed to be true

Hypothesis (H)	$H_0$
Gender bias is decreasing over time in books	Gender bias remains the <b>same</b> over time in books
Reframing a tweet can increase its retweet rate	Reframing a tweet has <b>no</b> effect on retweet rate
Institutional mistrust is predictive of hate speech	Institutional mistrust <b>is not</b> predictive of hate speech
Classifier A is better than classifier B	Classifier A is the <b>same</b> as classifier B

# HYPOTHESIS TESTING

# HYPOTHESIS TESTING

- $H_0$  helps us get an expected result

# HYPOTHESIS TESTING

- $H_0$  helps us get an expected result
- Testing if hypothesis  $H$  holds ==> if  $H_0$  is assumed to be true, how likely does the observed data match our expectation?

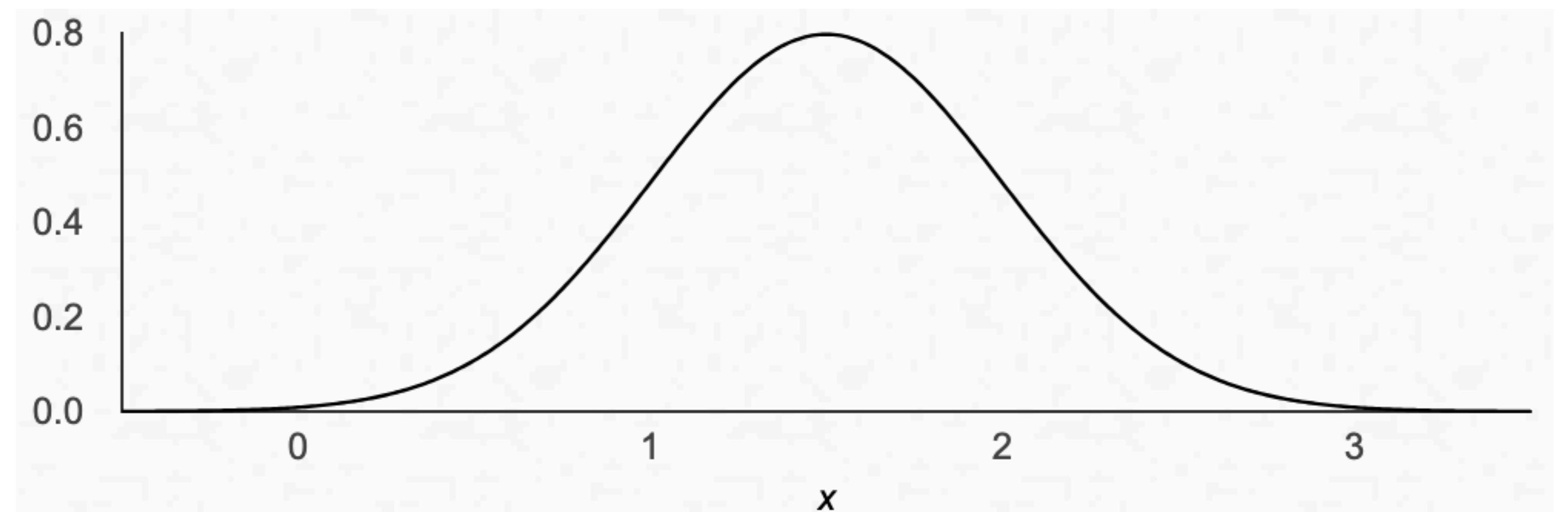
# HYPOTHESIS TESTING

# HYPOTHESIS TESTING

According to the null hypothesis, I expect the statistic to be normally distributed with some mean and variance

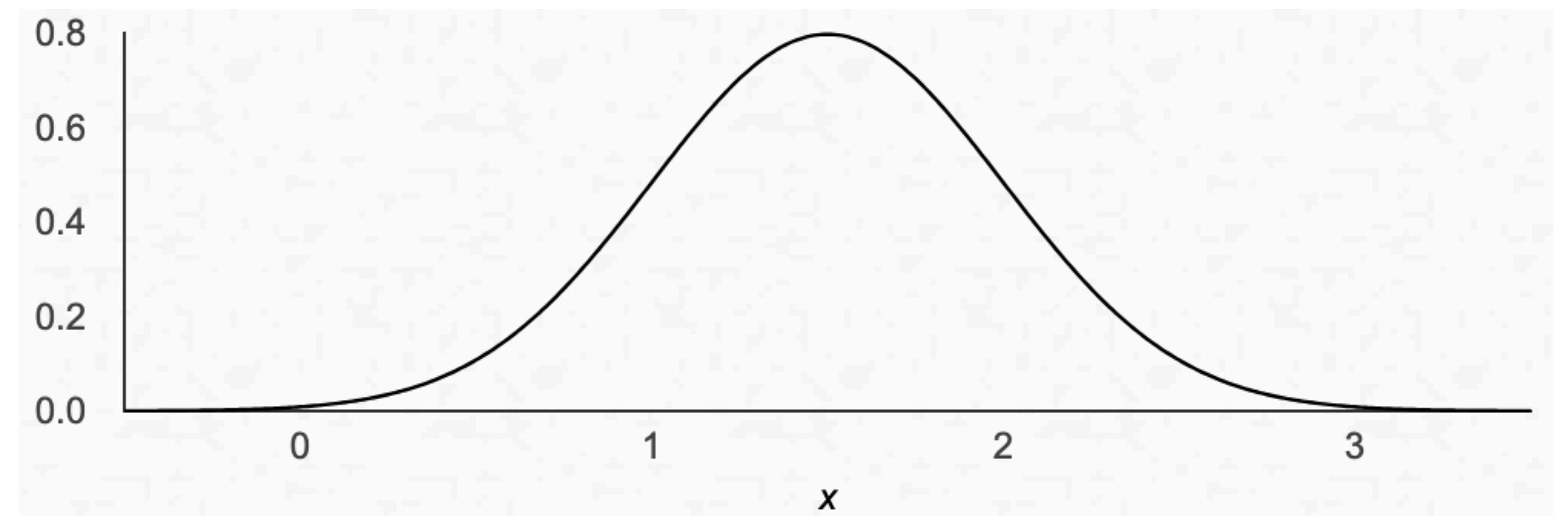
# HYPOTHESIS TESTING

According to the null hypothesis, I expect the statistic to be normally distributed with some mean and variance



# HYPOTHESIS TESTING

According to the null hypothesis, I expect the statistic to be normally distributed with some mean and variance



This sets up expectation about the shape and position of the distribution

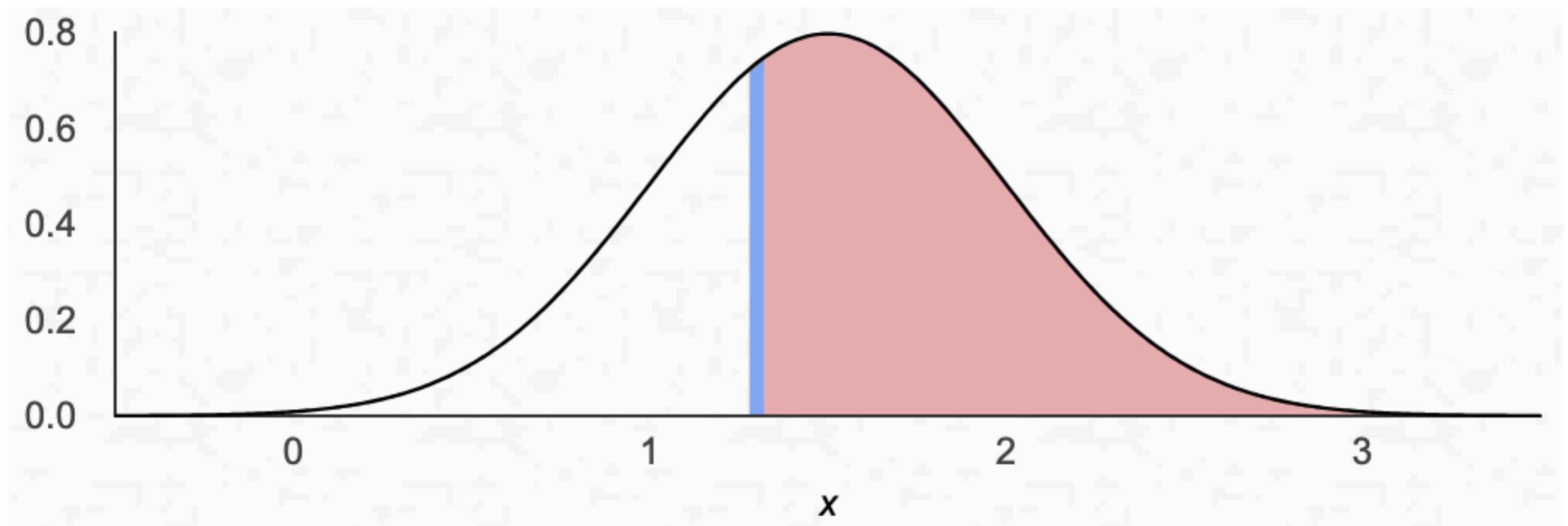
# HYPOTHESIS TESTING

# HYPOTHESIS TESTING

If we observe  
 $x$  to be 1.3,  
how likely are  
we to see that  
if the null  
hypothesis  
holds?

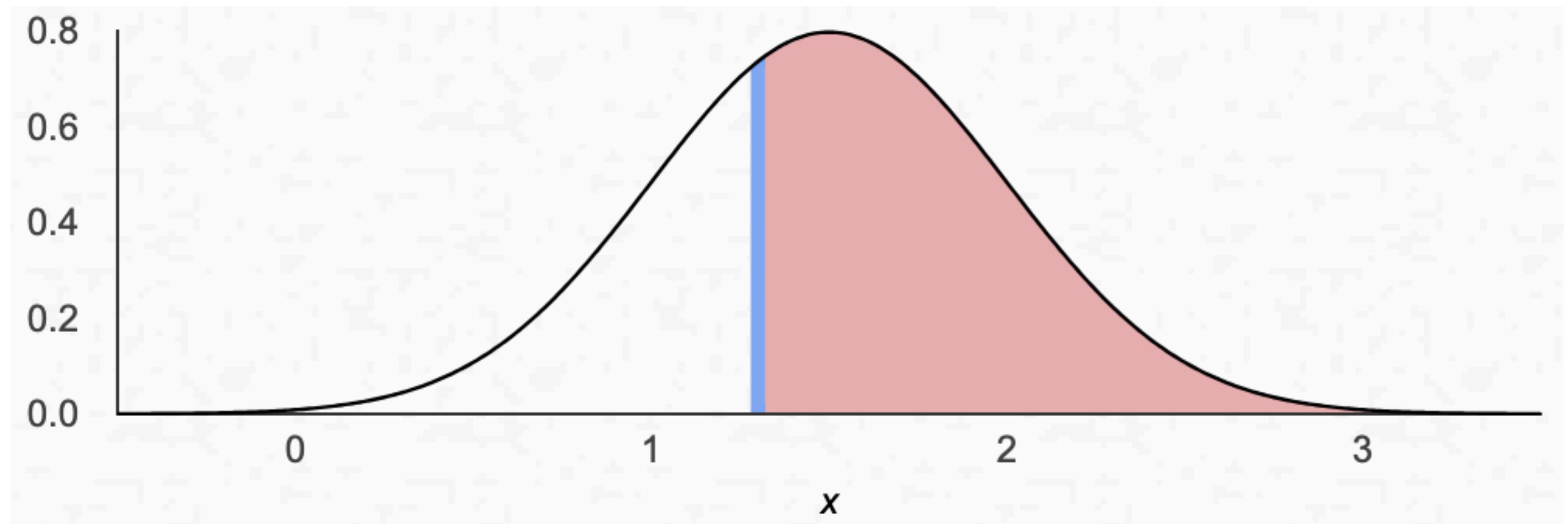
# HYPOTHESIS TESTING

If we observe  $x$  to be 1.3, how likely are we to see that if the null hypothesis holds?



# HYPOTHESIS TESTING

If we observe  $x$  to be 1.3, how likely are we to see that if the null hypothesis holds?



We can quantify this by calculating the probability of the shaded area

# HYPOTHESIS TESTING

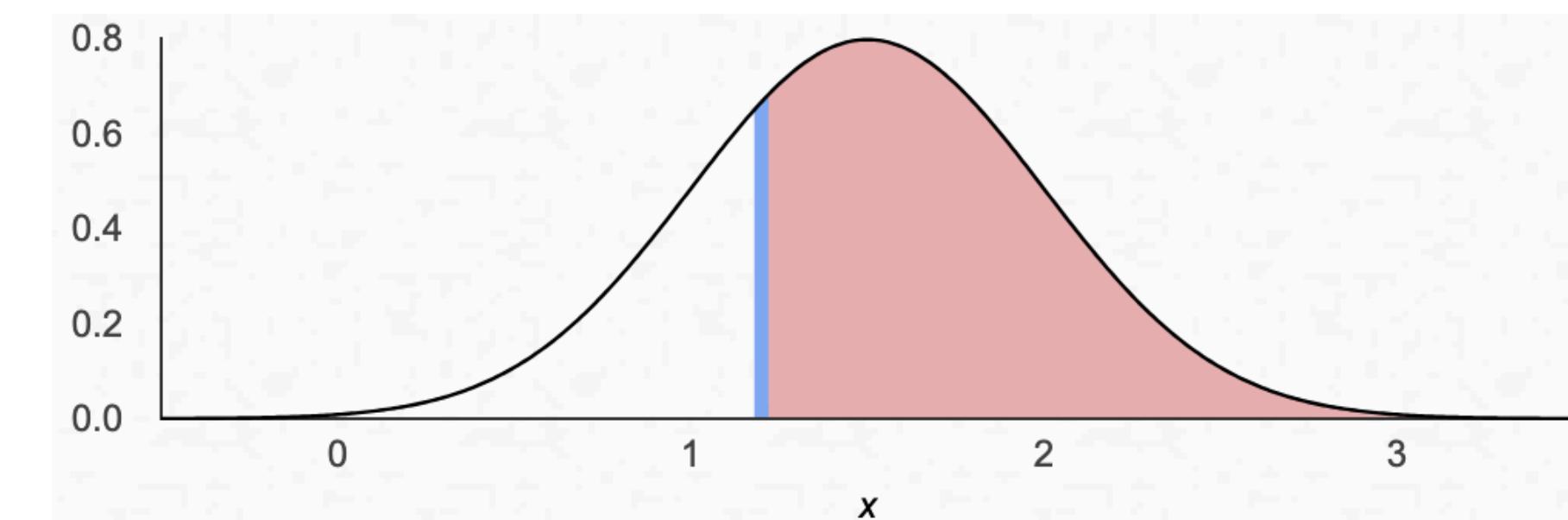
# HYPOTHESIS TESTING

Which  
observation  
would you say  
is surprising  
if the null  
distribution  
holds?

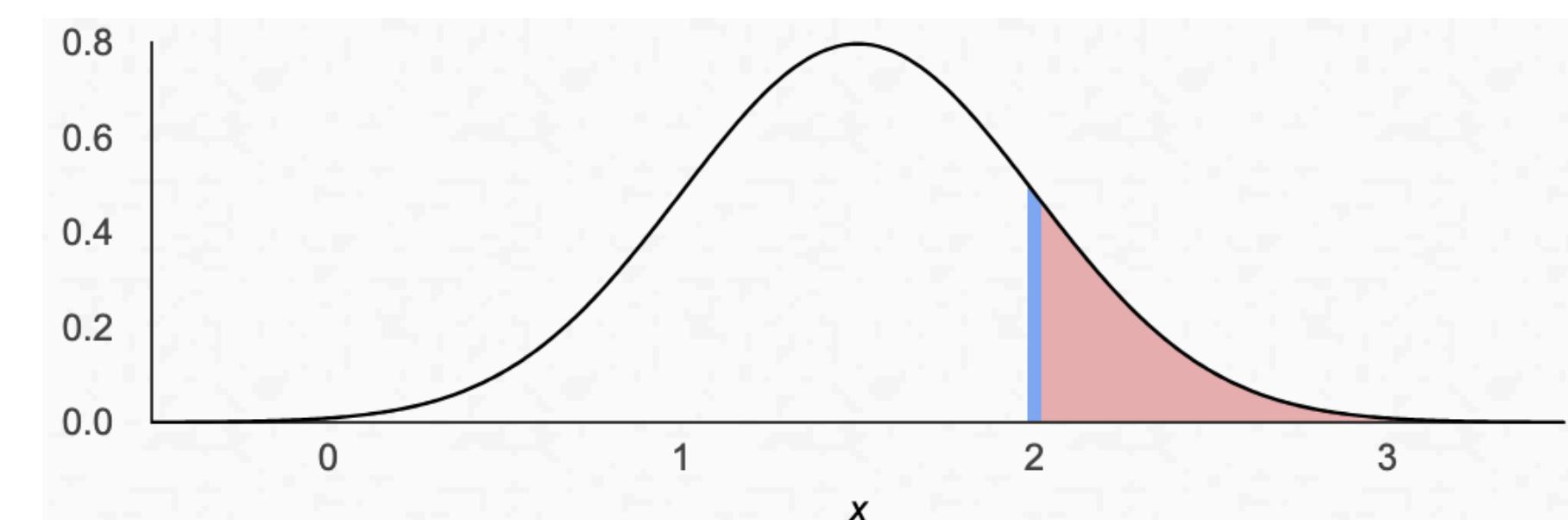
# HYPOTHESIS TESTING

Which observation would you say is surprising if the null distribution holds?

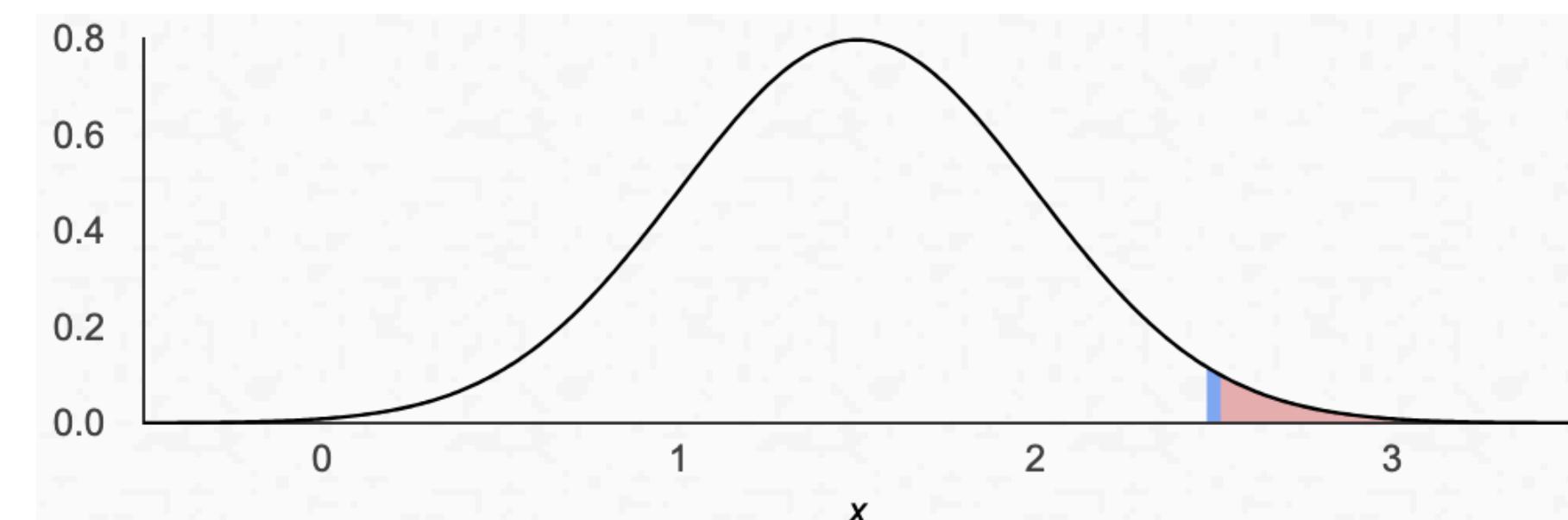
$x=1.25$



$x=2$



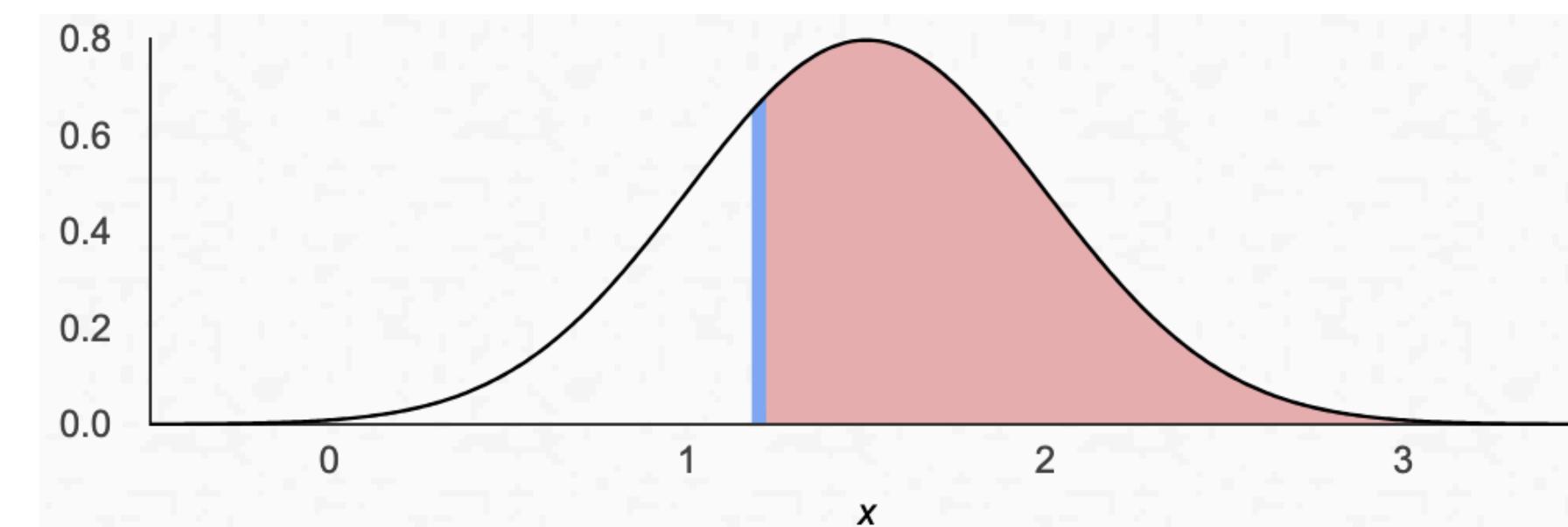
$x=2.5$



# HYPOTHESIS TESTING

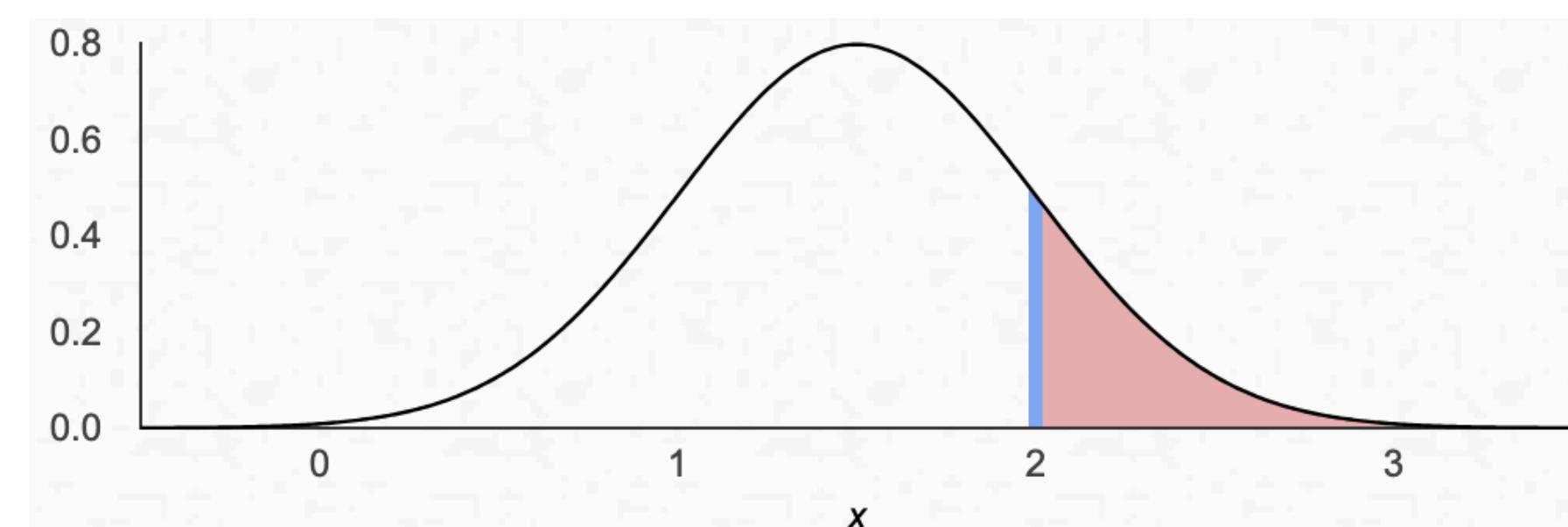
Which observation would you say is surprising if the null distribution holds?

