



# NON PARAMETRIC TESTING

Sandeep Soni

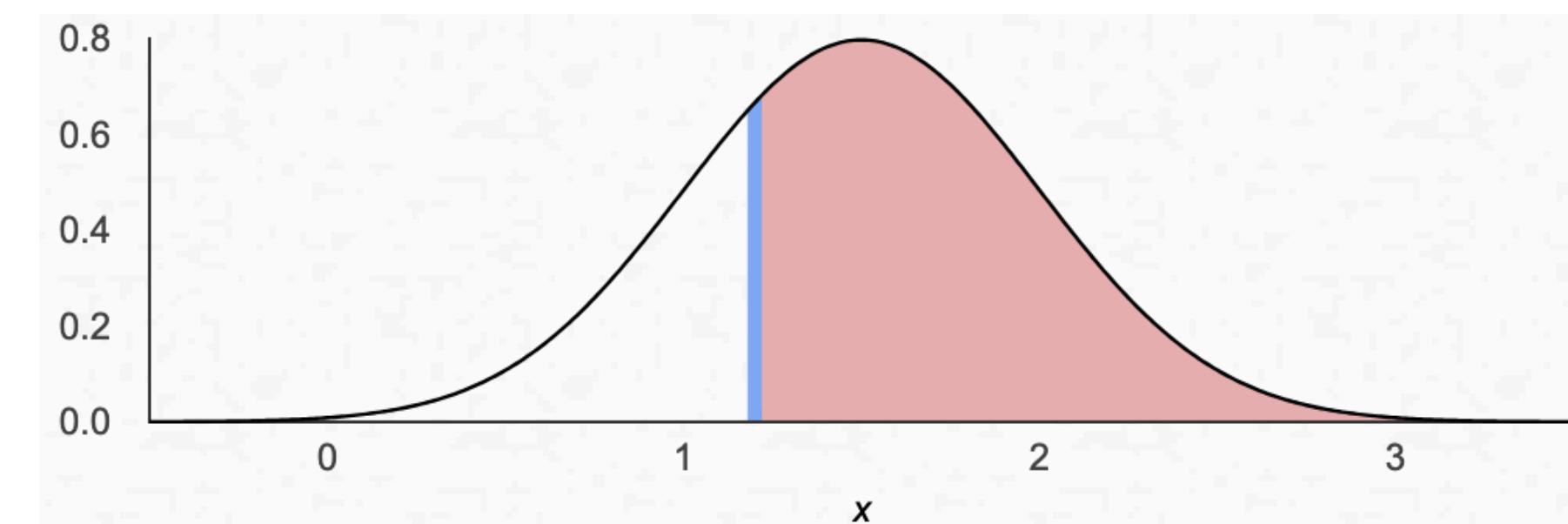
---

02/15/2024

# STORY SO FAR

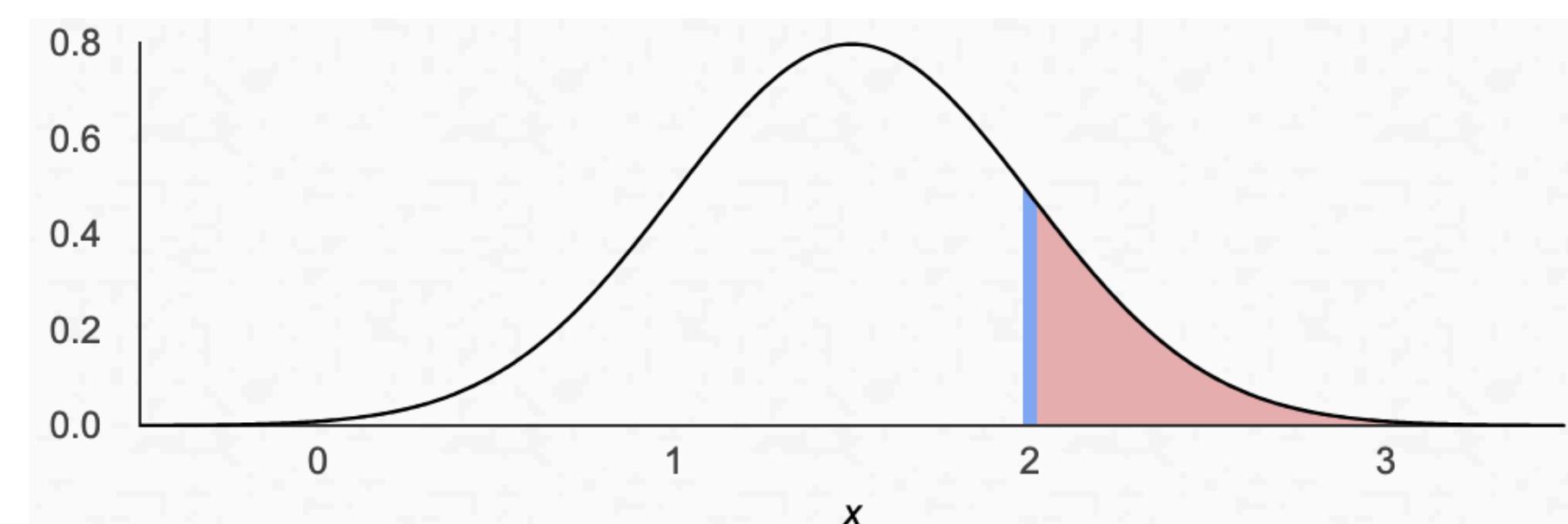