x=1.25



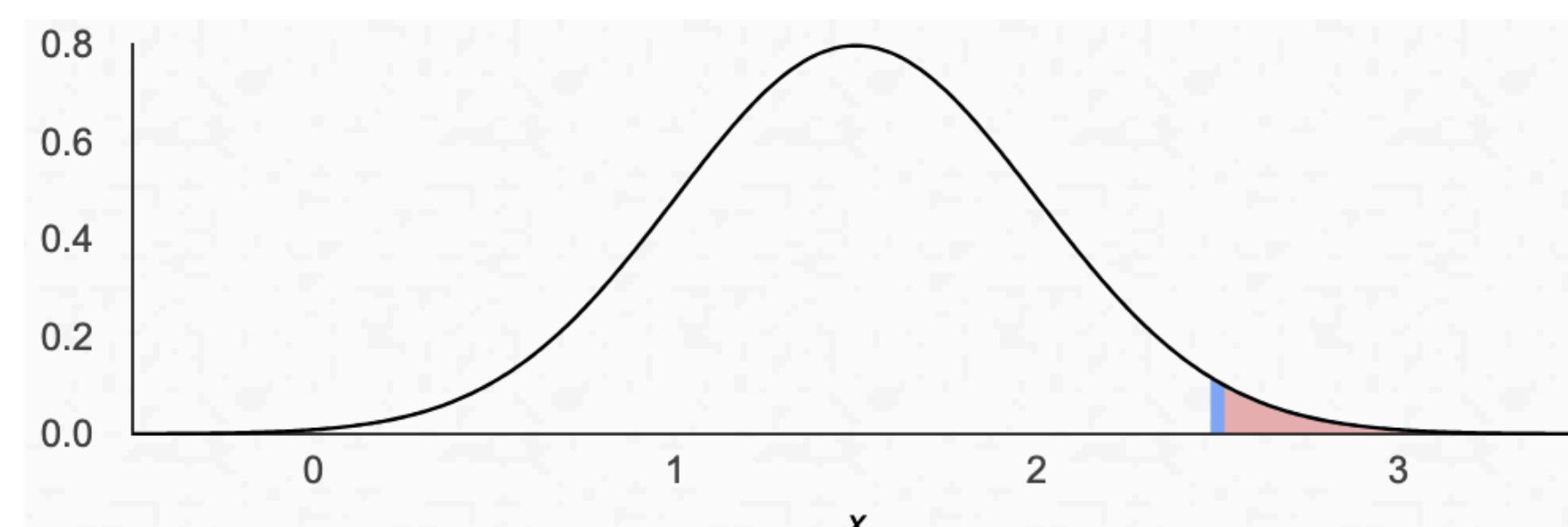
p=0.69

x=2



p=0.16

x=2.5



p=0.02

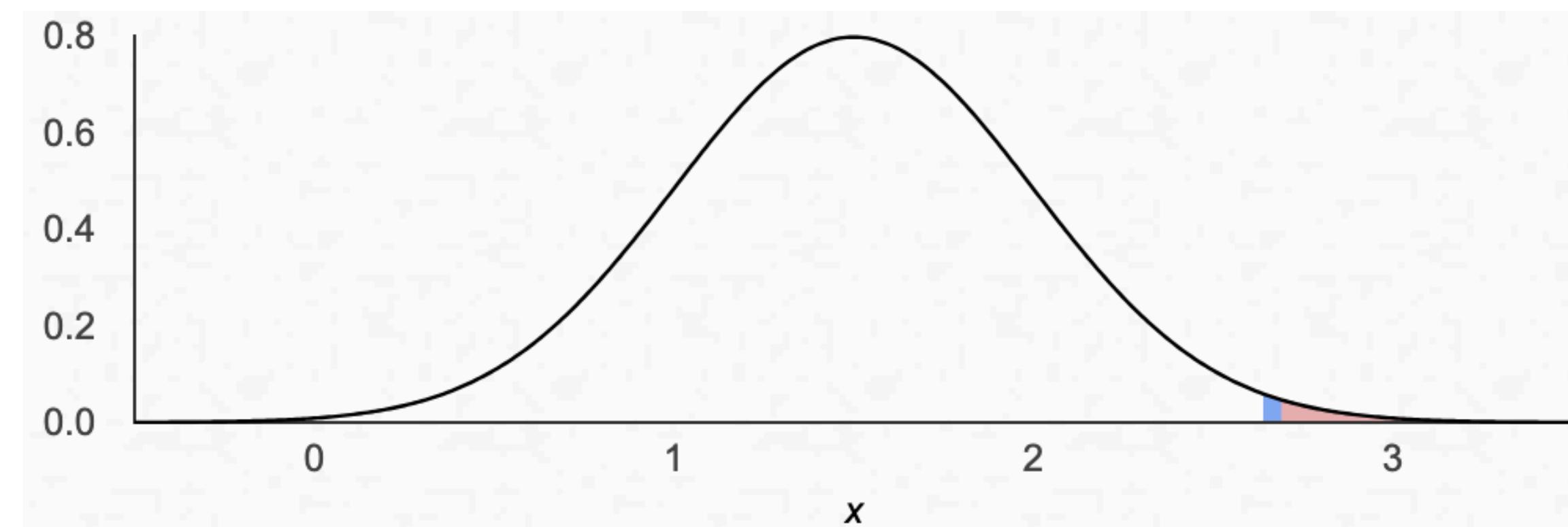
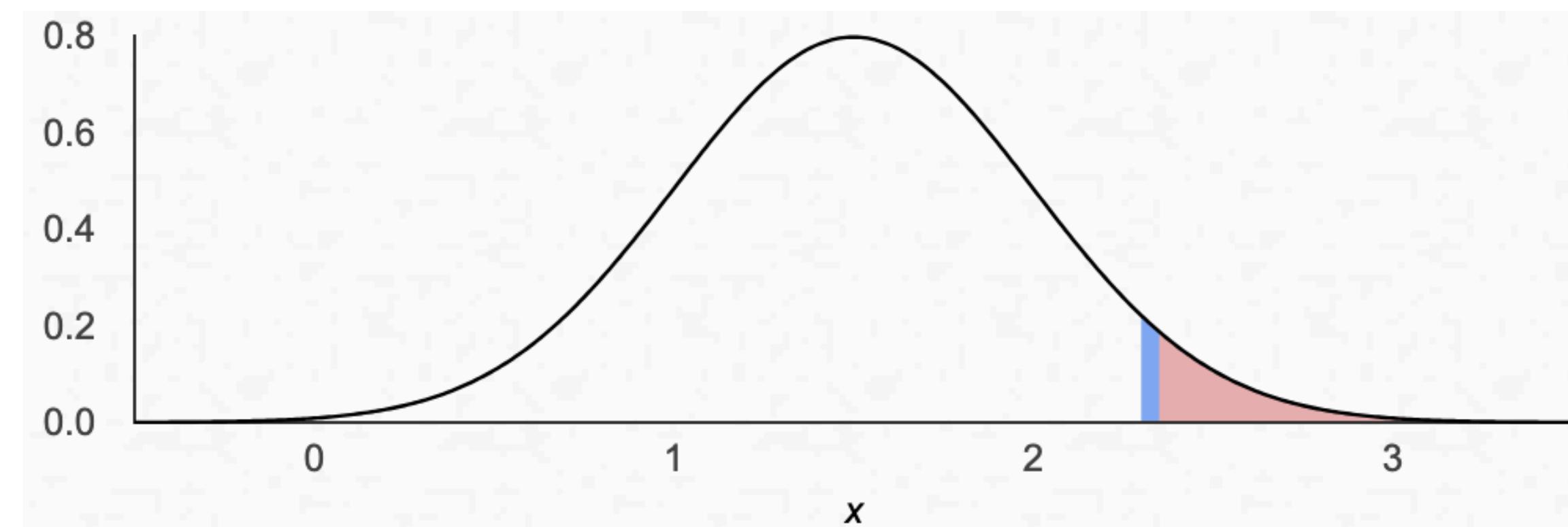
# SIGNIFICANCE LEVEL

# SIGNIFICANCE LEVEL

To make a decision, set up a rejection region by choosing the value of  $\alpha$

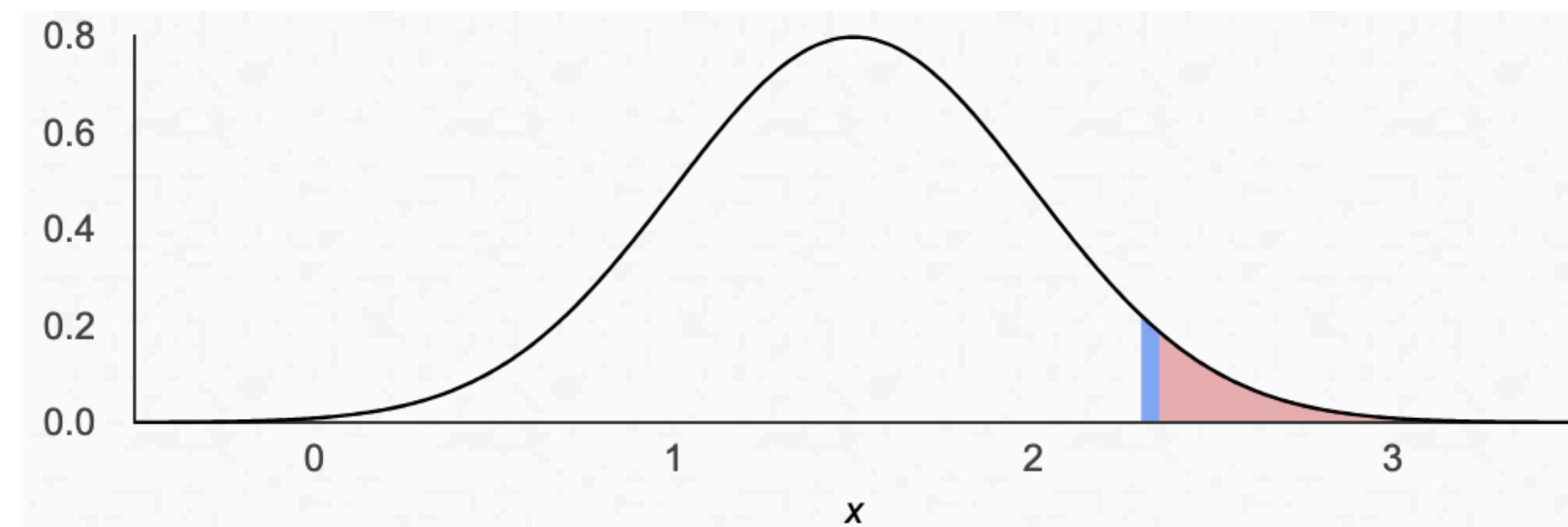
# SIGNIFICANCE LEVEL

To make a decision, set up a rejection region by choosing the value of  $\alpha$

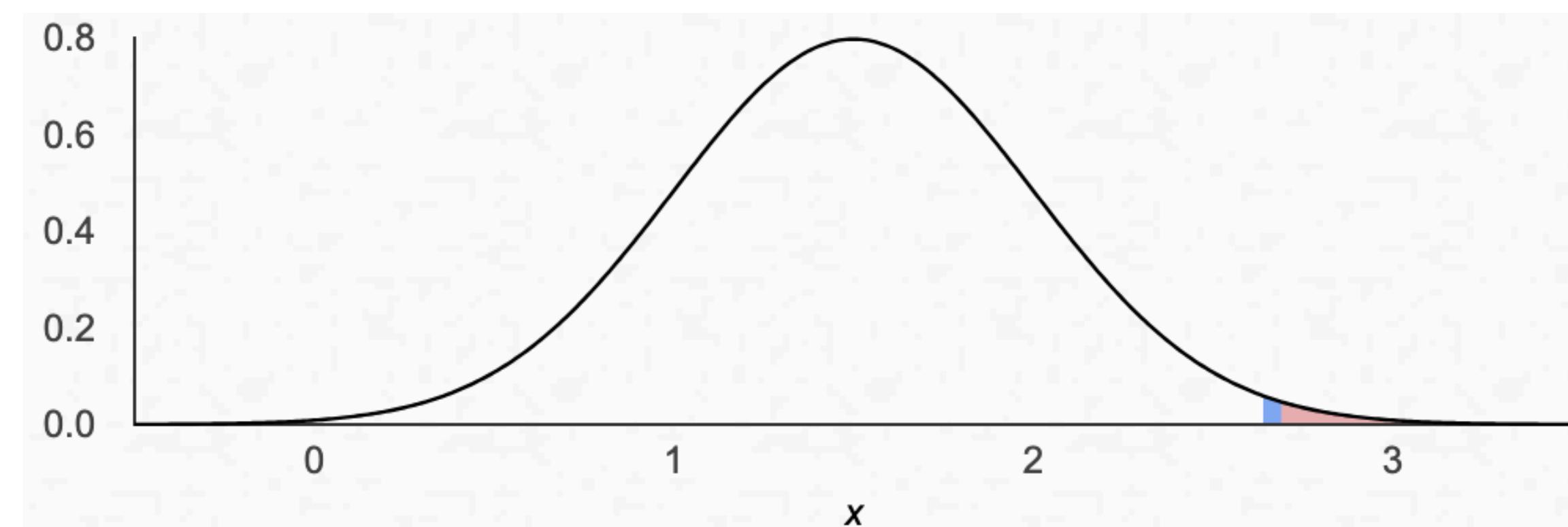


# SIGNIFICANCE LEVEL

To make a decision, set up a rejection region by choosing the value of  $\alpha$

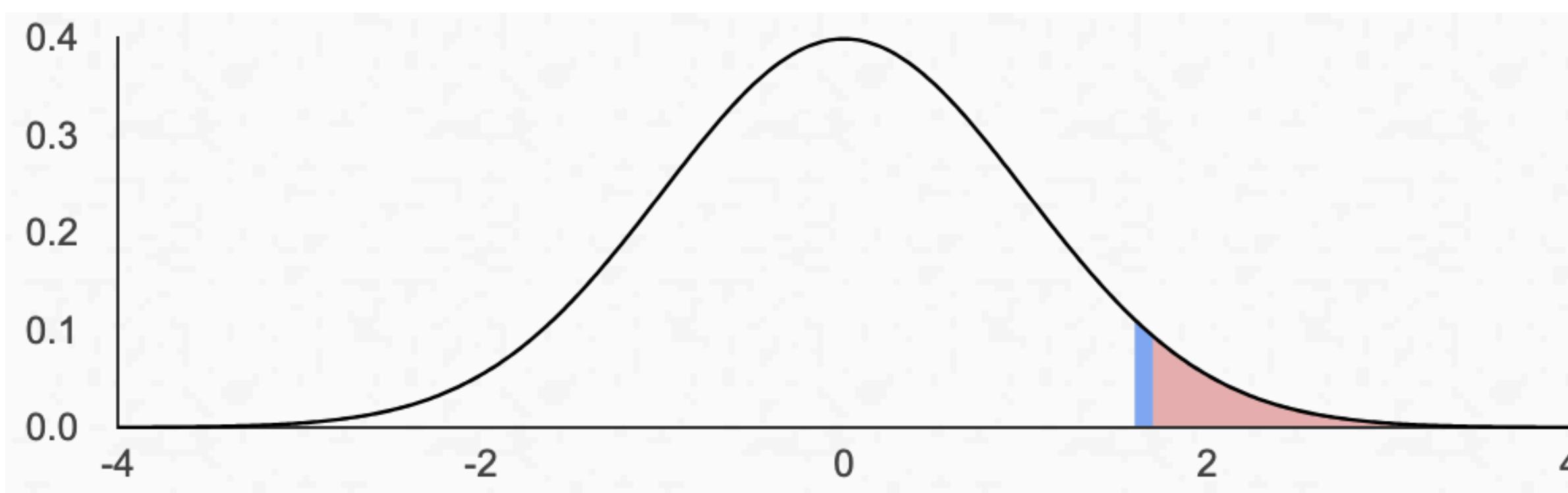


$$\alpha = 0.05$$



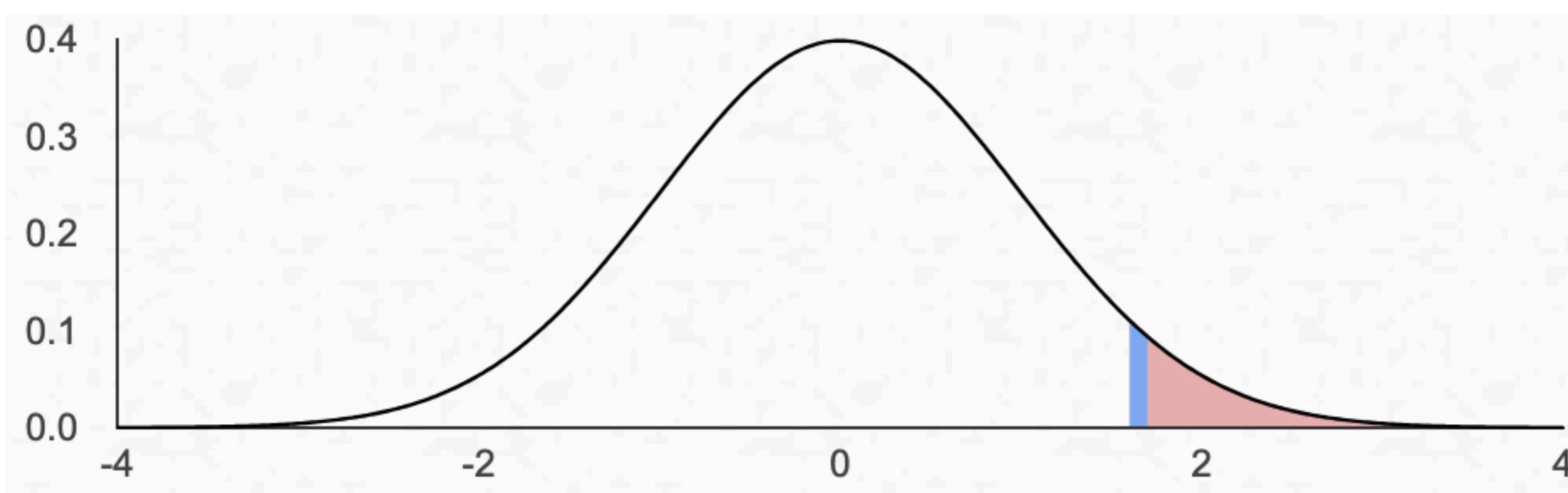
$$\alpha = 0.01$$

# Z-SCORE



# Z-SCORE

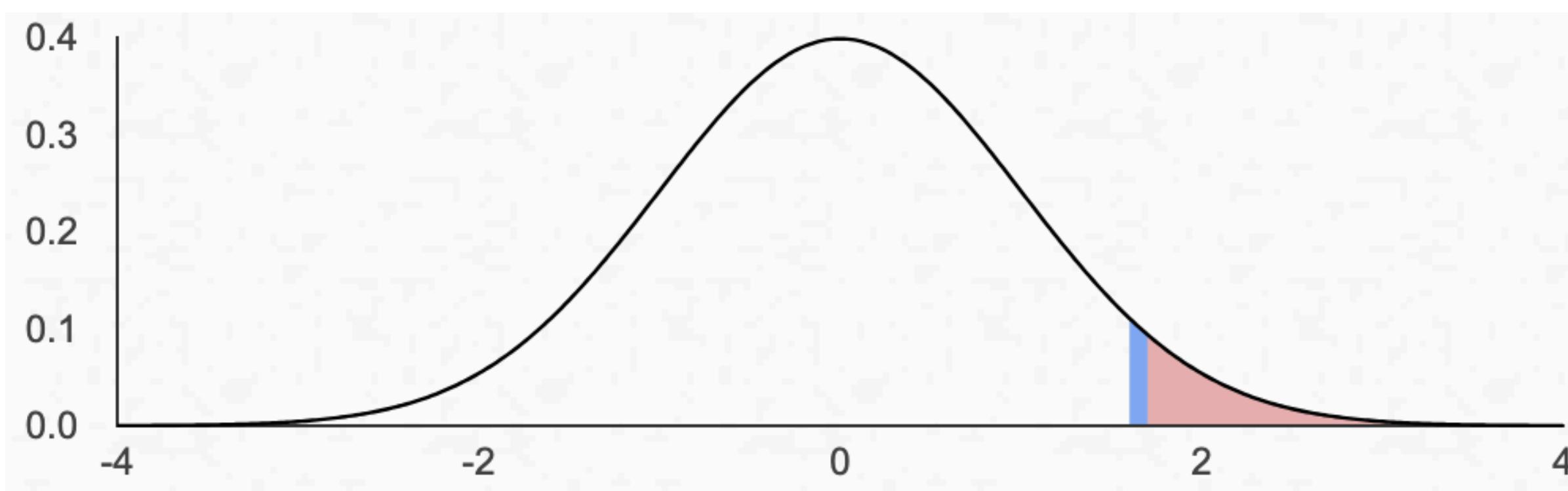
Calculating probabilities and comparing with the significance level can be tedious



# Z-SCORE

Calculating probabilities and comparing with the significance level can be tedious

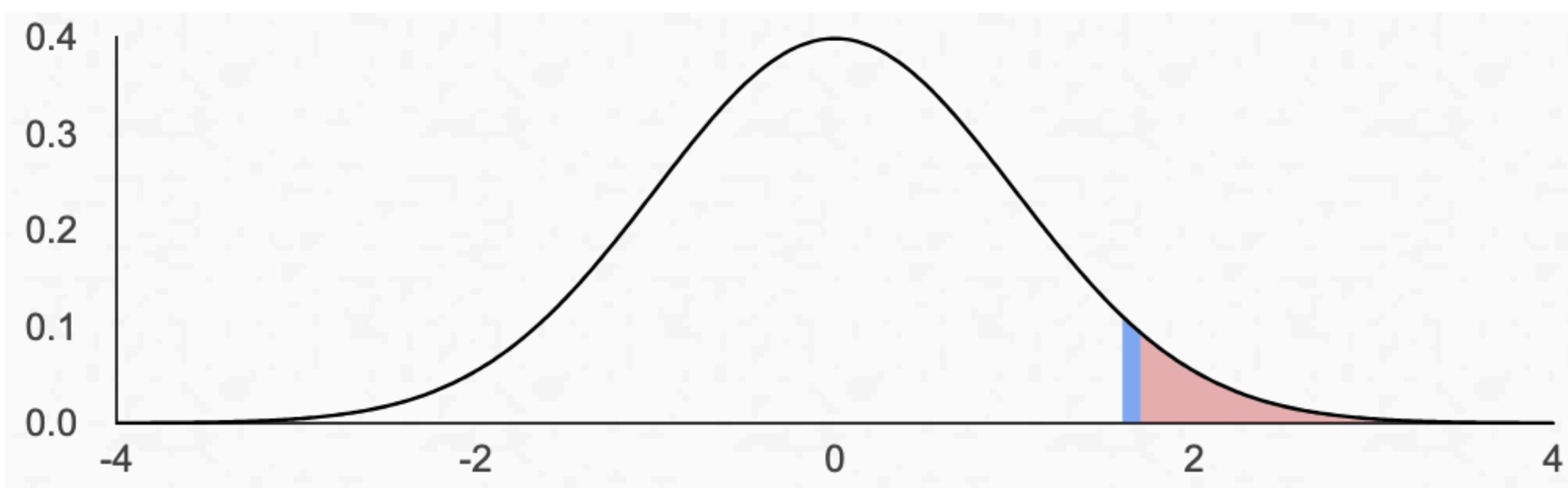
$$z = \frac{x - \mu}{\sigma/\sqrt{n}}$$



# Z-SCORE

Calculating probabilities and comparing with the significance level can be tedious

$$z = \frac{x - \mu}{\sigma/\sqrt{n}}$$

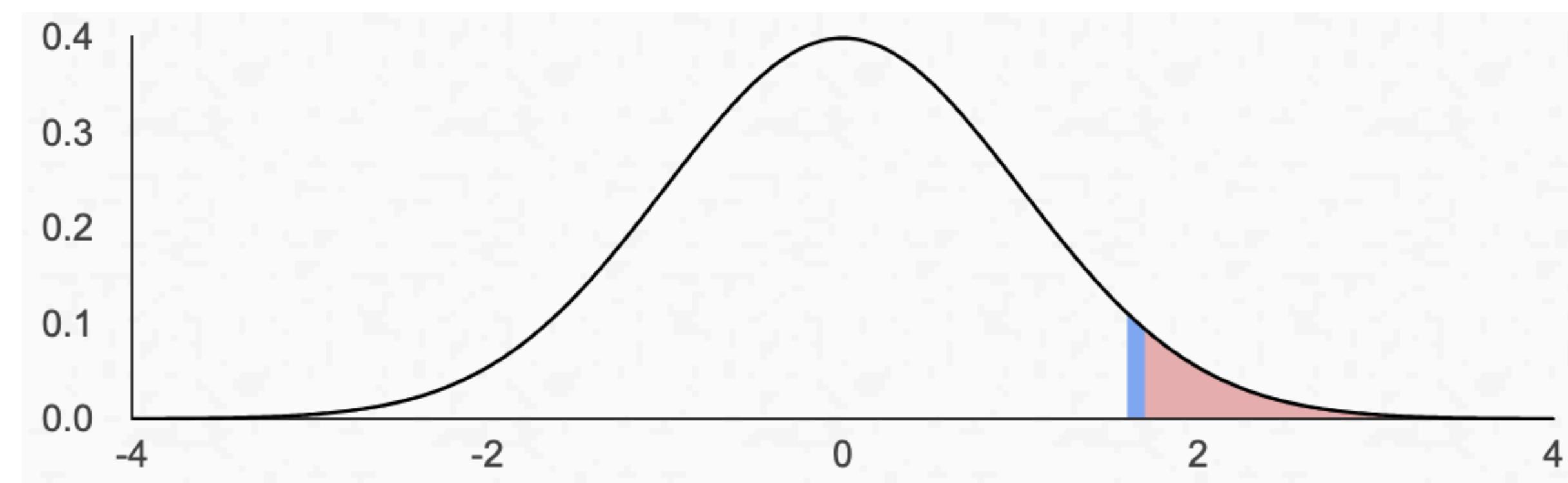


$z = 1.65$

# Z-SCORE

Calculating probabilities and comparing with the significance level can be tedious

$$z = \frac{x - \mu}{\sigma/\sqrt{n}}$$



$$z = 1.65$$

We can simply calculate the z-score of the statistic and compare it to the z-score for the significance levels

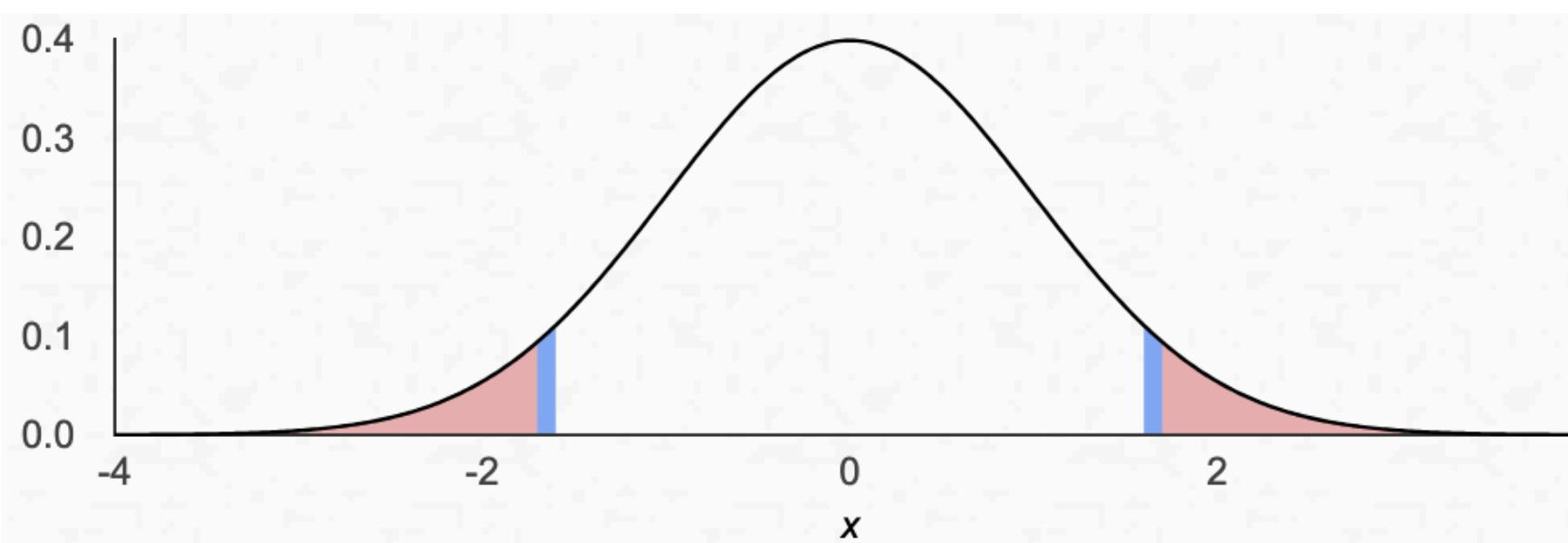
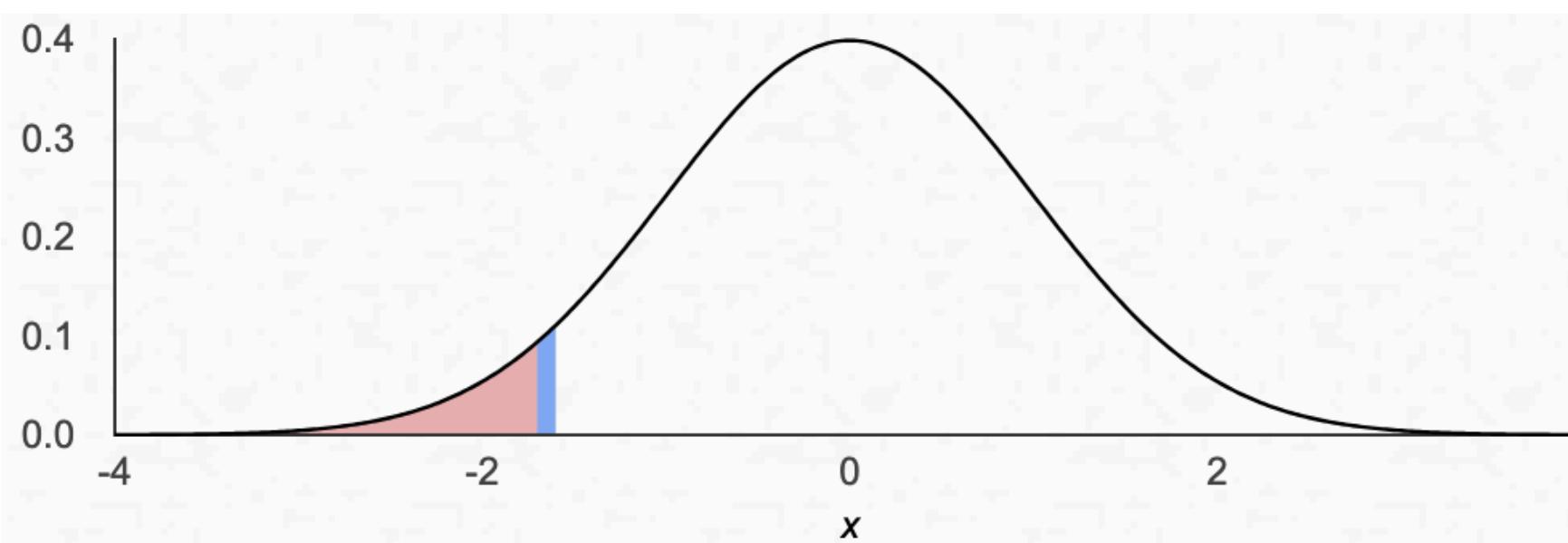
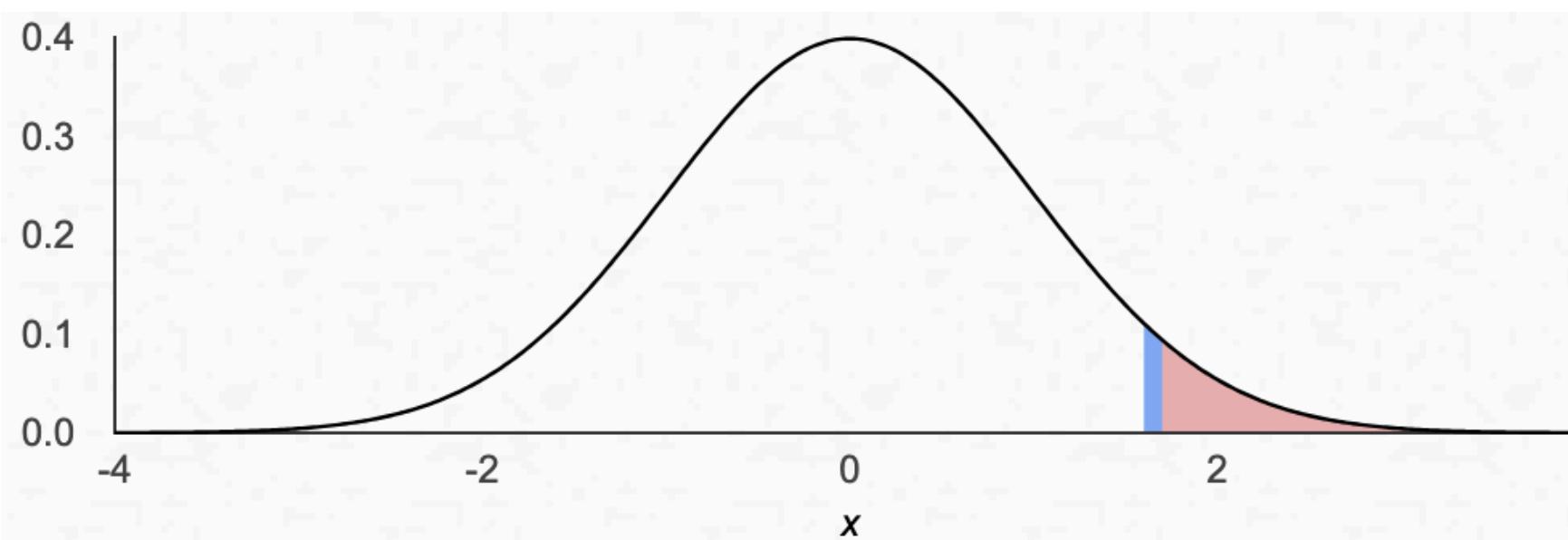
# TAILS

# TAILS

What is  
considered  
surprising is  
dependent on  
the problem at  
hand

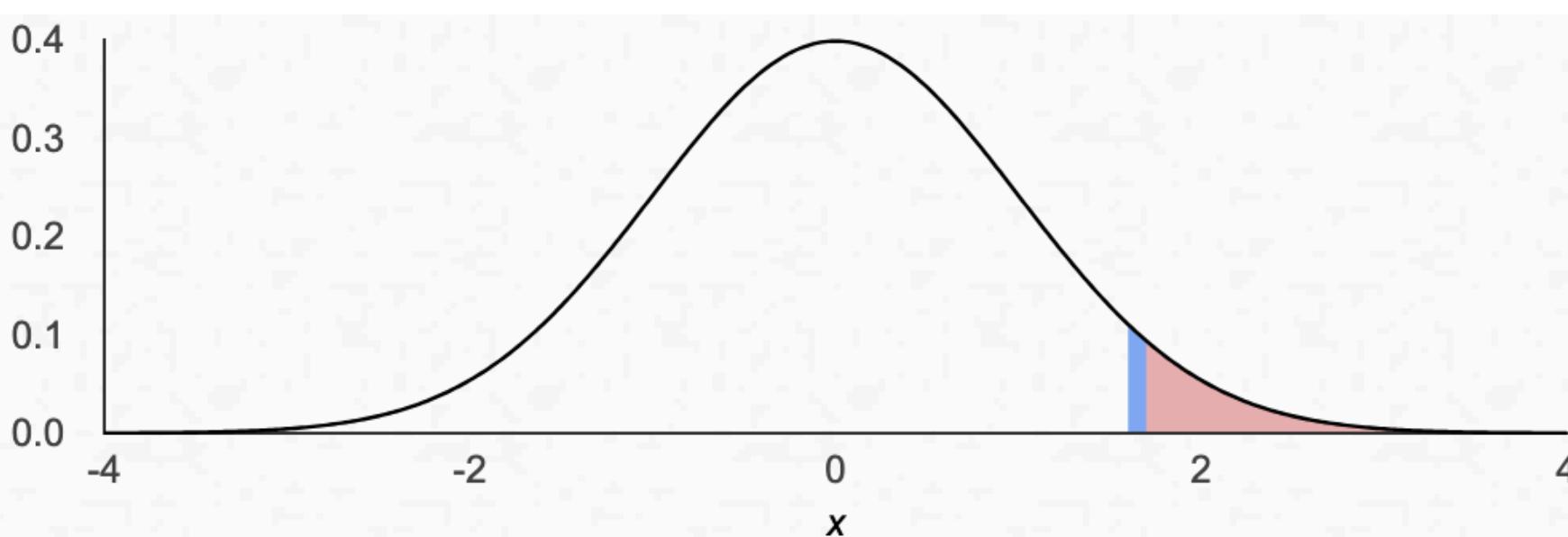
# TAILS

What is  
considered  
surprising is  
dependent on  
the problem at  
hand

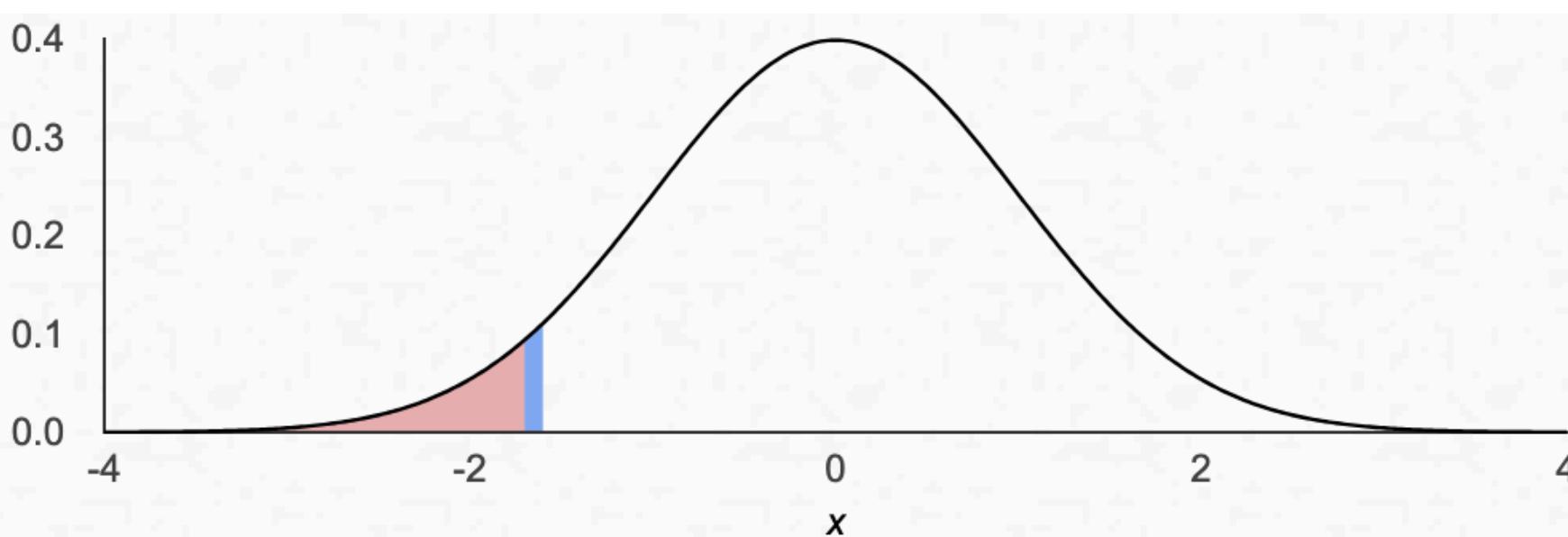


What is  
considered  
surprising is  
dependent on  
the problem at  
hand

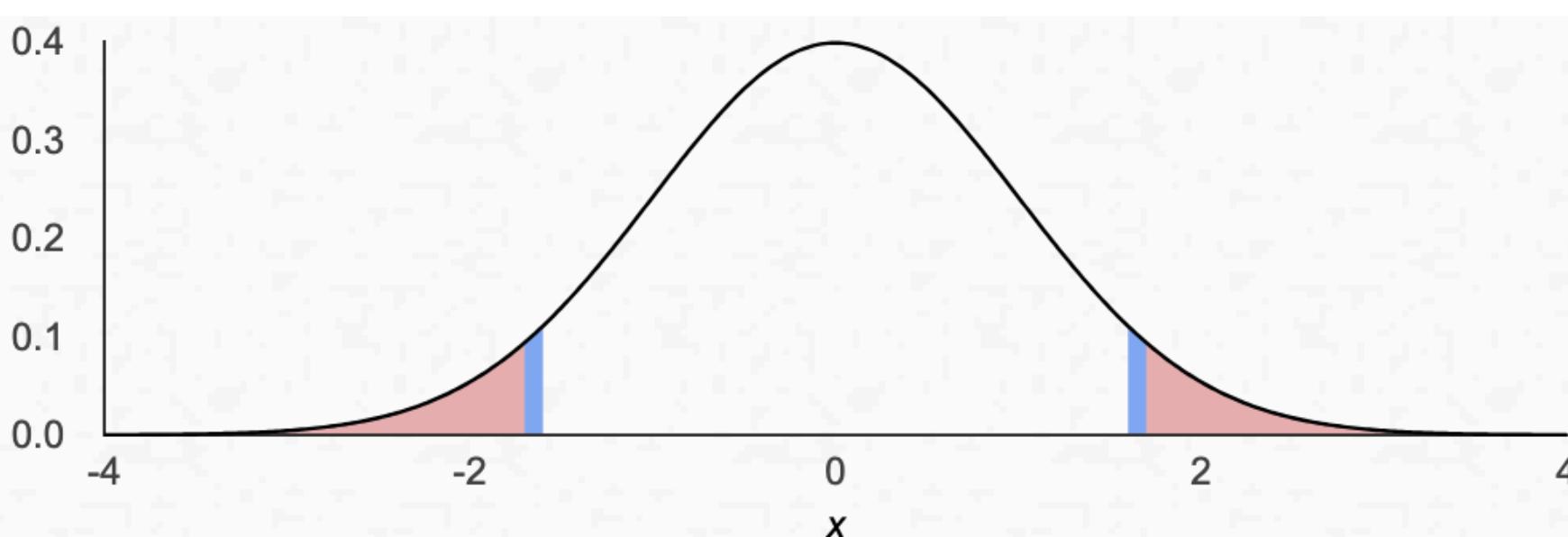
## TAILS



upper-tail: is  $x$   
significantly  
greater?



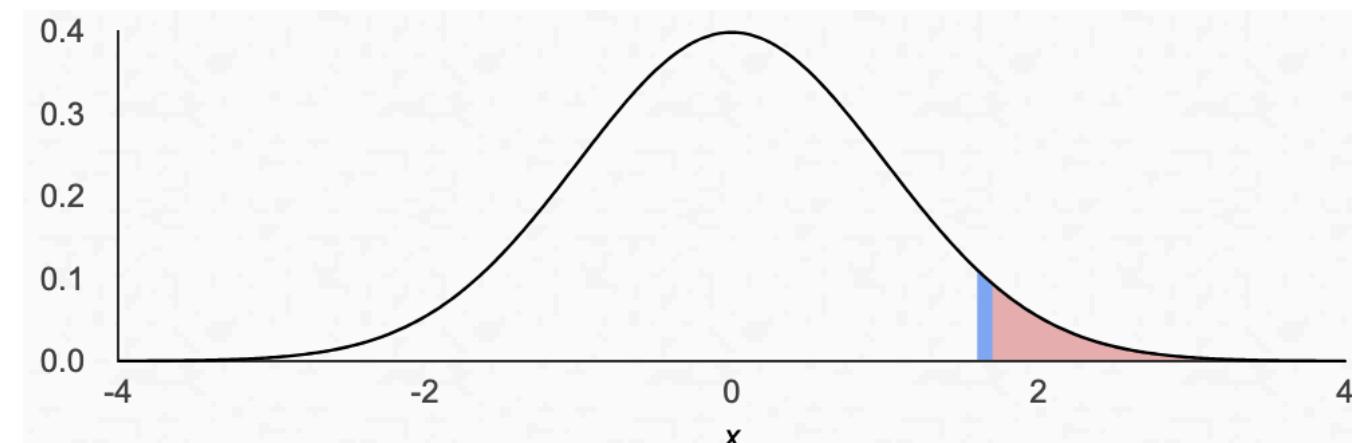
lower-tail: is  $x$   
significantly  
smaller?



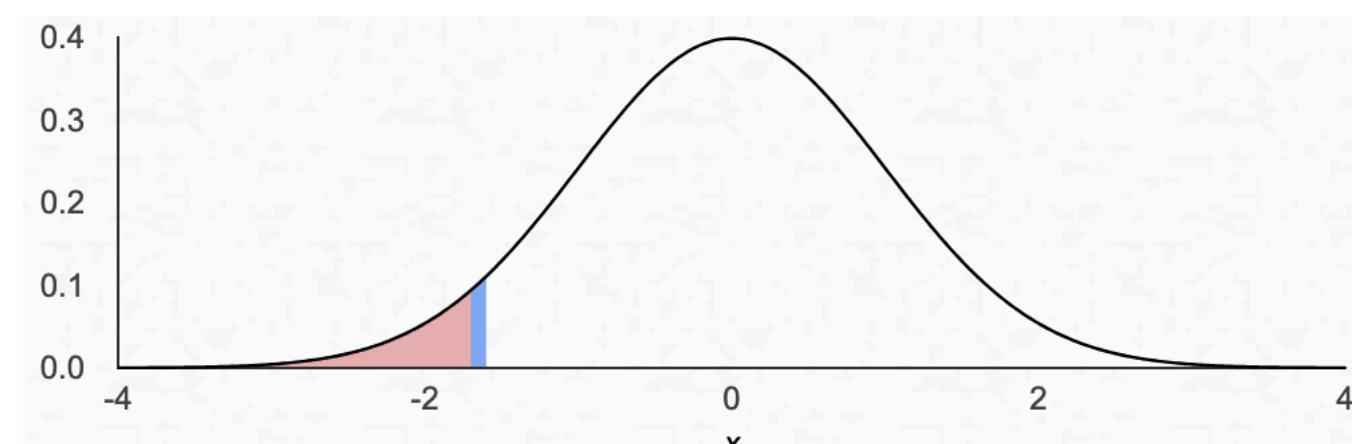
two-tail: is  $x$   
significantly  
different?

# P VALUES

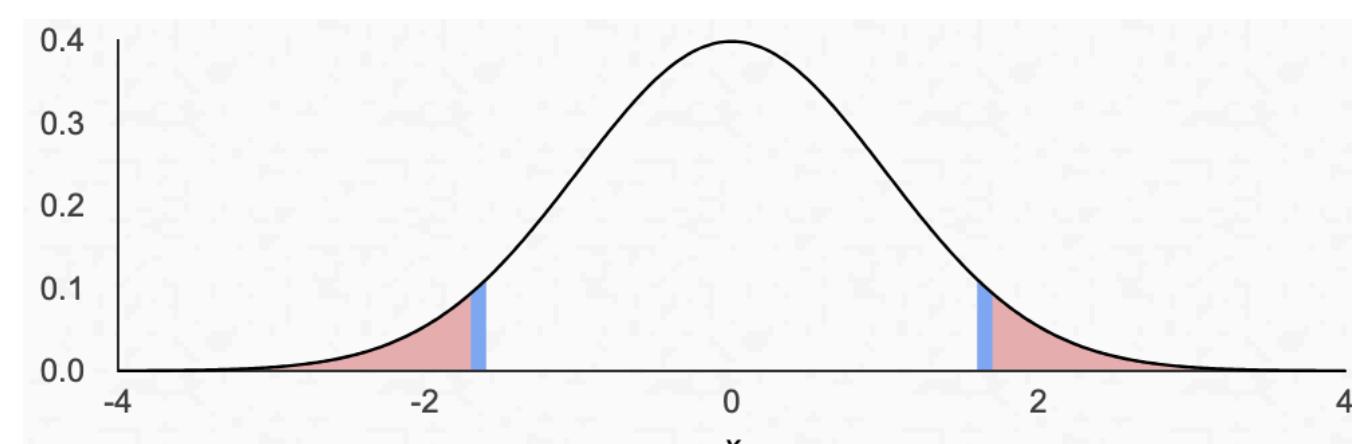
A p value is the probability of observing a statistic at least as extreme as the one we did if the null hypothesis were true.



$$\text{p-value } (x) = P(X \geq x | H_0) = 1 - P(X \leq x | H_0)$$



$$\text{p-value } (x) = 1 - P(X \leq x | H_0)$$



$$\text{p-value } (x) = 2 \times P(X \leq -|x| | H_0)$$

# RECIPE FOR HYPOTHESIS TESTING

# RECIPE FOR HYPOTHESIS TESTING

- Assume significance level  $\alpha$ ; state  $H_0$  and  $H_{\text{alternative}}$

# RECIPE FOR HYPOTHESIS TESTING

- Assume significance level  $\alpha$ ; state  $H_0$  and  $H_{\text{alternative}}$
- Calculate some statistic  $x$  whose conditional distribution is given by  $P(X|H_0)$

# RECIPE FOR HYPOTHESIS TESTING

- Assume significance level  $\alpha$ ; state  $H_0$  and  $H_{\text{alternative}}$
- Calculate some statistic  $x$  whose conditional distribution is given by  $P(X|H_0)$
- Calculate p-value

# RECIPE FOR HYPOTHESIS TESTING

- Assume significance level  $\alpha$ ; state  $H_0$  and  $H_{\text{alternative}}$
- Calculate some statistic  $x$  whose conditional distribution is given by  $P(X|H_0)$
- Calculate p-value
- If p-value falls in rejection region then null hypothesis can be rejected; else null hypothesis cannot be rejected

# ERRORS

# ERRORS

- Type I error: We incorrectly rejected the null hypothesis

# ERRORS

- Type I error: We incorrectly rejected the null hypothesis
- Type II error: We incorrectly failed to reject the null hypothesis

# ERRORS

- Type I error: We incorrectly rejected the null hypothesis
- Type II error: We incorrectly failed to reject the null hypothesis

		Test results	
		keep null	reject null
Truth	keep null	Type I error $\alpha$	Power
	reject null	Type II error $\beta$	



The Boy who Cried Wolf

- In the Aesop's fable, assuming that there typically is no wolf, the villagers make two types of errors
  - They believe there was a wolf when there was none
  - They believe there was no wolf when there was one



The Boy who Cried Wolf



The Boy who Cried Wolf

- In the Aesop's fable – The Boy who Cried Wolf – the villagers make two types of errors
  - They believe there was a wolf when there was none [TYPE I]
  - They believe there was no wolf when there was one [TYPE II]



The Boy who Cried Wolf

# ERRORS

# ERRORS

- A well-calibrated statistical test should have acceptable Type I error rate and high statistical power

# ERRORS

- A well-calibrated statistical test should have acceptable Type I error rate and high statistical power
- If  $\alpha = 0.05$  and we do 100 tests, we expect to make 5 mistakes

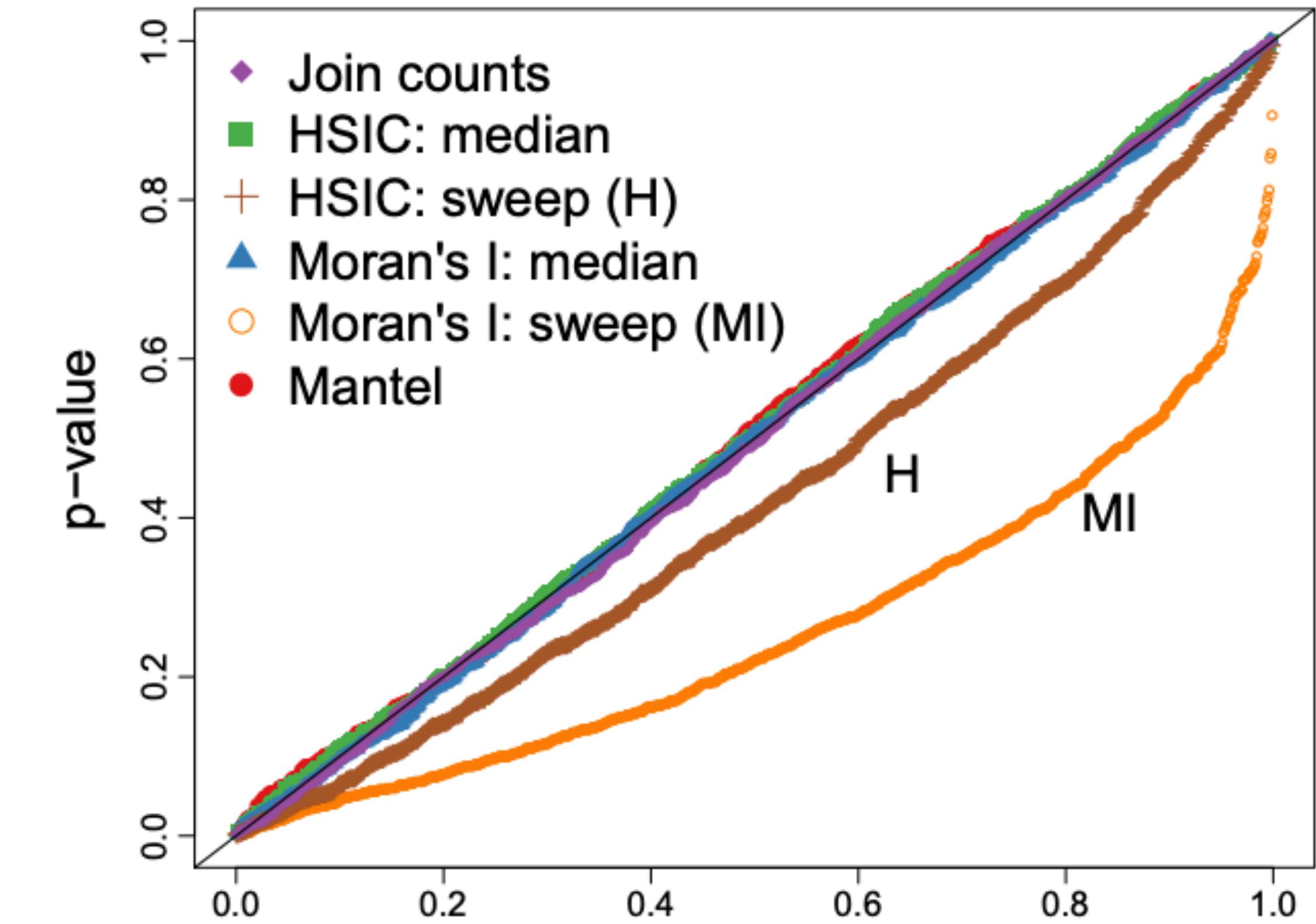
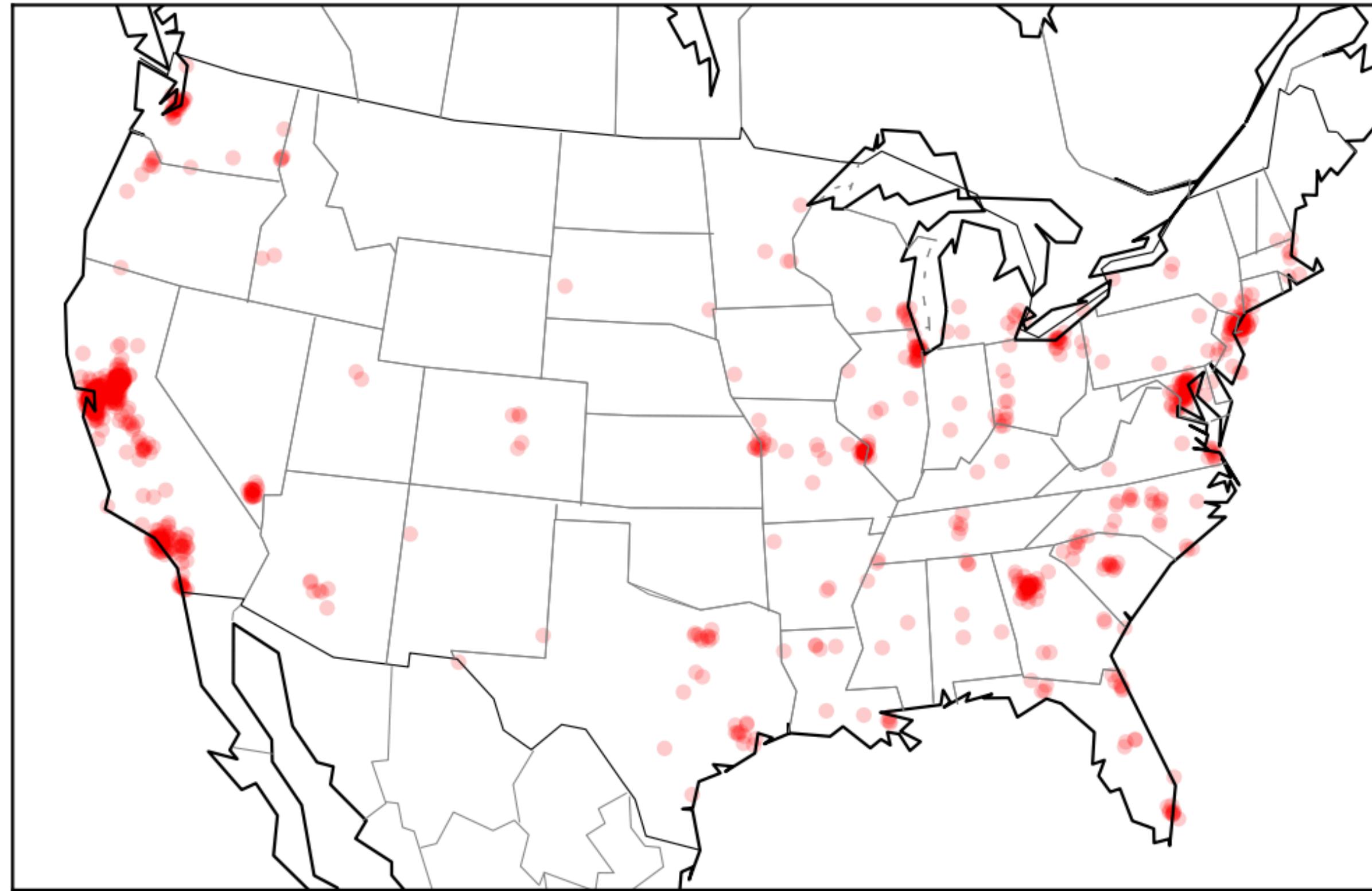
# ERRORS

- A well-calibrated statistical test should have acceptable Type I error rate and high statistical power
- If  $\alpha = 0.05$  and we do 100 tests, we expect to make 5 mistakes

		Test results	
		keep null	reject null
Truth	keep null	Type I error $\alpha$	Type II error $\beta$
	reject null	Power	



# Test if lexical variation is independent of geography



# MULTIPLE HYPOTHESIS CORRECTIONS

# MULTIPLE HYPOTHESIS CORRECTIONS

- When we do multiple tests, we want to correct for the likelihood of getting statistical significance by chance

# MULTIPLE HYPOTHESIS CORRECTIONS

- When we do multiple tests, we want to correct for the likelihood of getting statistical significance by chance
- Apply bonferroni correction which conservatively sets a lower significance threshold based on number of tests

# MULTIPLE HYPOTHESIS CORRECTIONS

- When we do multiple tests, we want to correct for the likelihood of getting statistical significance by chance
- Apply bonferroni correction which conservatively sets a lower significance threshold based on number of tests

$$\alpha = \frac{\alpha_0}{n}$$

Here the significance level  $\alpha_0$  is adjusted

# CONFIDENCE INTERVALS

# CONFIDENCE INTERVALS

- In many instances, instead of doing a statistical test, we want to bound the error of the metric

# CONFIDENCE INTERVALS

- In many instances, instead of doing a statistical test, we want to bound the error of the metric
- Confidence intervals helps us quantify the range in which the observed metric will lie for the unobserved population

# CONFIDENCE INTERVALS

# CONFIDENCE INTERVALS

- The CI is statistically determined with some parametric assumptions

# CONFIDENCE INTERVALS

- The CI is statistically determined with some parametric assumptions
- The observed value is assumed to be a point estimate of the mean

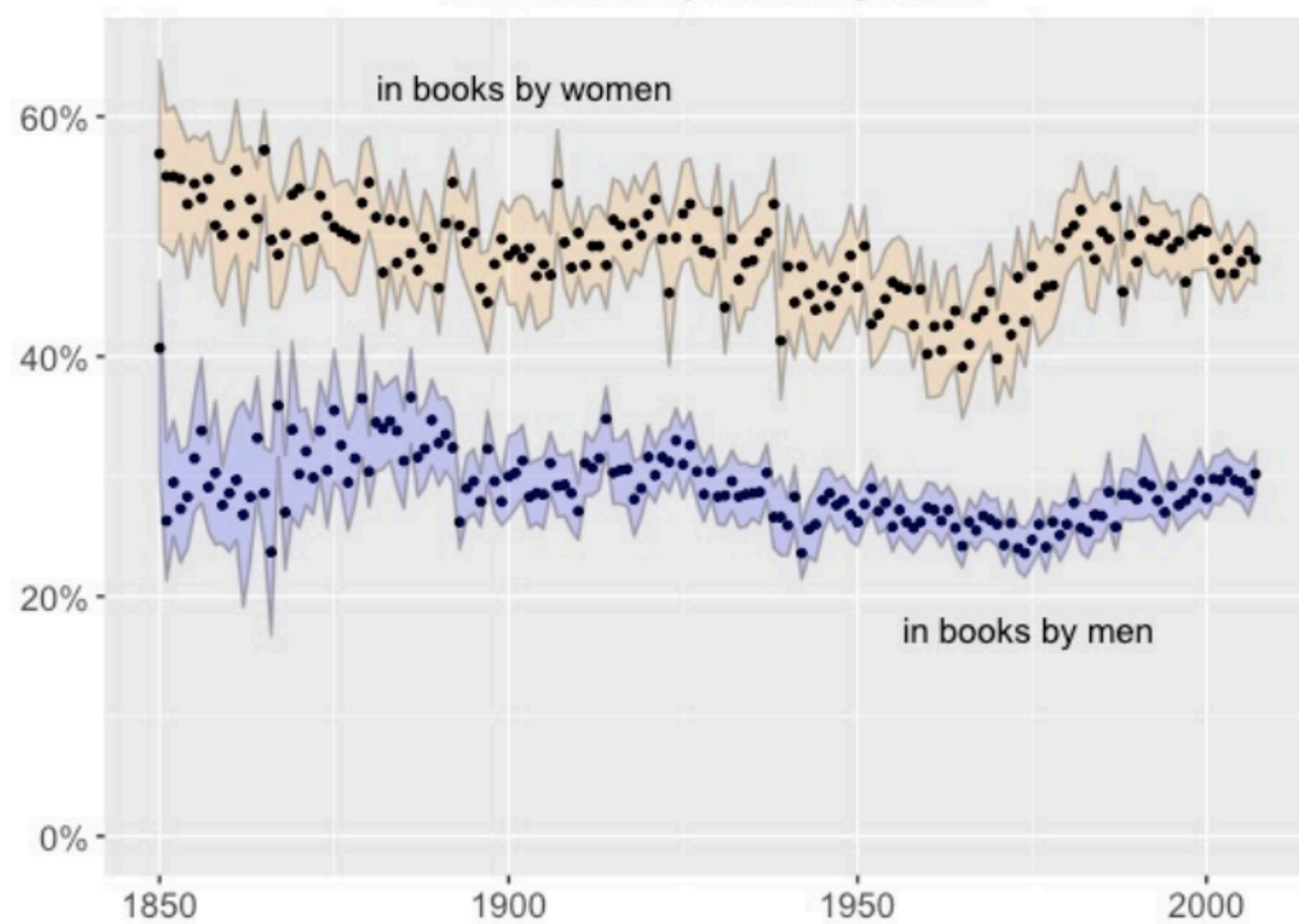
# CONFIDENCE INTERVALS

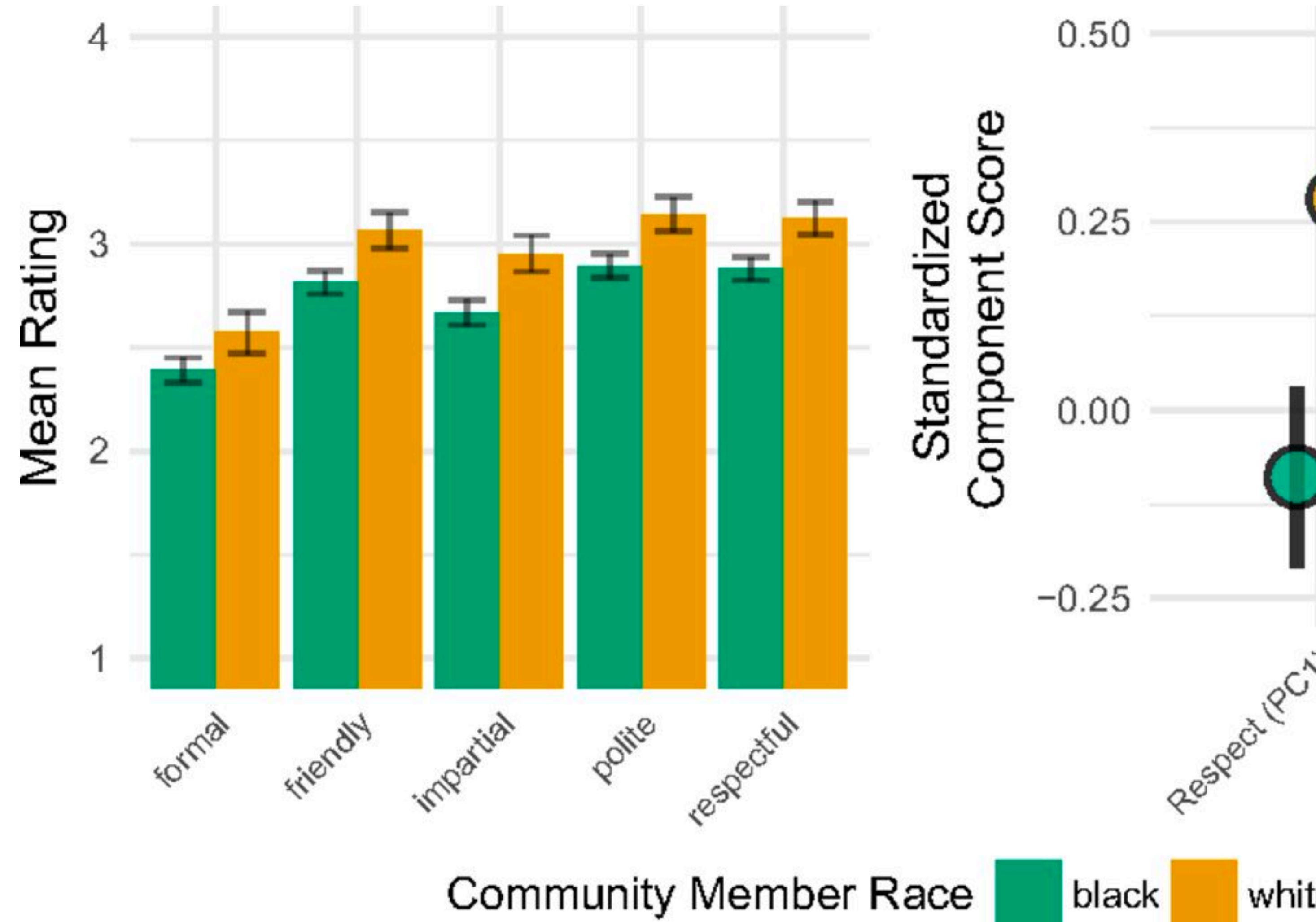
- The CI is statistically determined with some parametric assumptions
- The observed value is assumed to be a point estimate of the mean

$$CI = \bar{x} \pm z \cdot \frac{s}{\sqrt{n}}$$

Mean value      Lower/Upper limit      z-value for the confidence level  
Standard deviation      Sample size

Description of women, as a percentage of characterization,  
broken out by author gender





# ABLATION TESTING

# ABLATION TESTING

- Are a set of features really important for the model?

# ABLATION TESTING

- Are a set of features really important for the model?
- One can statistically test this:
  - Create model with all the features and evaluate
  - Create another model with test features removed and evaluate
  - Compare the difference in performance as a statistic and check for its statistical significance

	Features	# of features	frequency or presence?	NB	ME	SVM
(1)	unigrams	16165	freq.	<b>78.7</b>	N/A	72.8
(2)	unigrams	"	pres.	81.0	80.4	<b>82.9</b>
(3)	unigrams+bigrams	32330	pres.	80.6	80.8	<b>82.7</b>
(4)	bigrams	16165	pres.	77.3	<b>77.4</b>	77.1
(5)	unigrams+POS	16695	pres.	81.5	80.4	<b>81.9</b>
(6)	adjectives	2633	pres.	77.0	<b>77.7</b>	75.1
(7)	top 2633 unigrams	2633	pres.	80.3	81.0	<b>81.4</b>
(8)	unigrams+position	22430	pres.	81.0	80.1	<b>81.6</b>

# RECIPE FOR HYPOTHESIS TESTING

# RECIPE FOR HYPOTHESIS TESTING

- Assume significance level  $\alpha$ ; state  $H_0$  and  $H_{\text{alternative}}$

# RECIPE FOR HYPOTHESIS TESTING

- Assume significance level  $\alpha$ ; state  $H_0$  and  $H_{\text{alternative}}$
- Calculate some statistic  $x$  whose conditional distribution is given by  $P(X|H_0)$

# RECIPE FOR HYPOTHESIS TESTING

- Assume significance level  $\alpha$ ; state  $H_0$  and  $H_{\text{alternative}}$
- Calculate some statistic  $x$  whose conditional distribution is given by  $P(X|H_0)$
- Calculate p-value

# RECIPE FOR HYPOTHESIS TESTING

- Assume significance level  $\alpha$ ; state  $H_0$  and  $H_{\text{alternative}}$
- Calculate some statistic  $x$  whose conditional distribution is given by  $P(X|H_0)$
- Calculate p-value
- If p-value falls in rejection region then null hypothesis can be rejected; else null hypothesis cannot be rejected

# HYPOTHESIS TESTING

# HYPOTHESIS TESTING

- What should be  $P(X | H_0)$ ?

# HYPOTHESIS TESTING

- What should be  $P(X|H_0)$ ?
- In many situations, we can use parametric distributions to characterize  $P(X|H_0)$ 
  - Binomial (probability of success  $p$ , #trials  $n$ )
  - Normal (mean  $\mu$  and standard deviation  $\sigma$ )

# PARAMETRIC TESTS

# PARAMETRIC TESTS

- For these tests, probabilities can be calculated by plugging it in an equation

# PARAMETRIC TESTS

- For these tests, probabilities can be calculated by plugging it in an equation
- E.g. Assume we observe 75 heads from 100 trials of a fair coin, then calculate  $\Pr(X=75 | p=0.5, n=100)$  using a binomial distribution's parametric form

$$P(x) = \binom{n}{x} p^x (1-p)^{n-x}$$

# CENTRAL LIMIT THEOREM

# CENTRAL LIMIT THEOREM

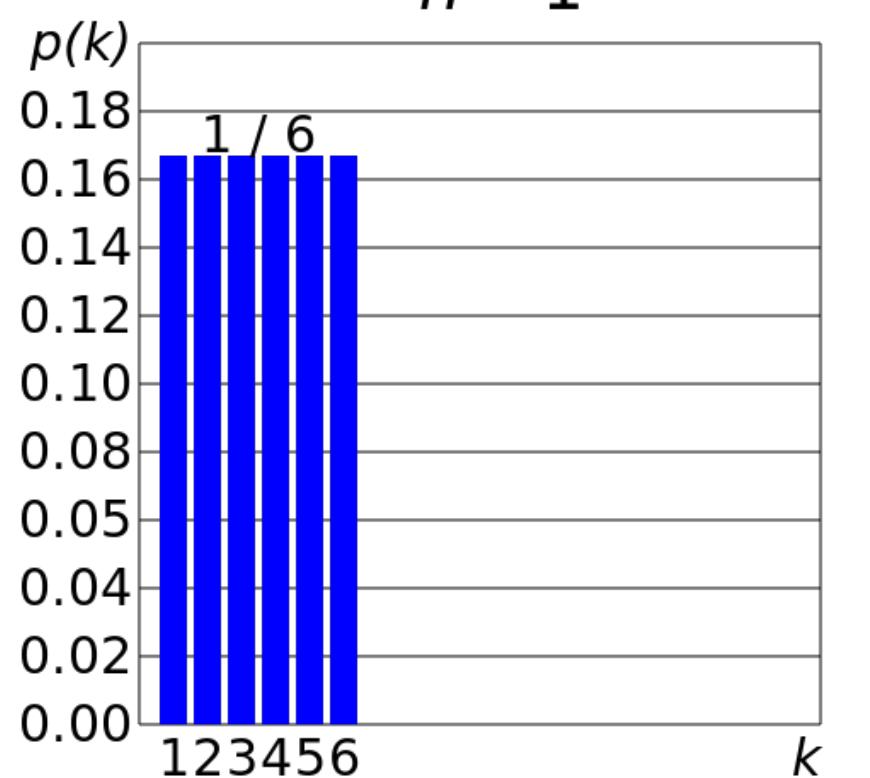
- Often our statistic is calculated by averaging over multiple independent instances (e.g., accuracy)

# CENTRAL LIMIT THEOREM

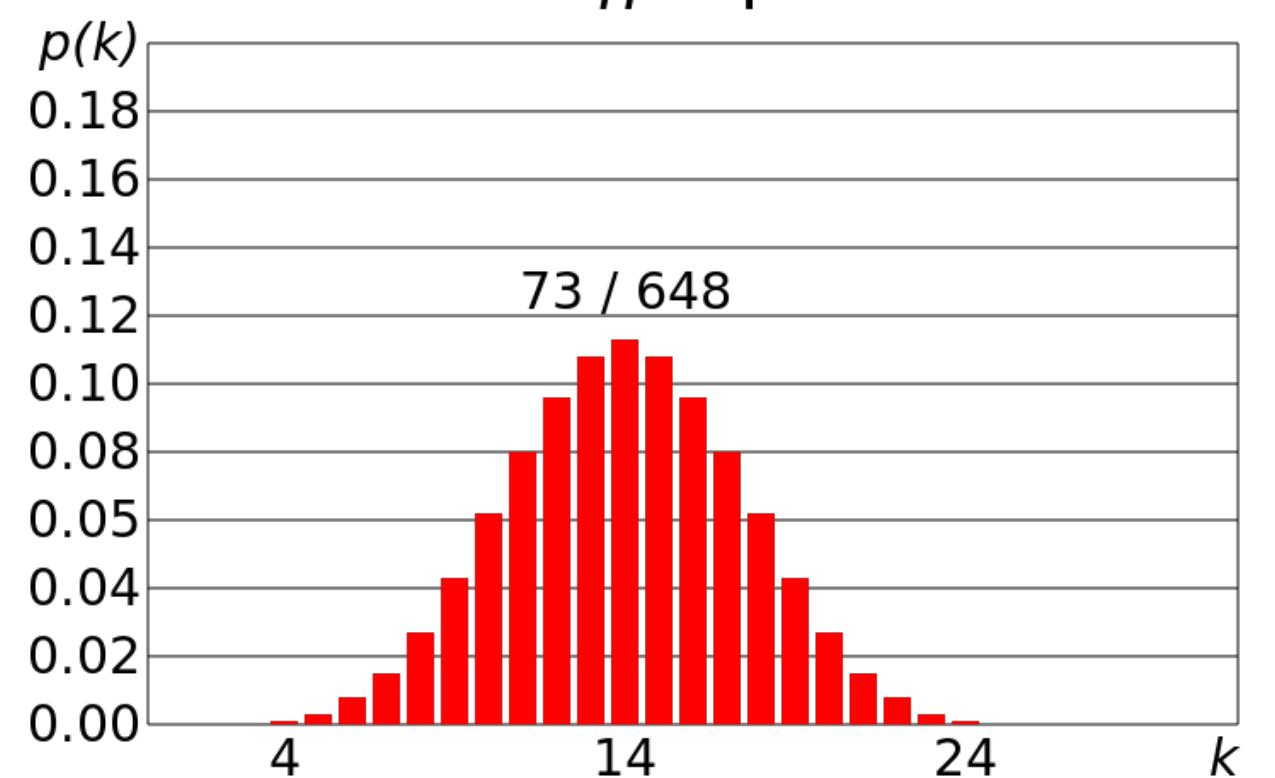
- Often our statistic is calculated by averaging over multiple independent instances (e.g., accuracy)
- According to CLT, the average of independent and identically distributed random variables tends to be a normal distribution, even if the random variables are not normally distributed



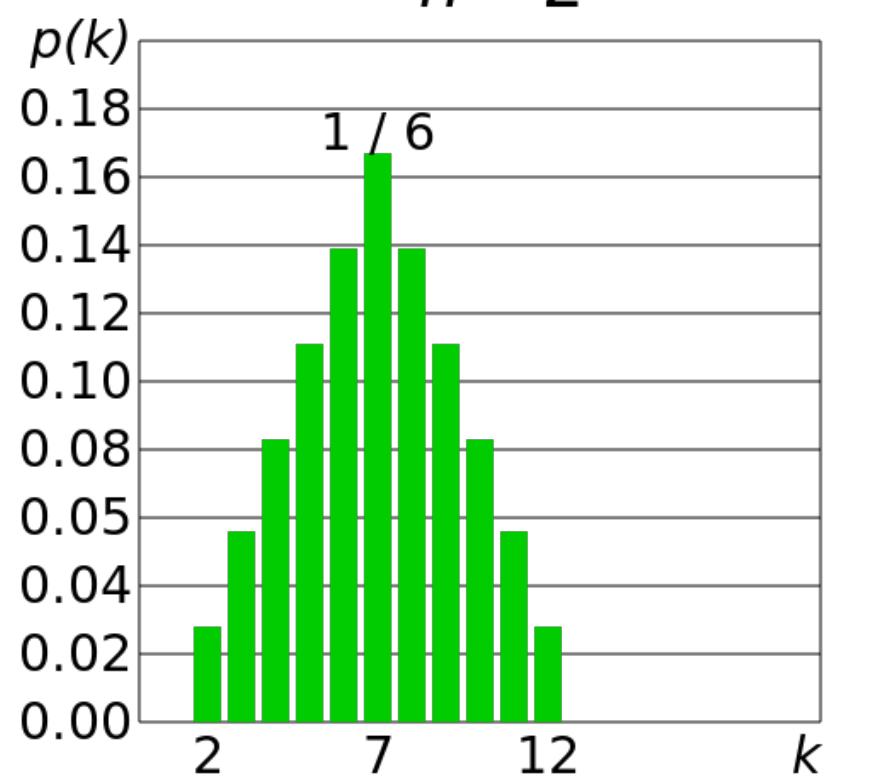
$n = 1$



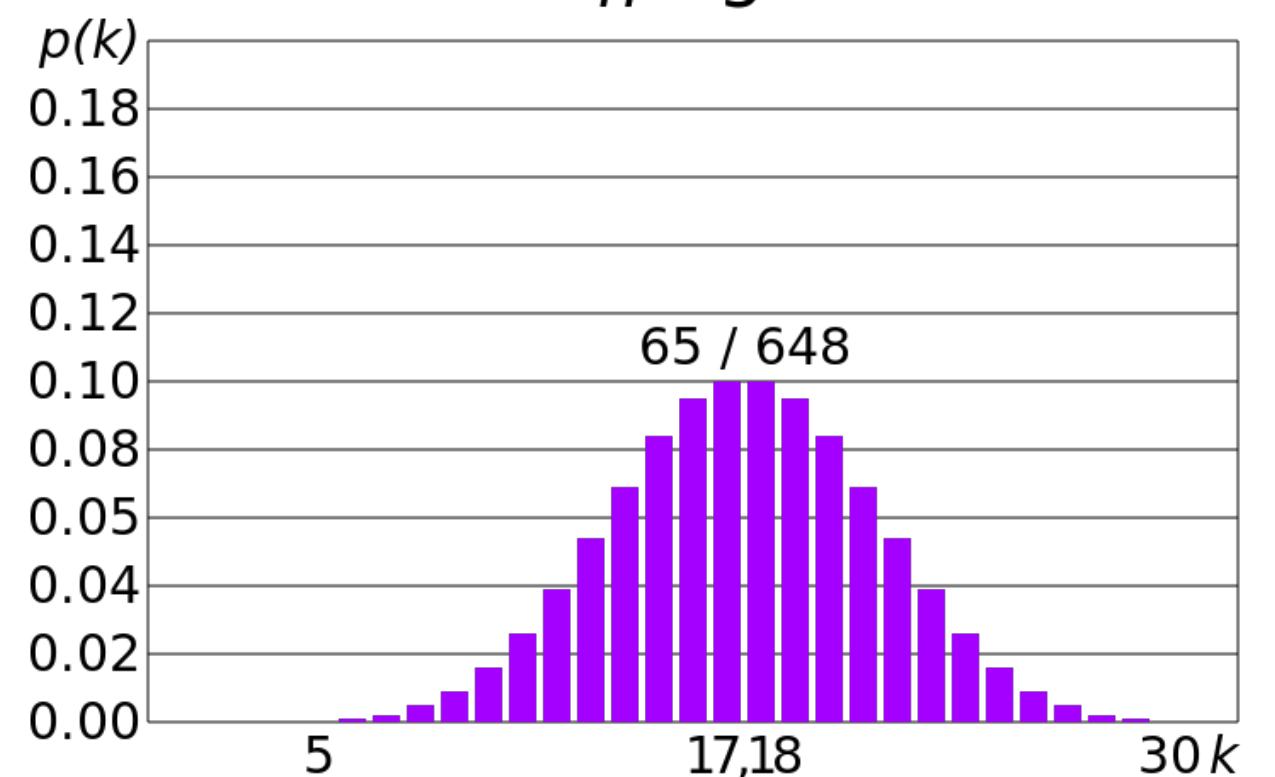
$n = 4$



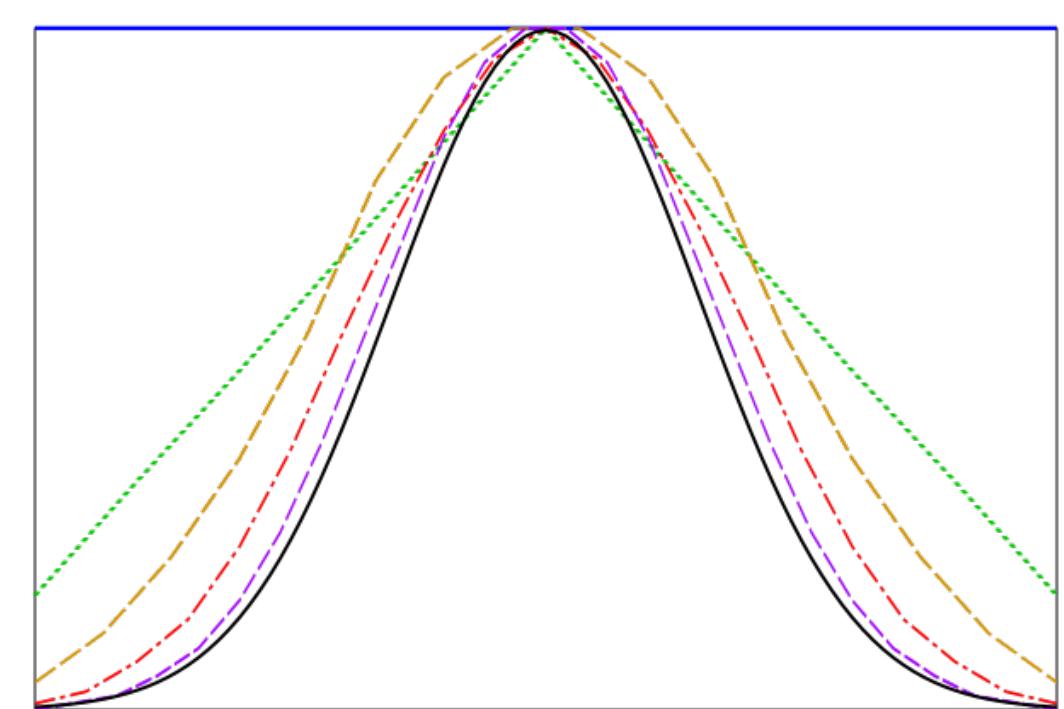
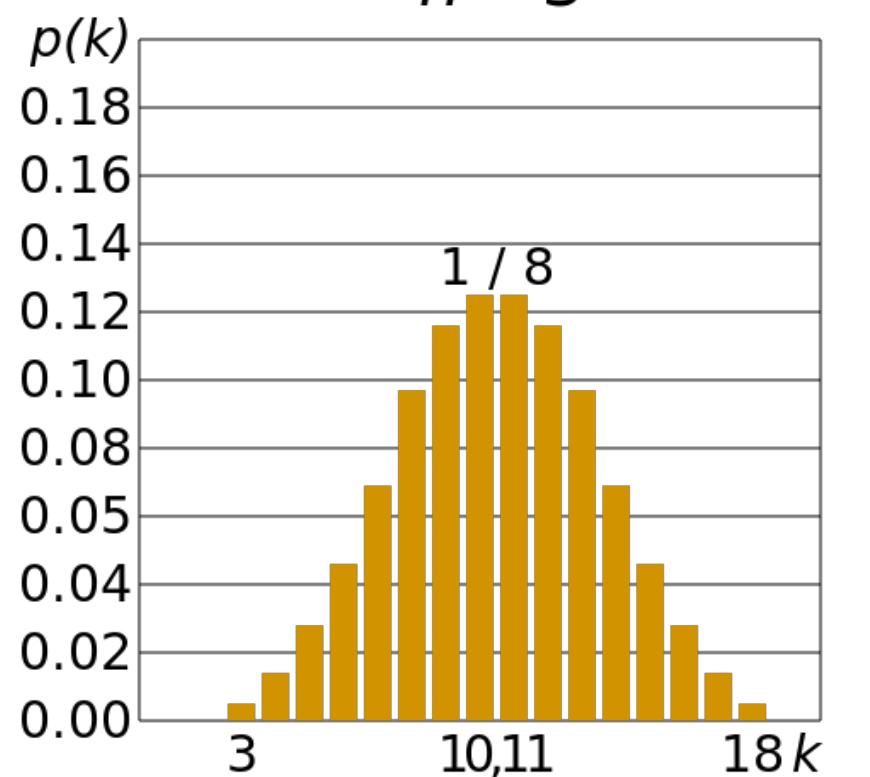
$n = 2$



$n = 5$



$n = 3$



CLT allows us to parametrize the distribution of many statistics that are like sample averages

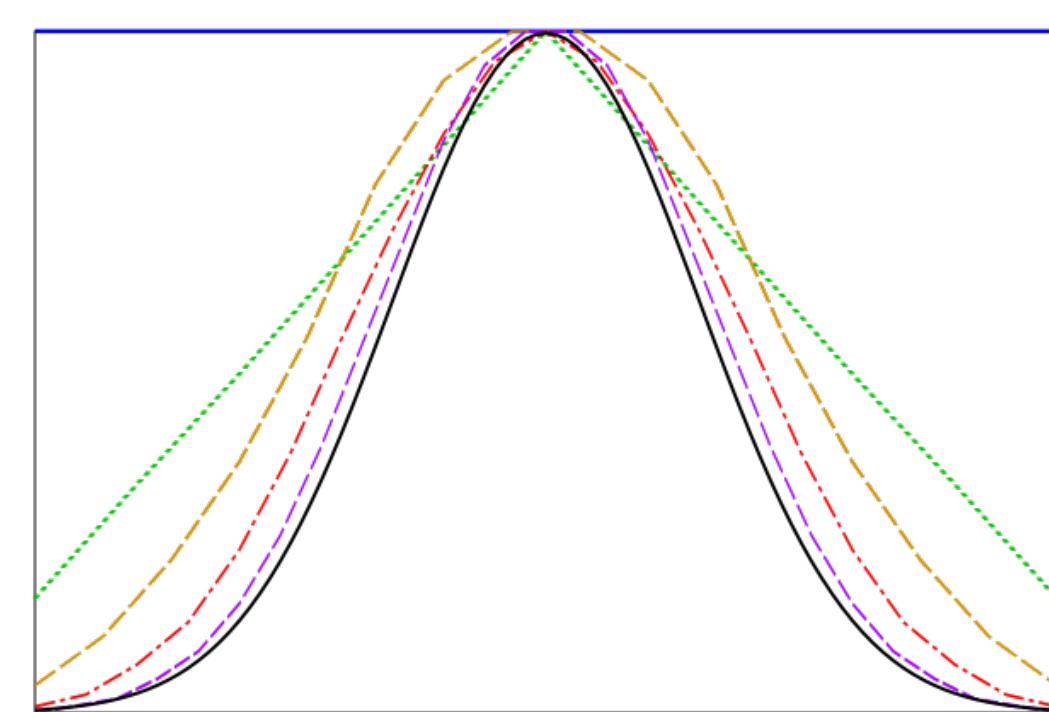
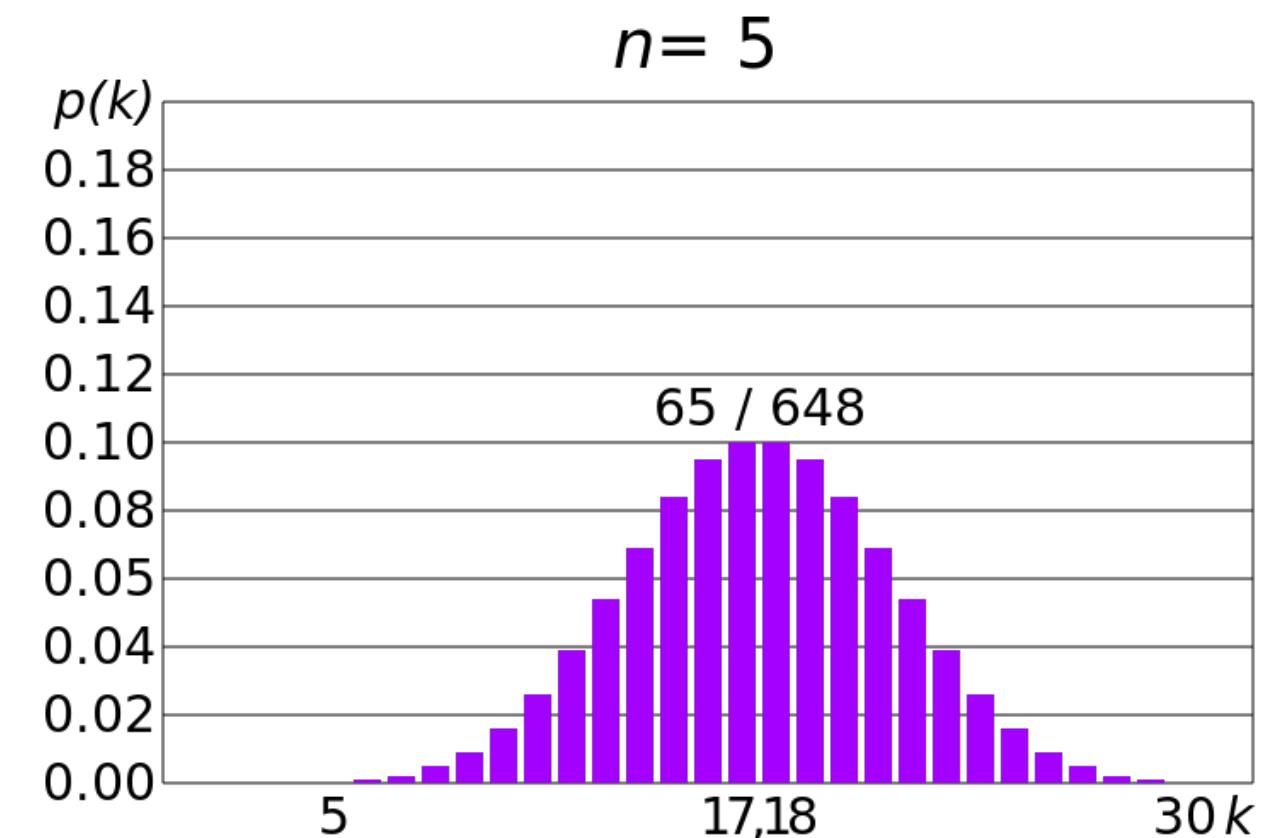
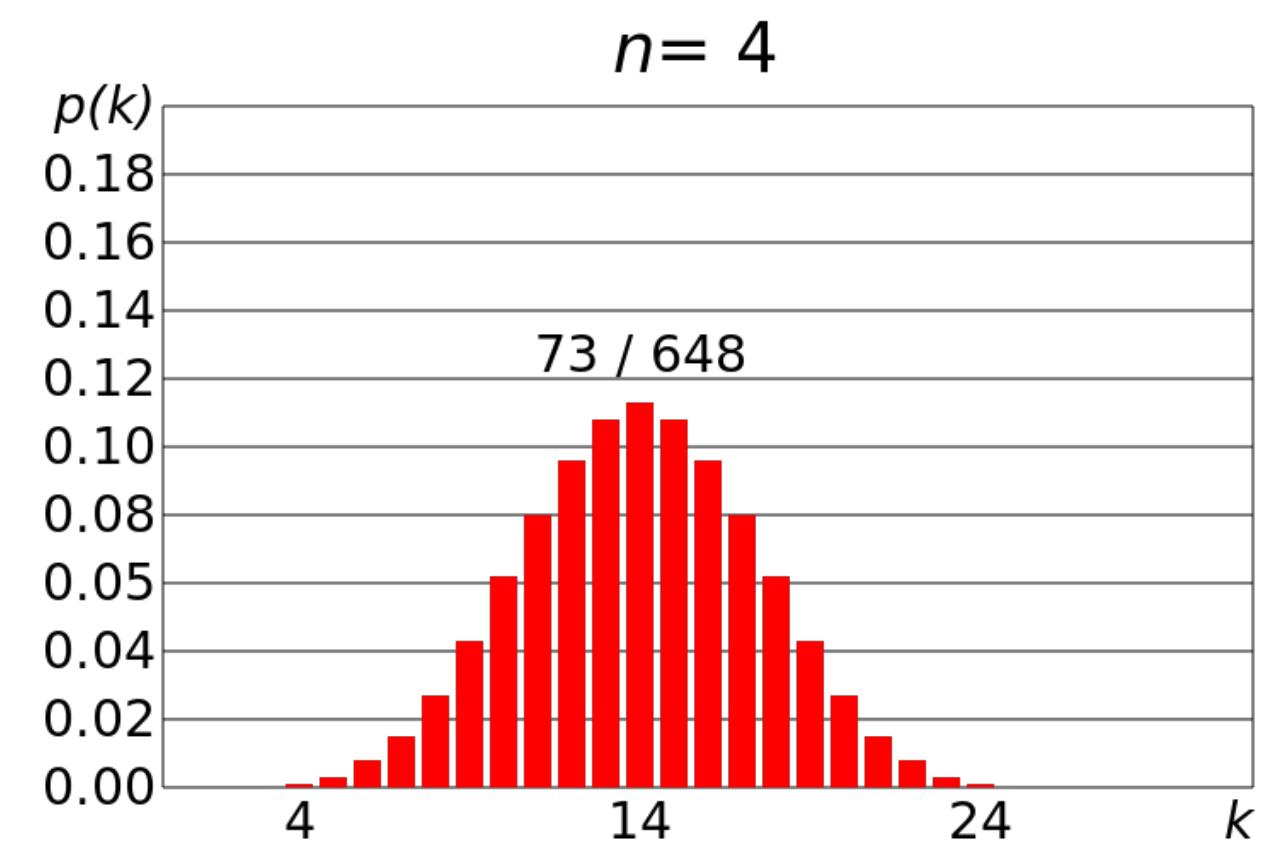
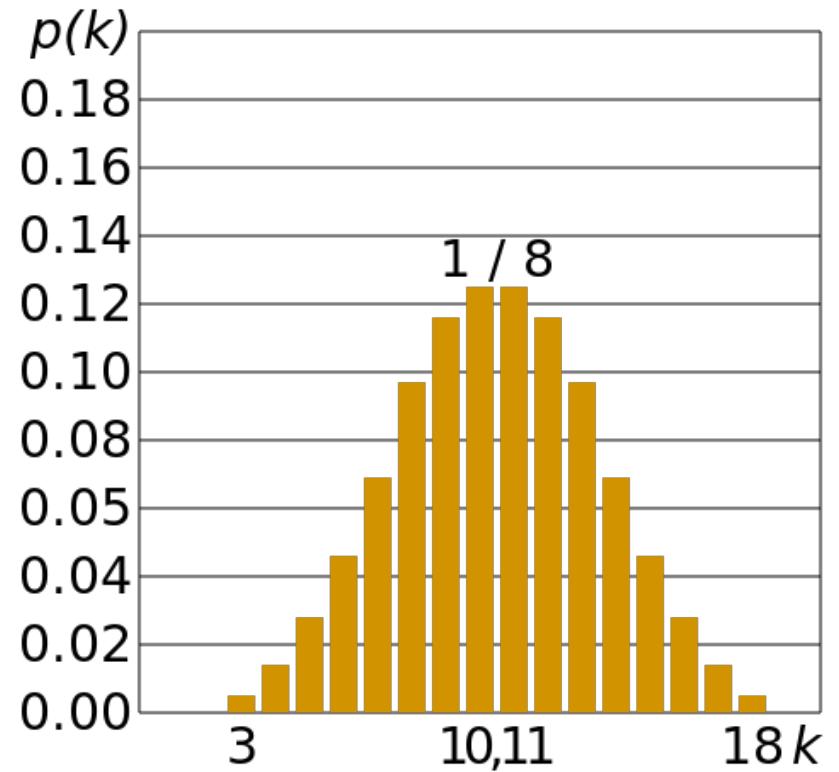
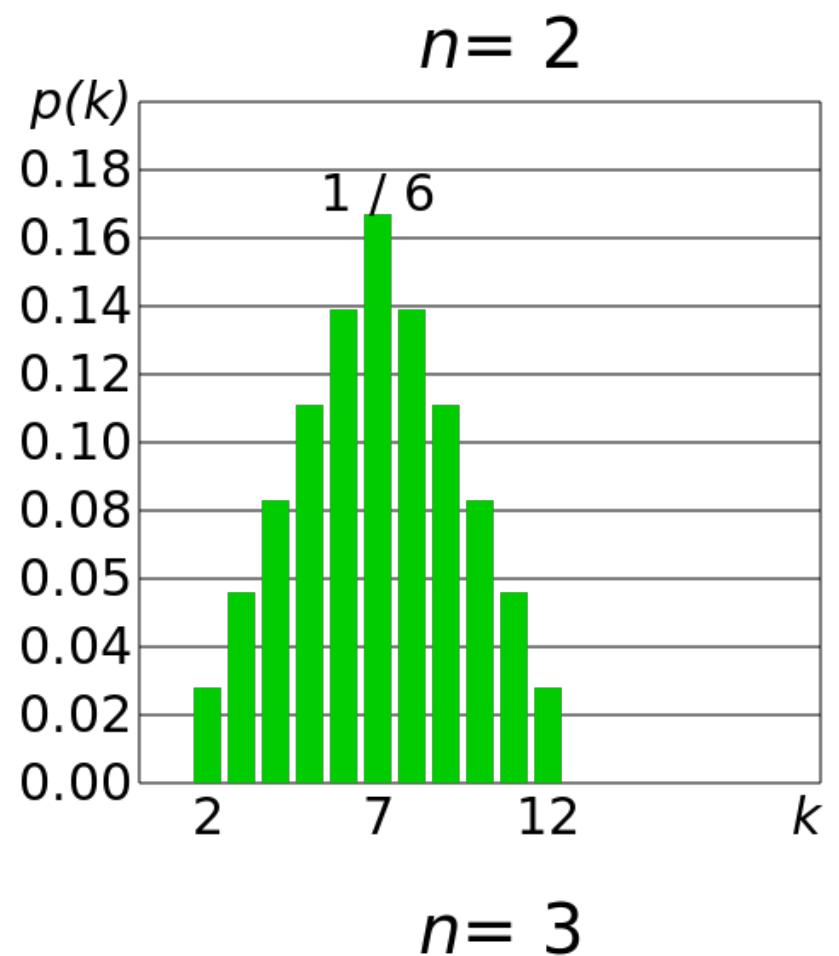
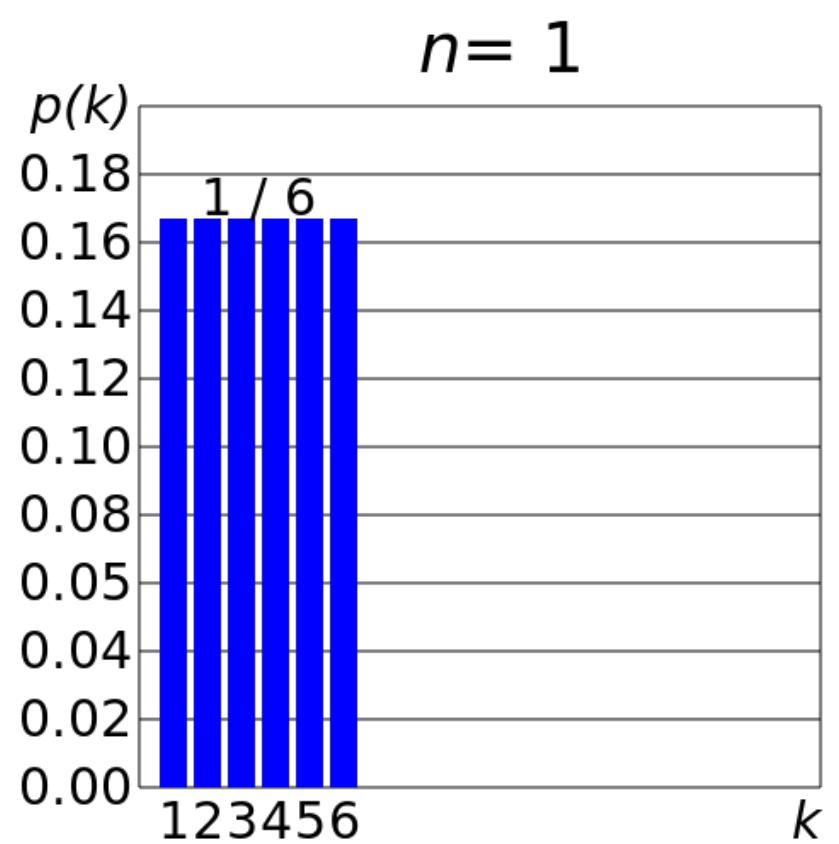


Image credit: Wikipedia

# IS CLT ALWAYS APPLICABLE?

# IS CLT ALWAYS APPLICABLE?

- Some statistics (e.g., median, reciprocal rank, f-score), are not simple averages.

# IS CLT ALWAYS APPLICABLE?

- Some statistics (e.g., median, reciprocal rank, f-score), are not simple averages.
- What should be  $P(X|H_0)$  in those cases?

# NON PARAMETRIC TESTS

# NON PARAMETRIC TESTS

- Can we construct the null distribution of the sample statistic?

# NON PARAMETRIC TESTS

- Can we construct the null distribution of the sample statistic?
- Two broadly applicable methods are:
  - Permutation tests
  - Bootstrap

# PERMUTATION TESTS

- Suppose we want to calculate any difference statistic between two categories
  - E.g. Median difference of % female characters between author genders
- $H_0$  posits no difference (i.e labels do not matter)
- **Idea!** If labels don't matter then repeatedly randomizing the label assignments to examples should yield us the null distribution of the statistic.



Book name

Gone with the  
wind

A tale of two  
cities

War and peace

• • •

• • •

Harry Potter I

Book name	%female characters	Author gender
Gone with the wind	33.4	F
A tale of two cities	45.6	M
War and peace	12.3	M
...	...	...
...	...	...
...	...	...
...	...	...
Harry Potter I	64.1	F

Book name	%female characters	Author gender	Permutation 1
Gone with the wind	33.4	F	M
A tale of two cities	45.6	M	M
War and peace	12.3	M	F
...	...	...	...
...	...	...	...
Harry Potter I	64.1	F	F

Book name	%female characters	Author gender	Permutation 1	Permutation 2
Gone with the wind	33.4	F	M	F
A tale of two cities	45.6	M	M	M
War and peace	12.3	M	F	F
...	...	...	...	...
...	...	...	...	...
Harry Potter I	64.1	F	F	M

Book name	%female characters	Author gender	Permutation 1	Permutation 2	...	Permutation n
Gone with the wind	33.4	F	M	F	...	M
A tale of two cities	45.6	M	M	M	...	F
War and peace	12.3	M	F	F	...	M
...	...	...	...	...	...	...
...	...	...	...	...	...	...
Harry Potter I	64.1	F	F	M	...	F



Book name	%female characters	Author gender	Permutation 1	Permutation 2	...	Permutation n
Gone with the wind	33.4	F	M	F	...	M
A tale of two cities	45.6	M	M	M	...	F
War and peace	12.3	M	F	F	...	M
...	...	...	...	...	...	...
...	...	...	...	...	...	...
...	...	...	...	...	...	...
...	...	...	...	...	...	...
...	...	...	...	...	...	...
Harry Potter I	64.1	F	F	M	...	F

Median  
difference

3.9

-1.4

2.1

...

1.6

Book name	%female characters	Author gender	Permutation 1	Permutation 2	...	Permutation n
Gone with the wind	33.4	F	M	F	...	M
A tale of two cities	45.6	M	M	M	...	F
War and peace	12.3	M	F	F	...	M
...	...	...	...	...	...	...
...	...	...	...	...	...	...
...	...	...	...	...	...	...
...	...	...	...	...	...	...
...	...	...	...	...	...	...
...	...	...	...	...	...	...
Harry Potter I	64.1	F	F	M	...	F

Median difference

3.9

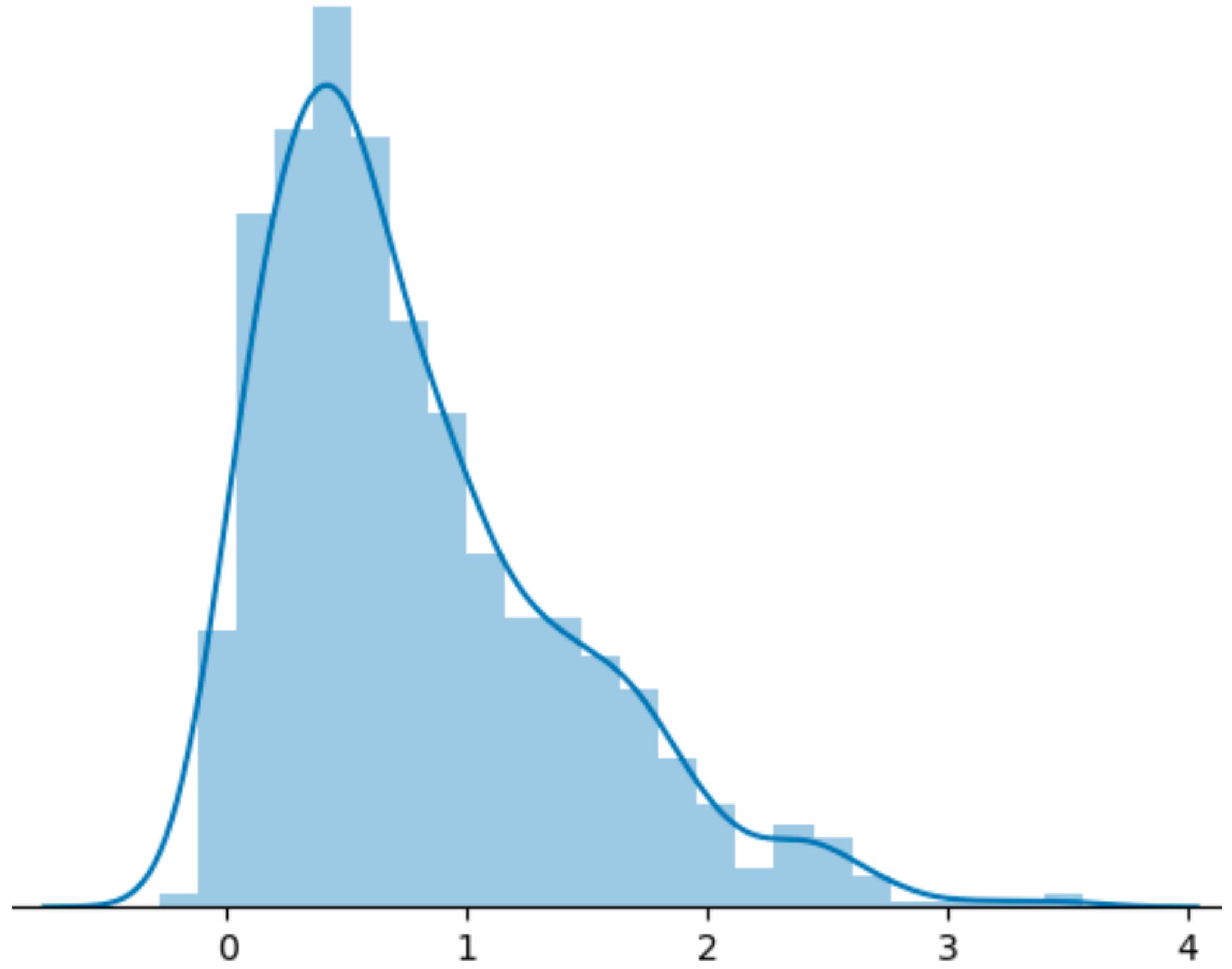
-1.4

2.1

...

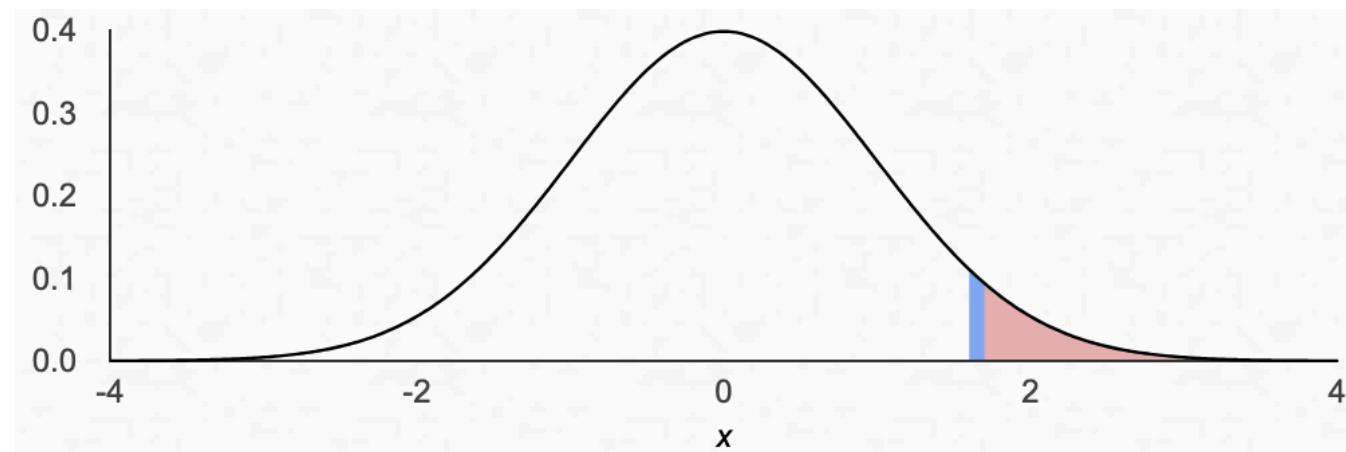
1.6

Is the observed statistic unusual compared to the statistics from all the permutations?

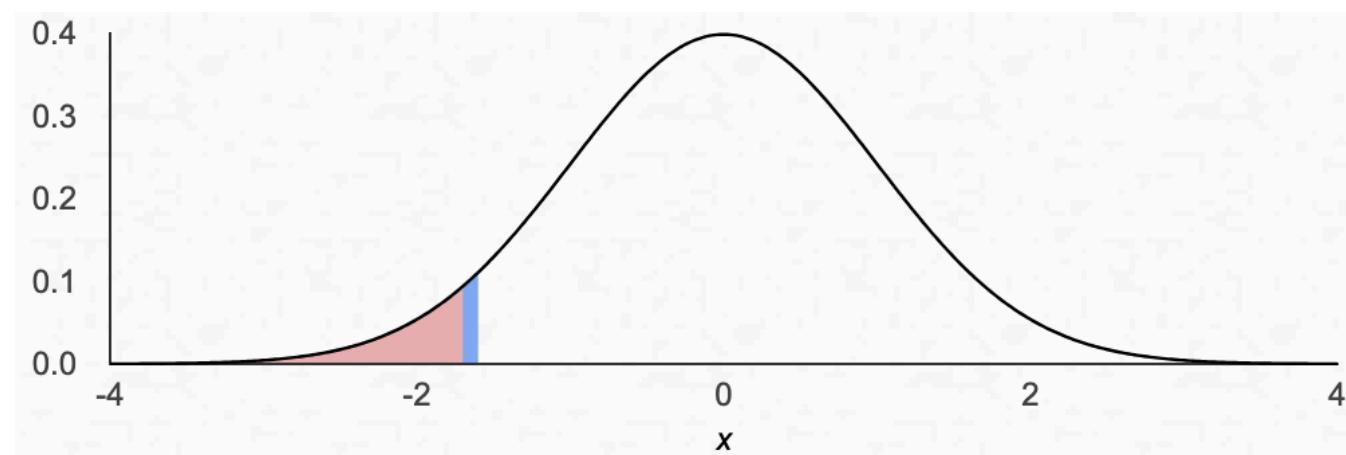


Is 3.9 unusual for this constructed null distribution?

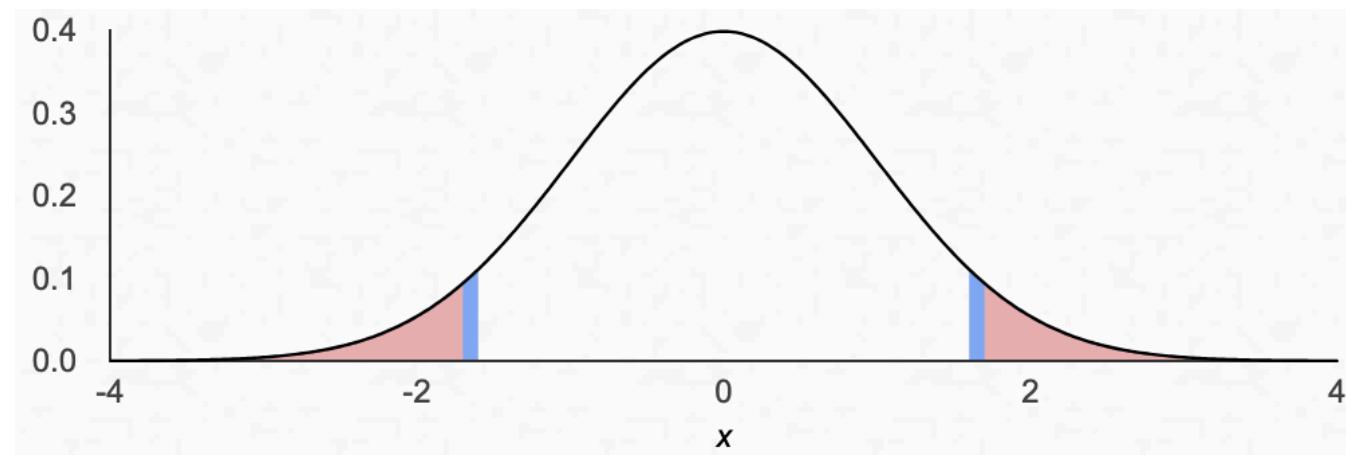
# P VALUE



$$\text{p-value}(x) = P(X \geq x | H_0) = 1 - P(X \leq x | H_0)$$



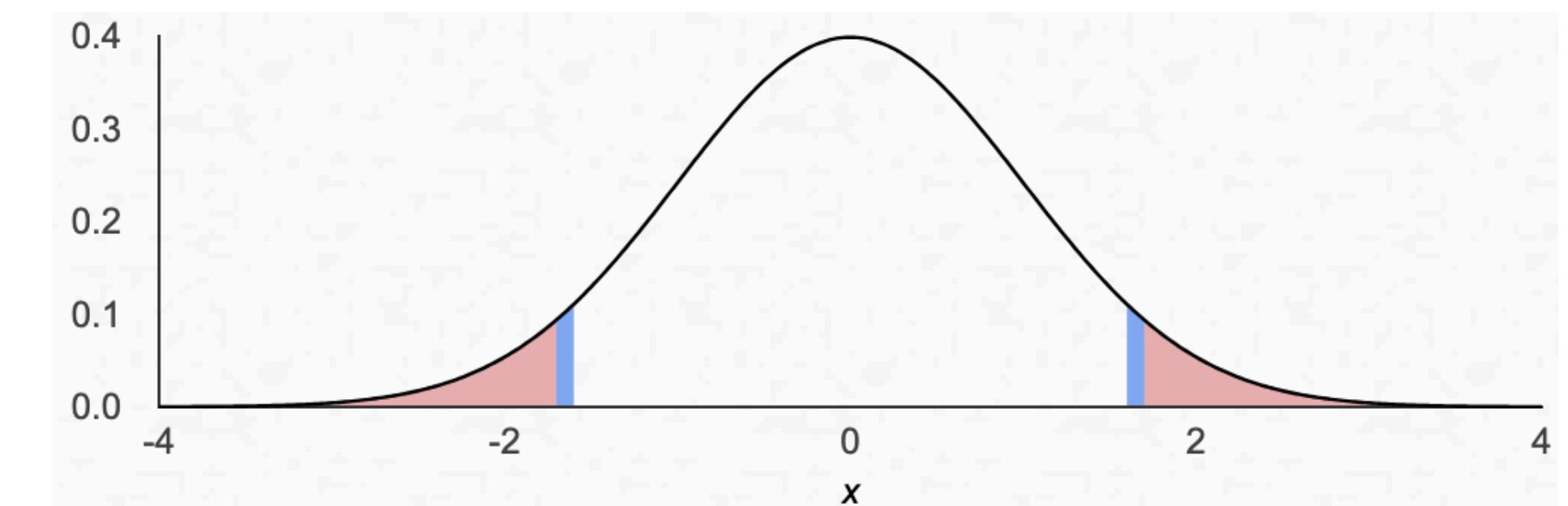
$$\text{p-value}(x) = 1 - P(X \leq x | H_0)$$



$$\text{p-value}(x) = 2 \times P(X \leq -|x| | H_0)$$

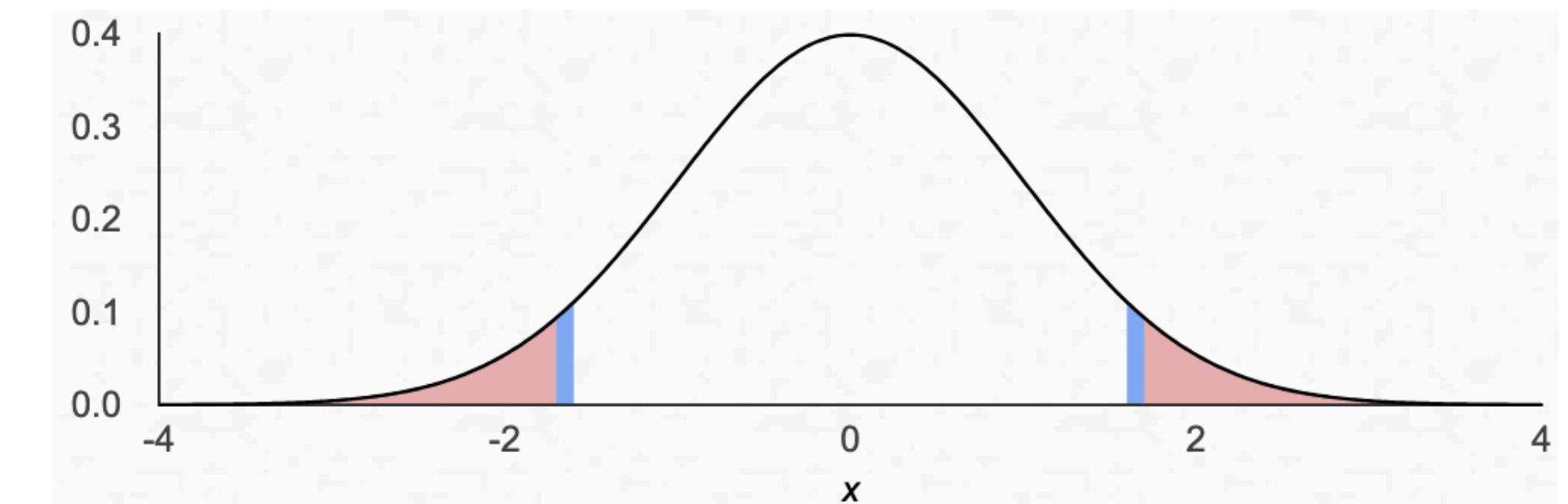
# EMPIRICAL P VALUE

$$p_{\text{empirical}} = \frac{1}{N} \sum_i^N I[abs(m_{\text{obs}}) < abs(m_{\text{perm}}^{(i)})]$$



# EMPIRICAL P VALUE

$$p_{\text{empirical}} = \frac{1}{N} \sum_i^N I[abs(m_{\text{obs}}) < abs(m_{\text{perm}}^{(i)})]$$



- $N$  = number of permutations
- $m_{\text{obs}}$  is the observed value of the statistic
- $m_{\text{perm}}$  is the value of the statistic under permutation
- $I[\cdot]$  is an indicator function

# PERMUTATION TESTS

# PERMUTATION TESTS

- Permutation tests have broad application

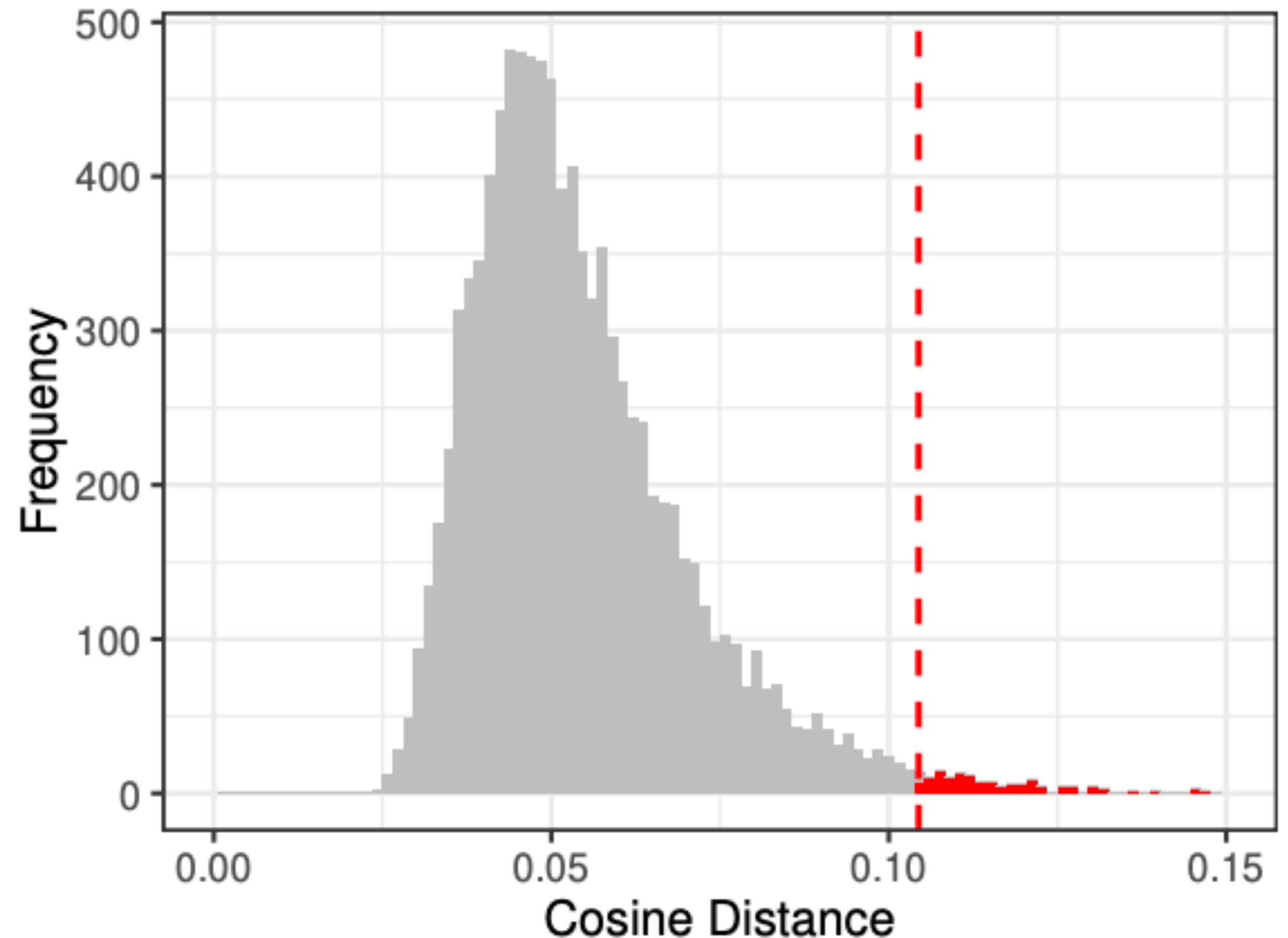
# PERMUTATION TESTS

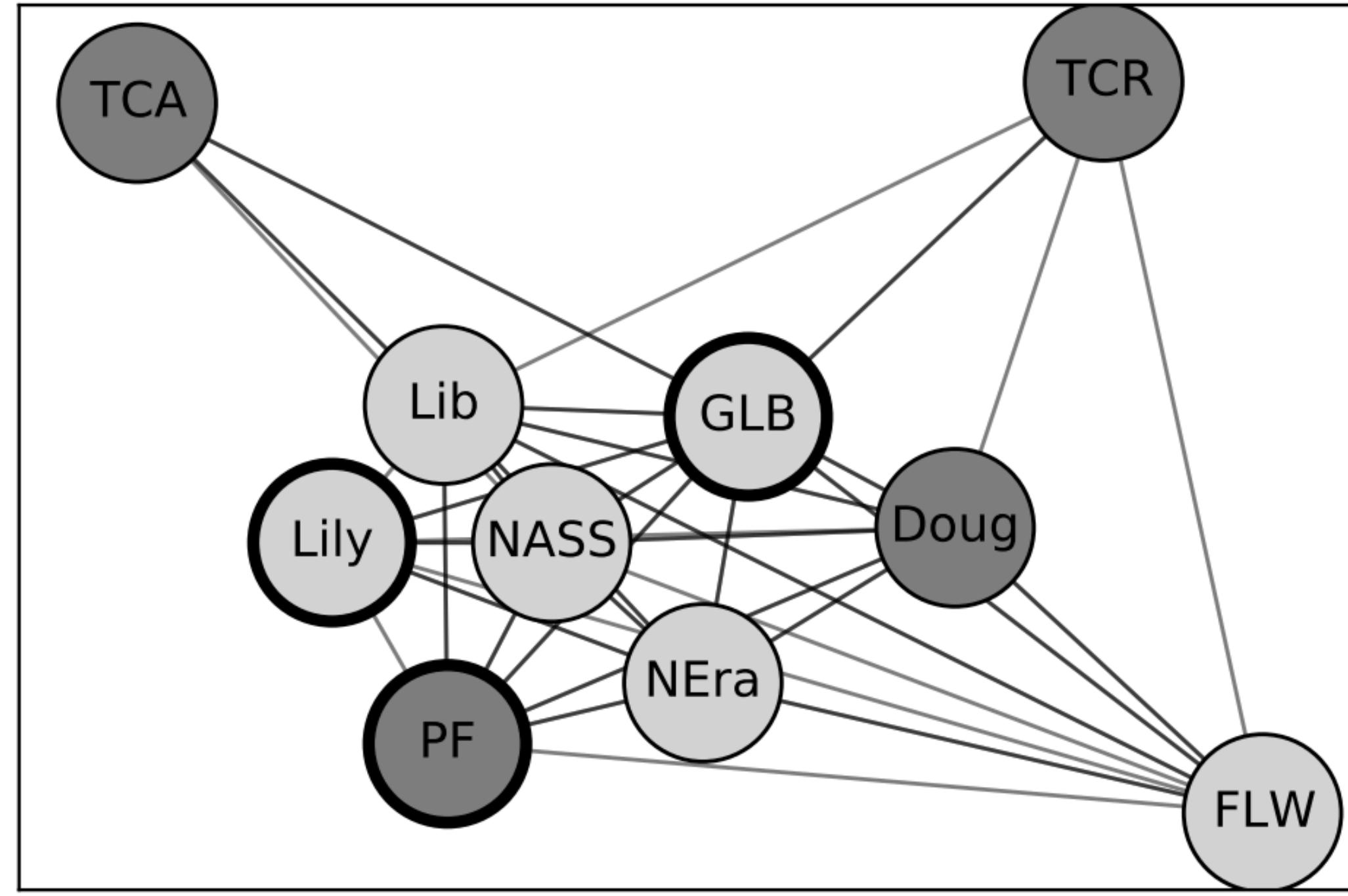
- Permutation tests have broad application
- You don't make any assumptions about the shape and form of the null distribution for the statistic

# PERMUTATION TESTS

- Permutation tests have broad application
- You don't make any assumptions about the shape and form of the null distribution for the statistic
- You can use pretty much any statistic!

- Is there a significant change in meaning of some word (e.g., shovel)?





Doug	DOUGLASS NEWSPAPERS
FLW	FRANK LESLIE'S WEEKLY
GLB	GODEY'S LADY'S BOOK
Lib	THE LIBERATOR
Lily	THE LILY
NASS	NATIONAL ANTI-SLAVERY STANDARD
NEra	THE NATIONAL ERA
PF	THE PROVINCIAL FREEMAN
TCA	THE COLORED AMERICAN
TCR	THE CHRISTIAN RECORDER

- Which newspaper is a consistent leader of some other newspaper?

# BOOTSTRAP

# BOOTSTRAP

- In permutation test, the data is not changed – just the assignment of labels

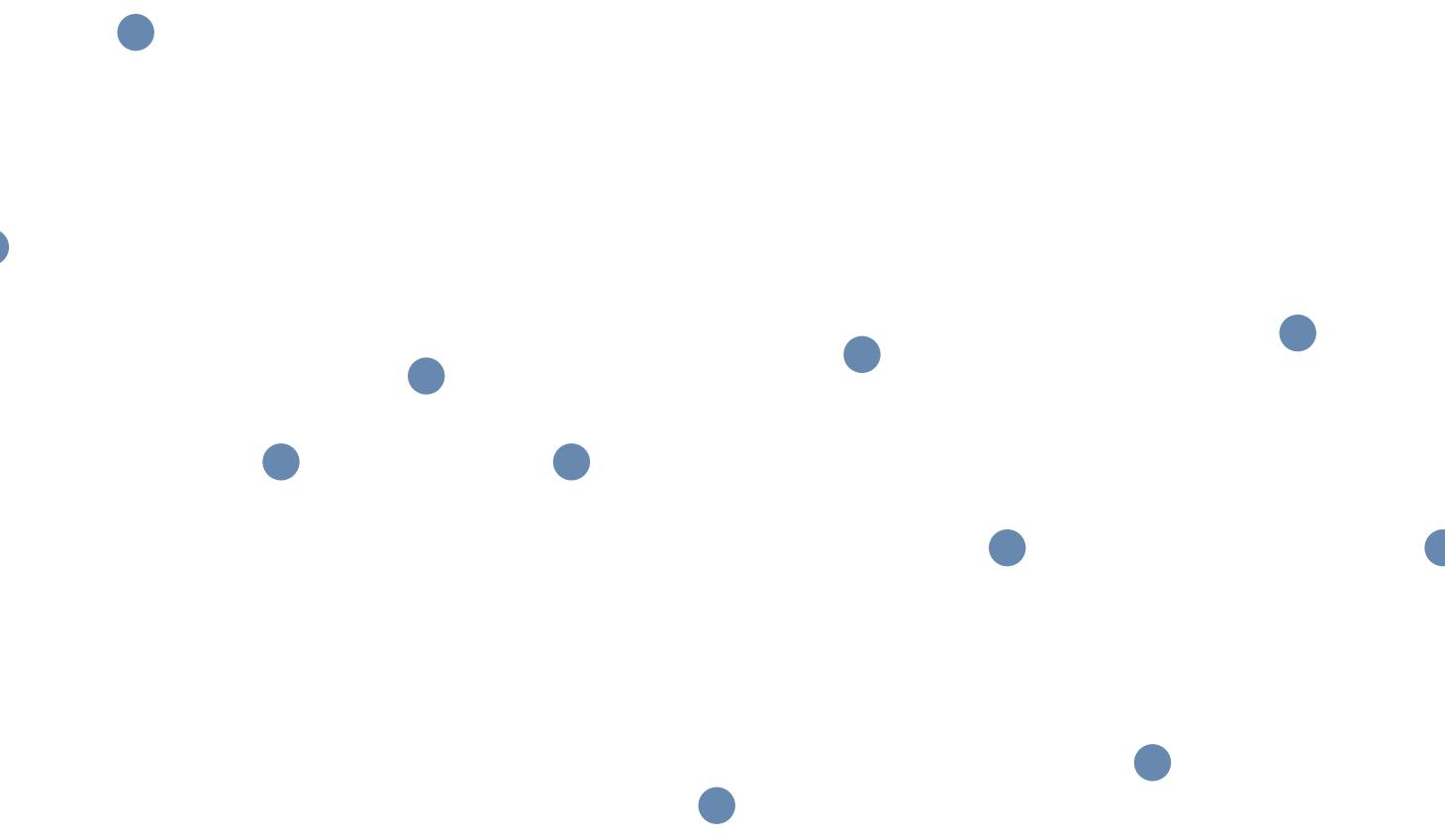
# BOOTSTRAP

- In permutation test, the data is not changed – just the assignment of labels
- The variability in the data can be quantified by constructing hypothetical datasets that follow the same distribution

# BOOTSTRAP

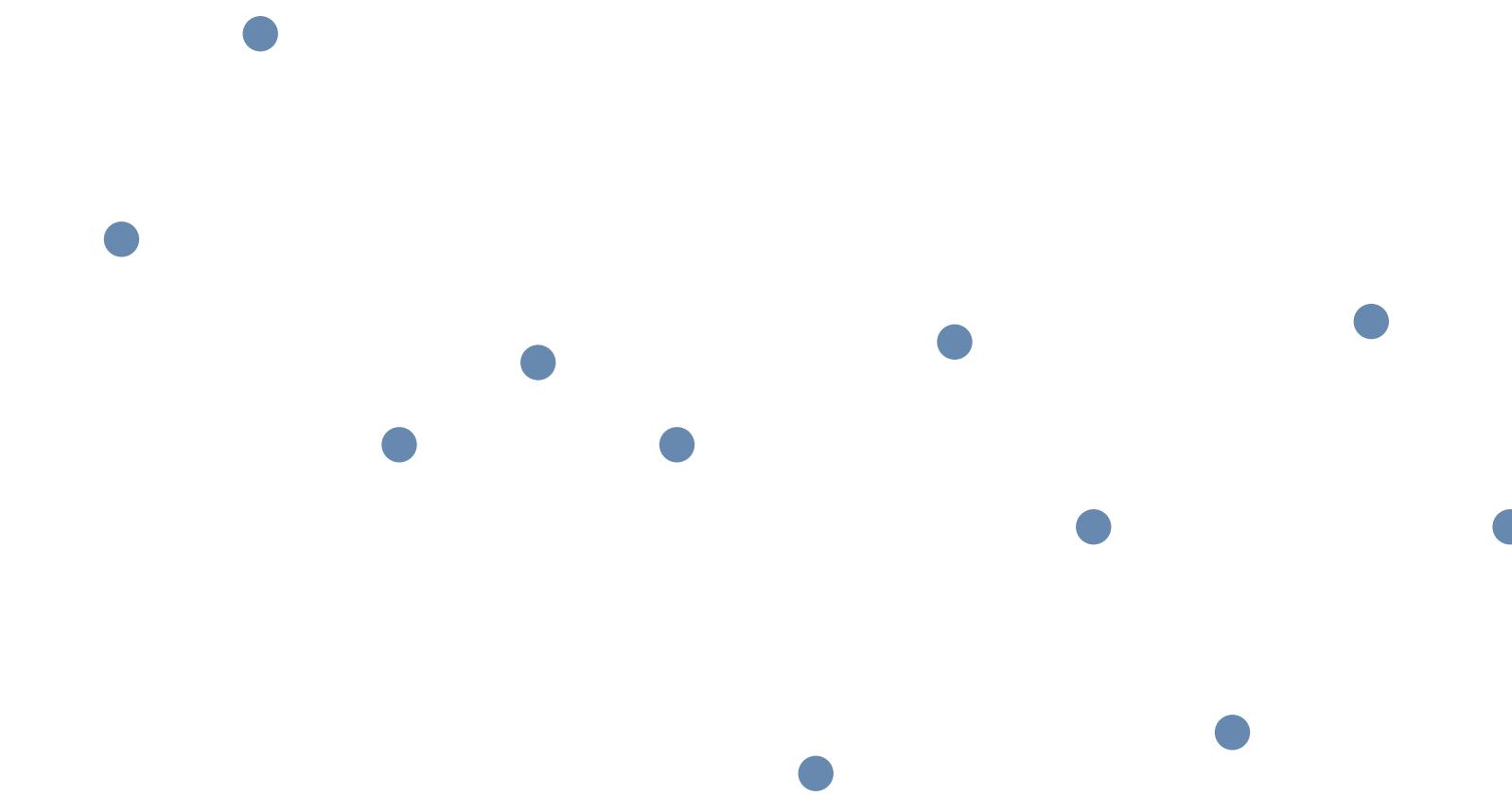
- In permutation test, the data is not changed – just the assignment of labels
- The variability in the data can be quantified by constructing hypothetical datasets that follow the same distribution
- This is the idea of bootstrapping!

# BOOTSTRAP



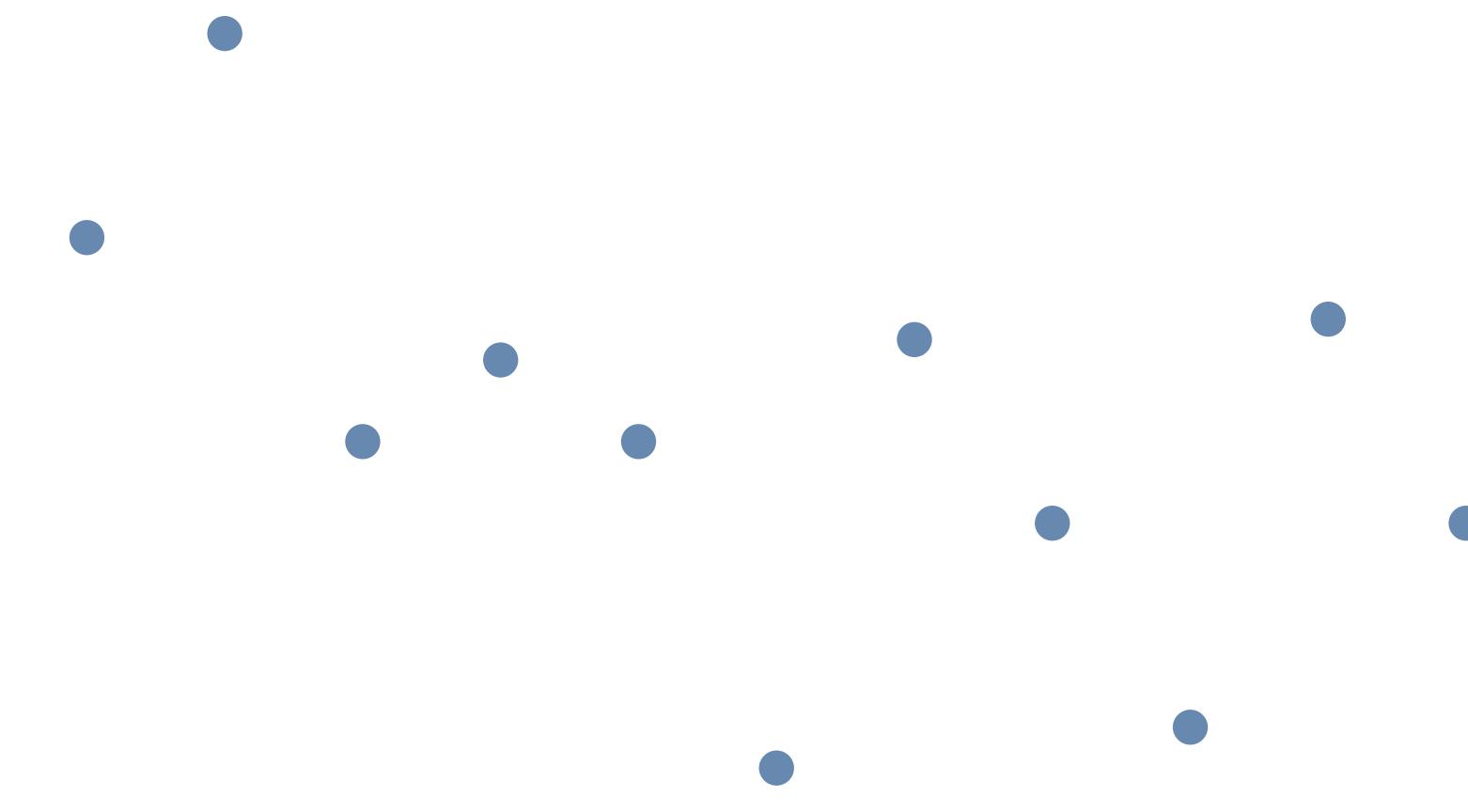
# BOOTSTRAP

- Our sample is assumed to be iid, and represents the population distribution

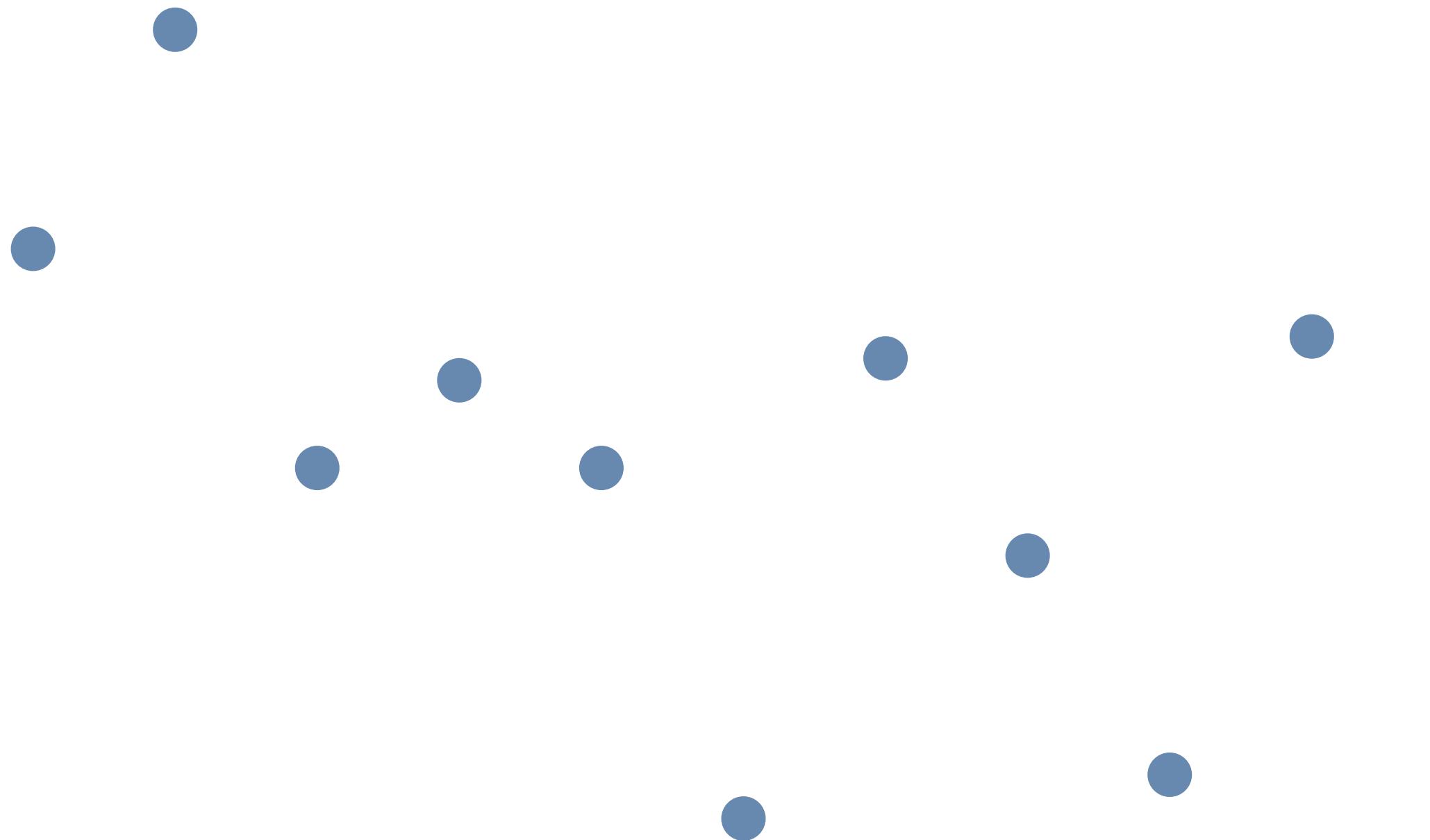


# BOOTSTRAP

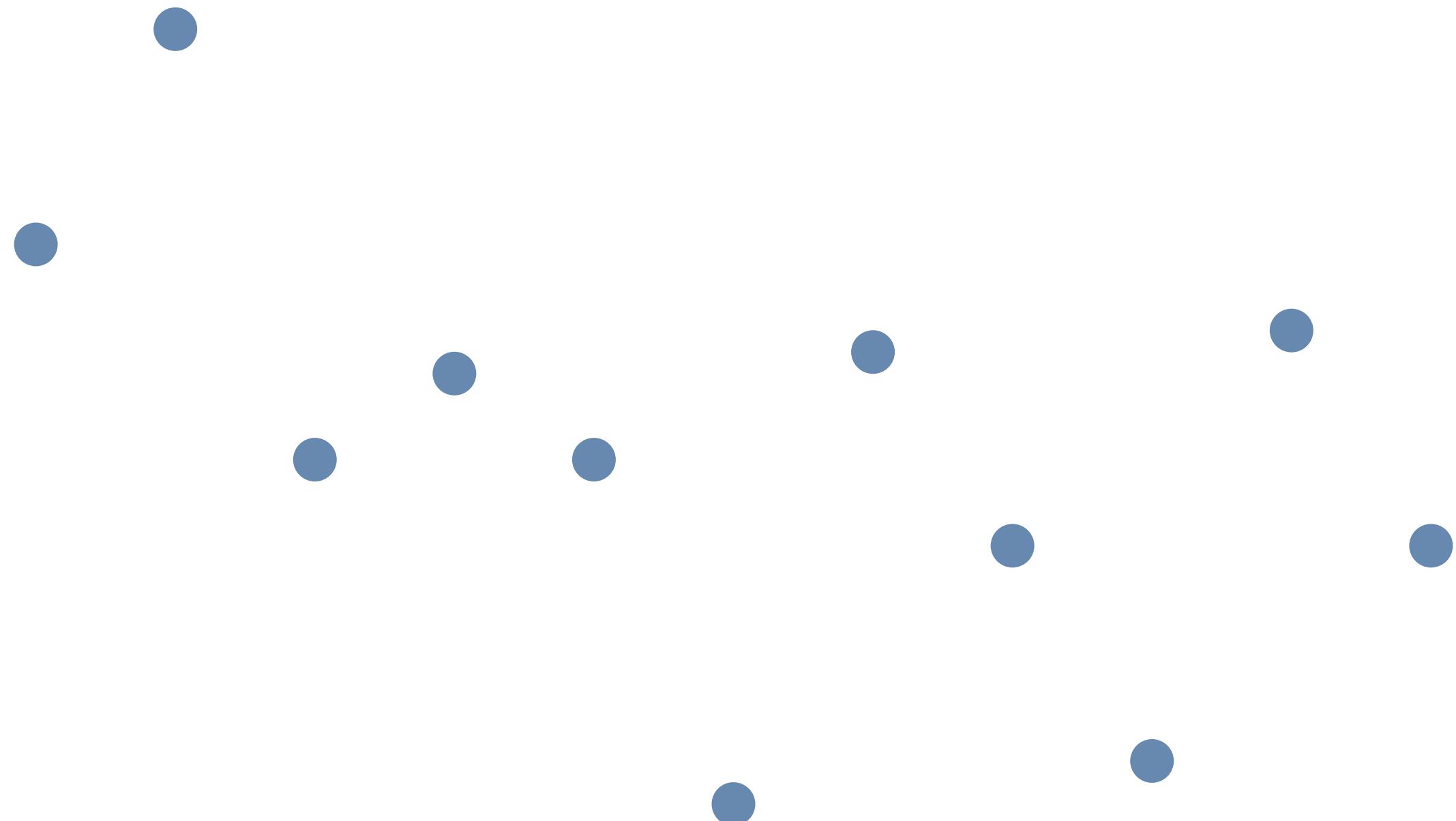
- Our sample is assumed to be iid, and represents the population distribution
- Repeatedly **sample** from our **sample** to generate many alternative datasets



# DISTRIBUTION OF OUR SAMPLE

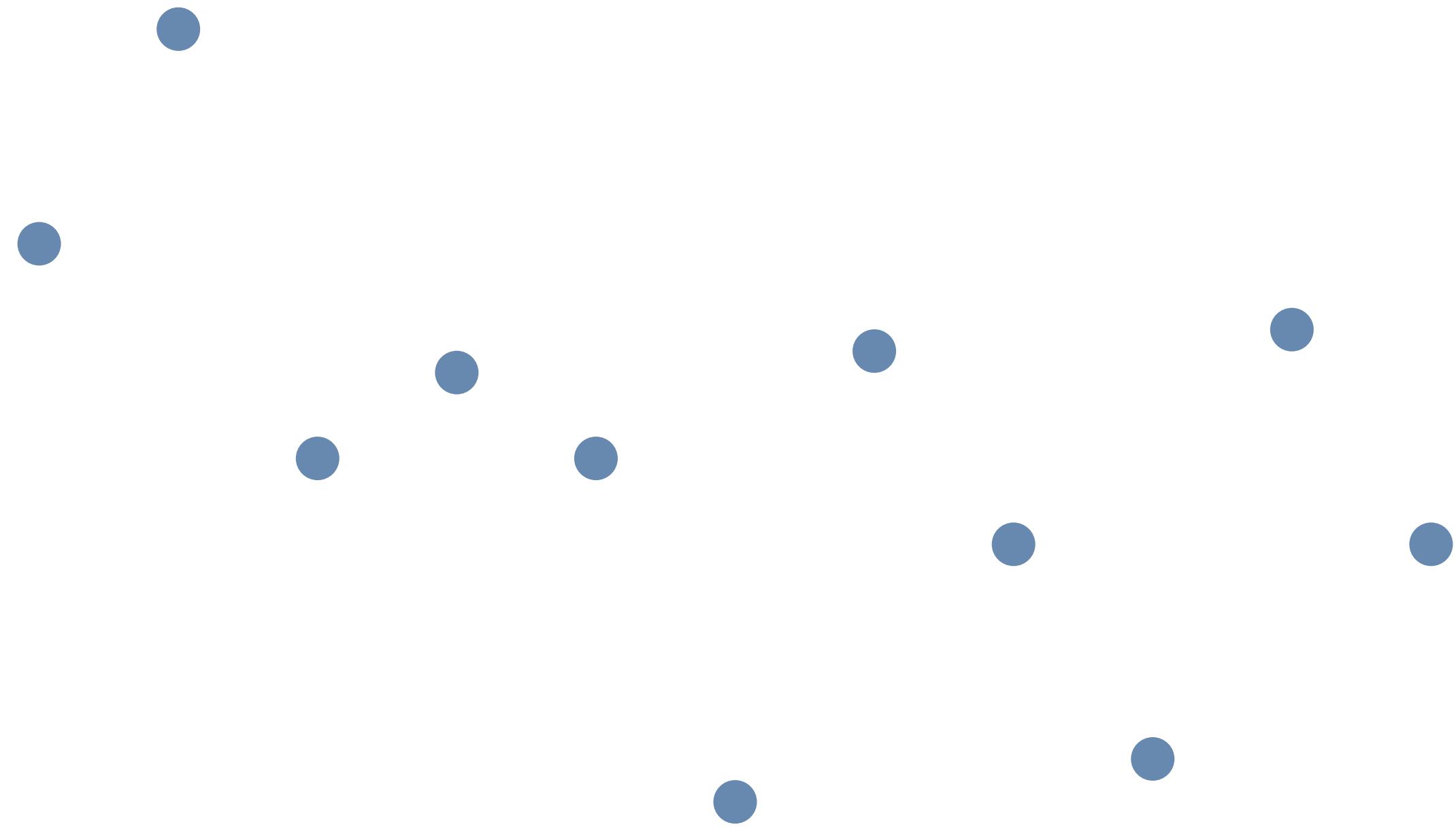


# DISTRIBUTION OF OUR SAMPLE



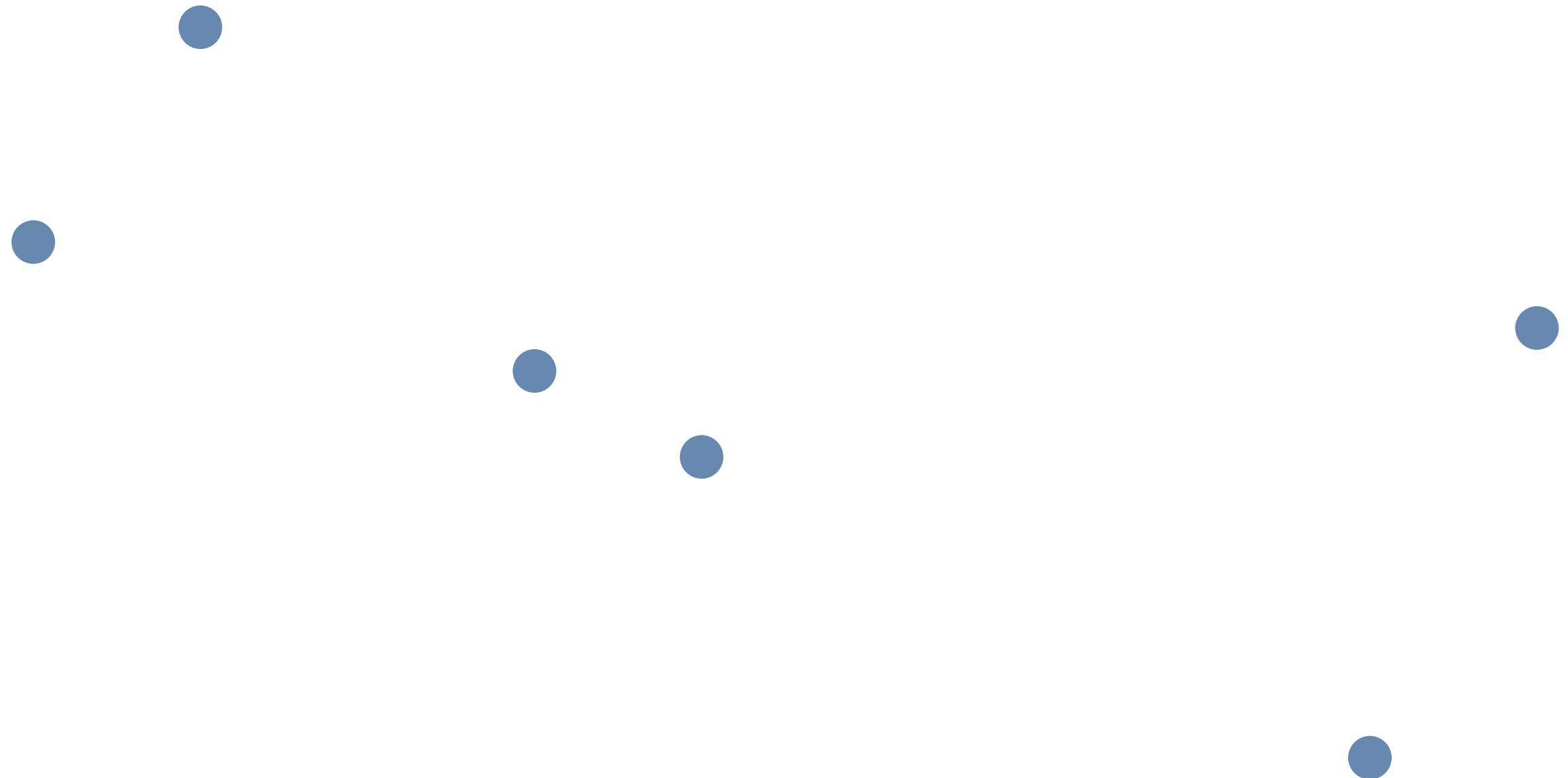
- We can calculate some metric on this sample such as accuracy

# DISTRIBUTION OF OUR SAMPLE

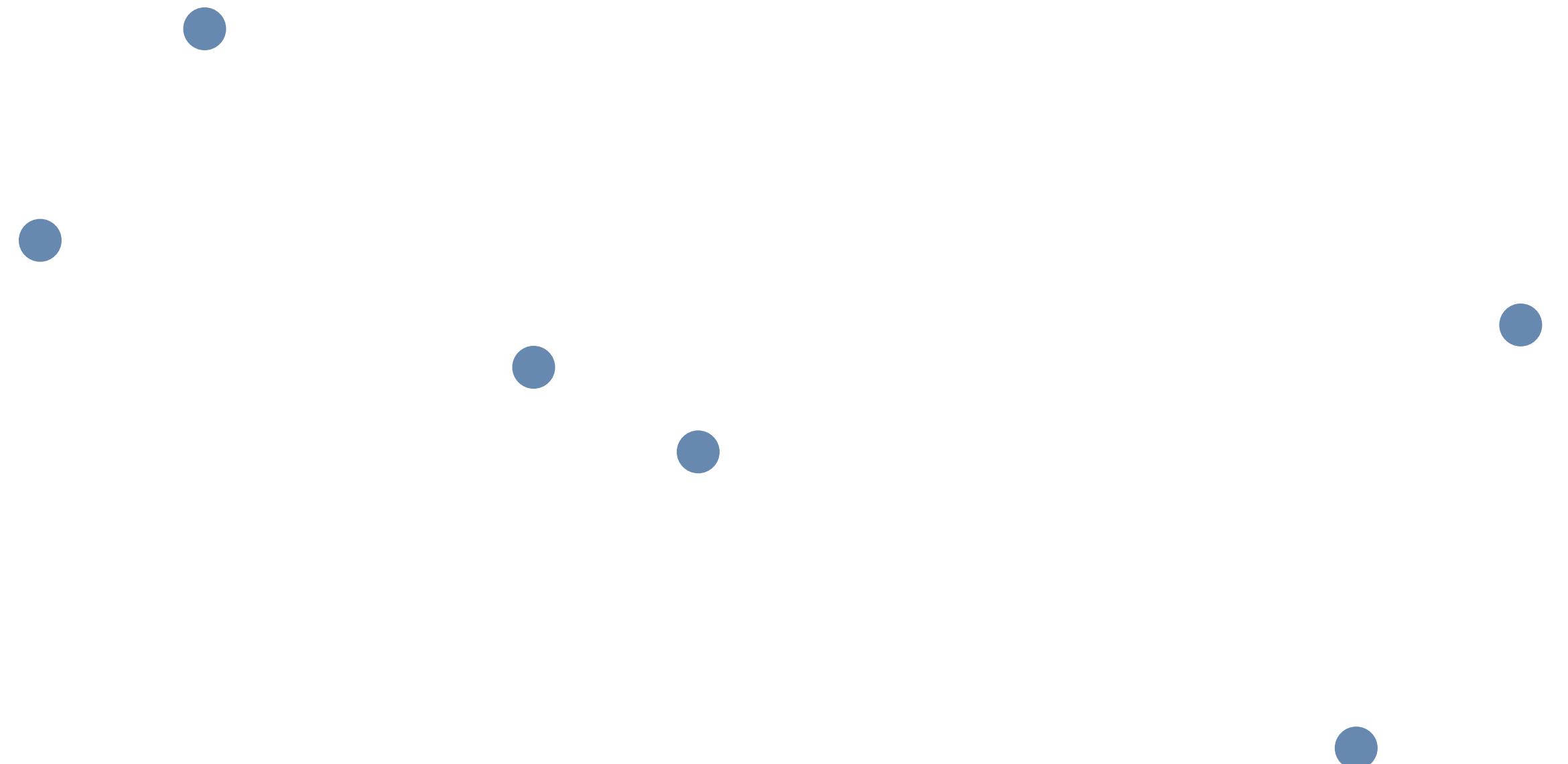


- We can calculate some metric on this sample such as accuracy
- $\text{acc}_{\text{obs}}$

# BOOTSTRAP 1

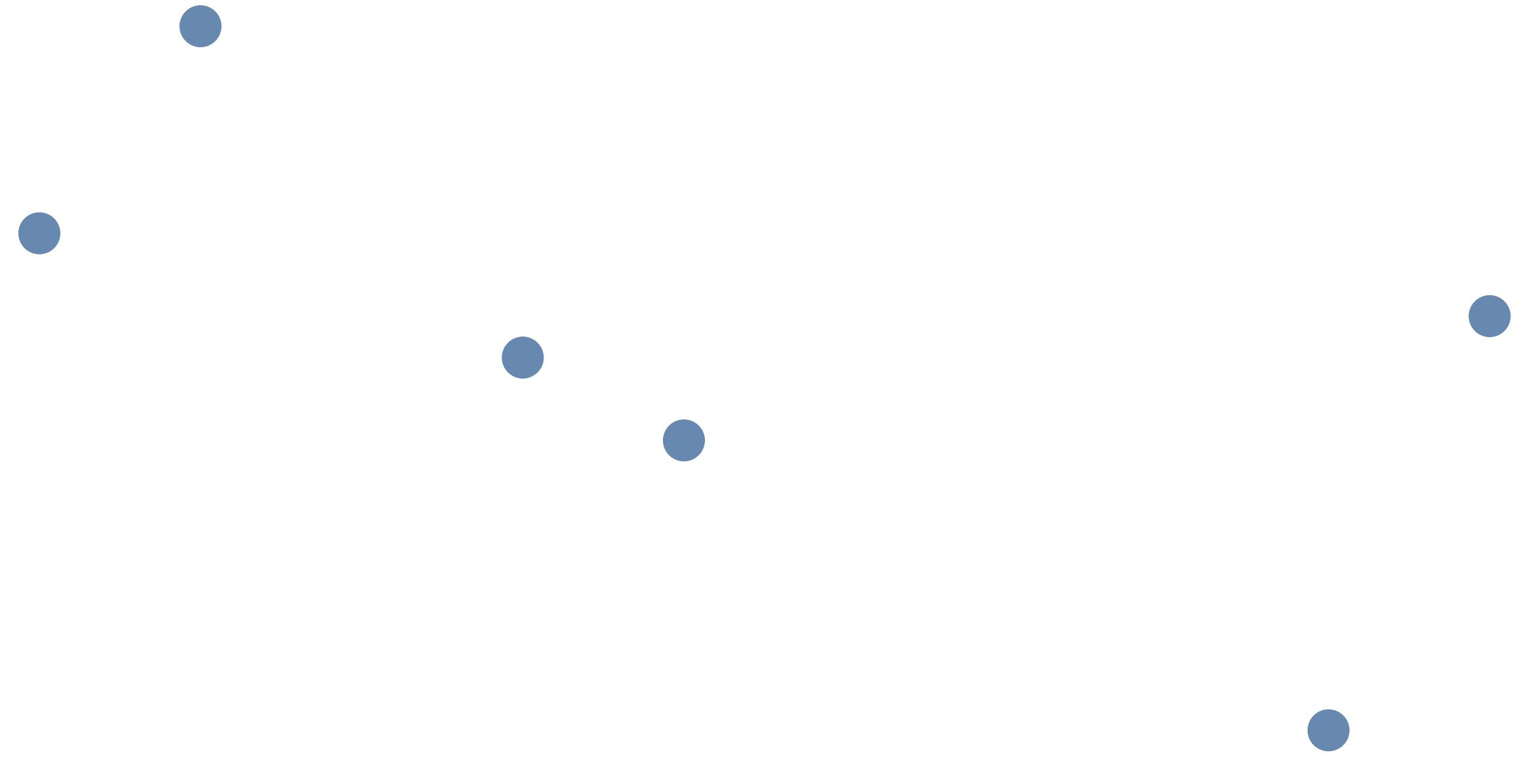


# BOOTSTRAP 1



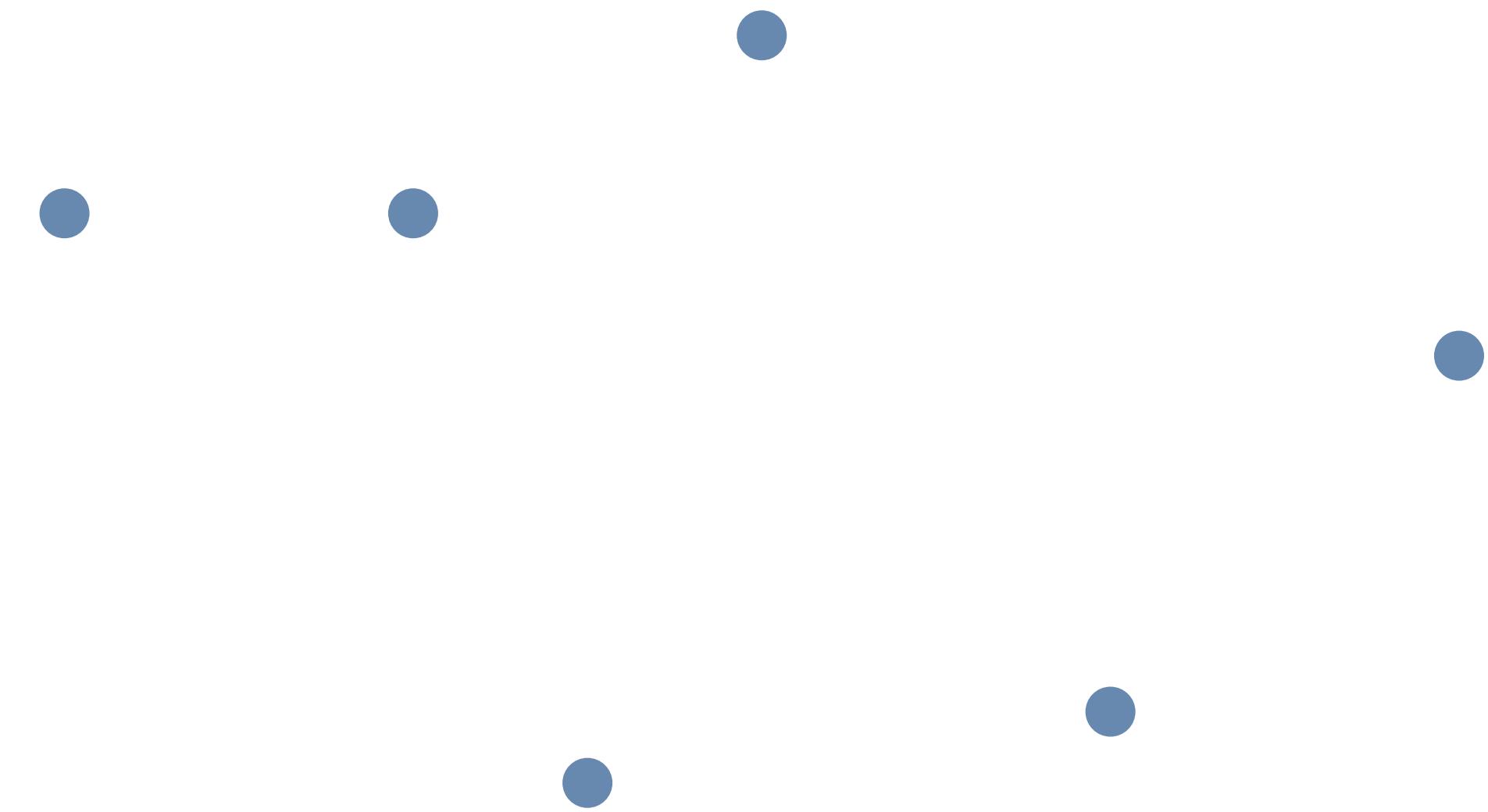
- Calculate accuracy on this bootstrapped sample

# BOOTSTRAP 1

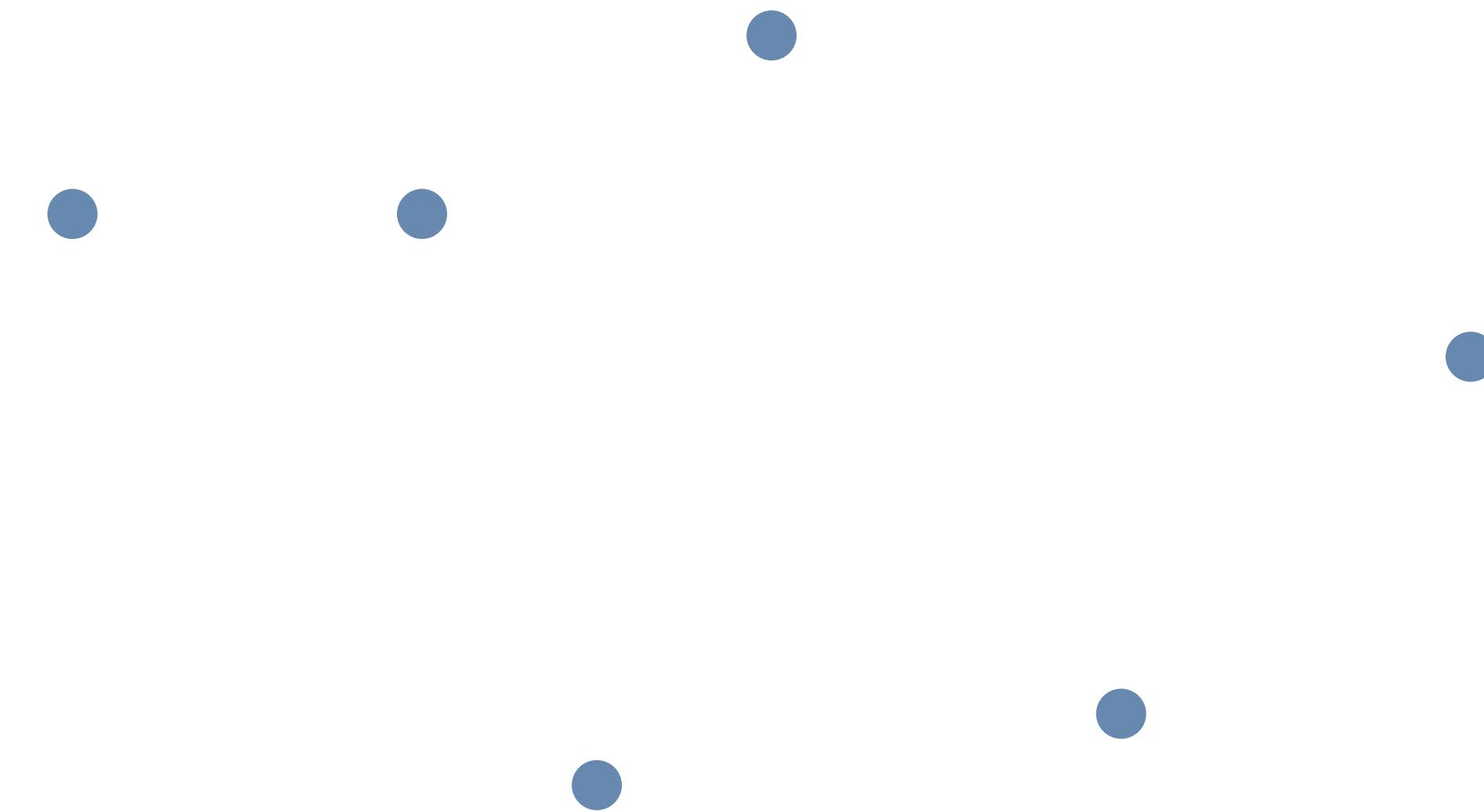


- Calculate accuracy on this bootstrapped sample
- $\text{acc}_{\text{boot}}$

# BOOTSTRAP 2



# BOOTSTRAP 2



- Again calculate the accuracy

# BOOTSTRAP 2

- Again calculate the accuracy
- $\text{acc}_{\text{boot}}$

**AND SO ON**

# BOOTSTRAPPED CONFIDENCE INTERVALS

# BOOTSTRAPPED CONFIDENCE INTERVALS

- From  $b$  bootstrap samples, we can get a size  $b$  array as estimates for accuracy

# BOOTSTRAPPED CONFIDENCE INTERVALS

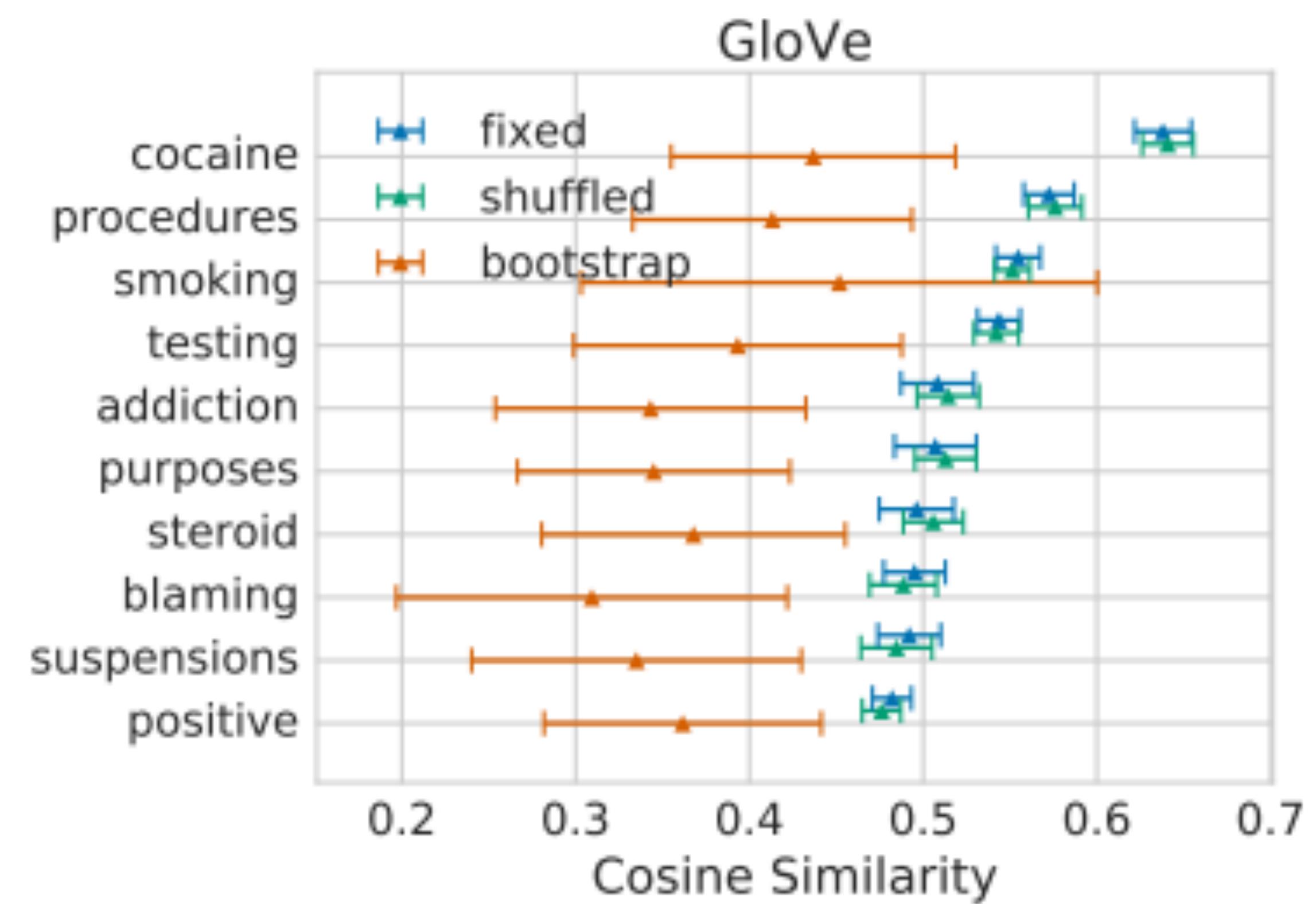
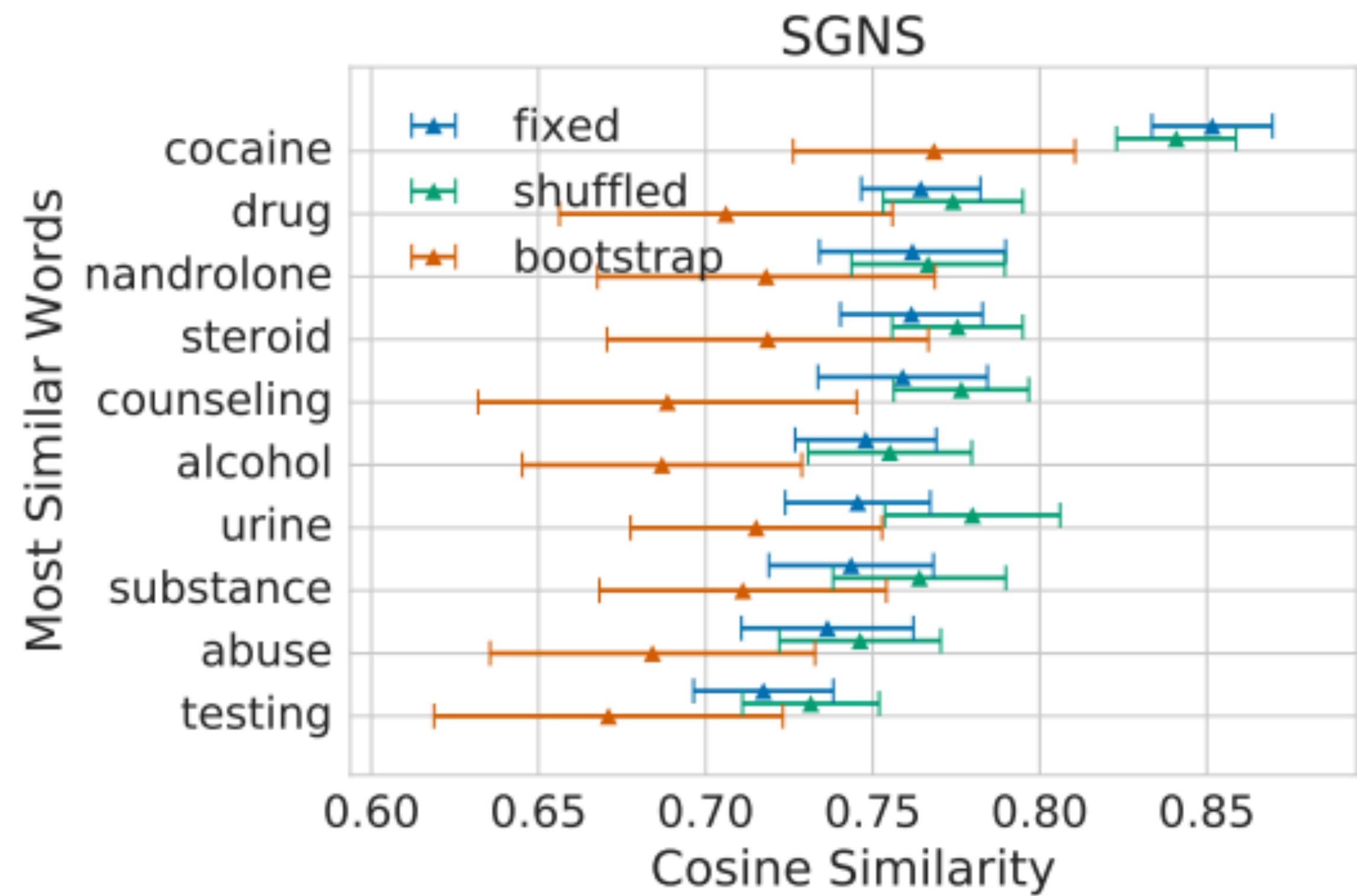
- From  $b$  bootstrap samples, we can get a size  $b$  array as estimates for accuracy
- We can rank these estimates and find the middle 95% to define a range in which the accuracy estimate falls

# BOOTSTRAPPED CONFIDENCE INTERVALS

- From  $b$  bootstrap samples, we can get a size  $b$  array as estimates for accuracy
- We can rank these estimates and find the middle 95% to define a range in which the accuracy estimate falls
- 95% CI is [2.5,97.5] percentile

# BOOTSTRAPPED CONFIDENCE INTERVALS

- From  $b$  bootstrap samples, we can get a size  $b$  array as estimates for accuracy
- We can rank these estimates and find the middle 95% to define a range in which the accuracy estimate falls
- 95% CI is [2.5,97.5] percentile
- For large  $b$  (e.g.,  $b=1000$ ), this gives a tight bound for the CI



Which words are most similar to marijuana?

# SUMMARY

# SUMMARY

- It's often a good idea to give a confidence interval for your estimated metric

# SUMMARY

- It's often a good idea to give a confidence interval for your estimated metric
- Statistical tests are useful to verify claims

# SUMMARY

- It's often a good idea to give a confidence interval for your estimated metric
- Statistical tests are useful to verify claims
- Use ablation testing to assess the importance of model components

# IN CLASS

- Parametric test
- Non parametric test