Which observation would you say is surprising if the null distribution holds?

x=1.25



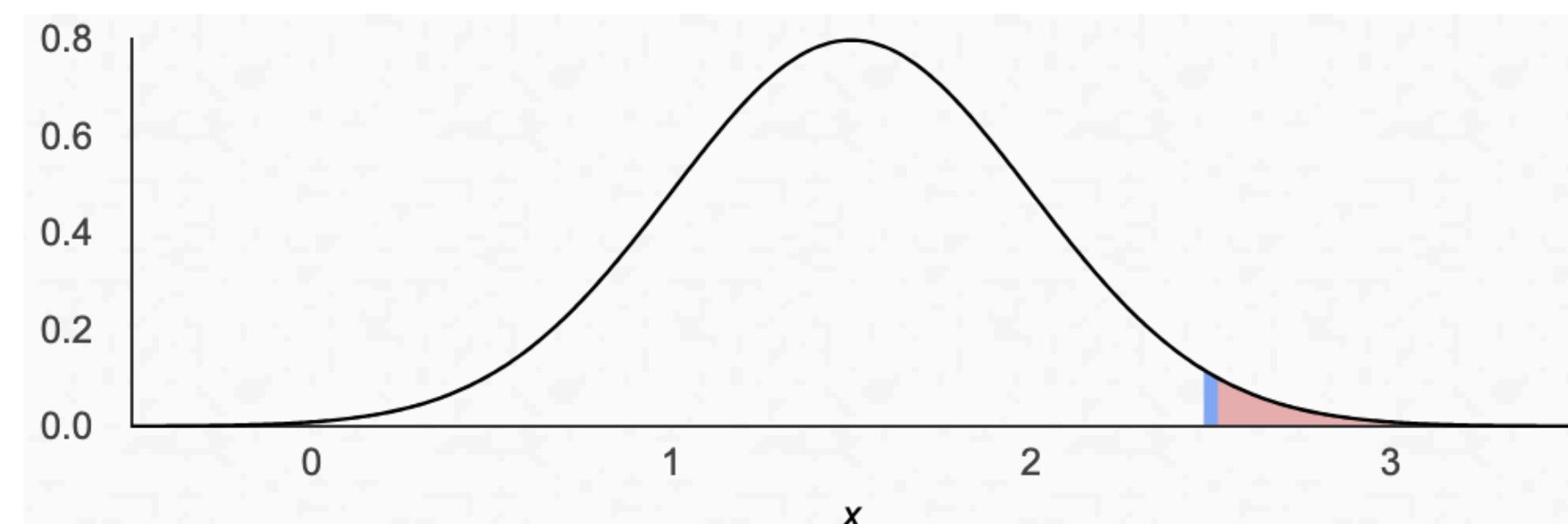
p=0.69

x=2



p=0.16

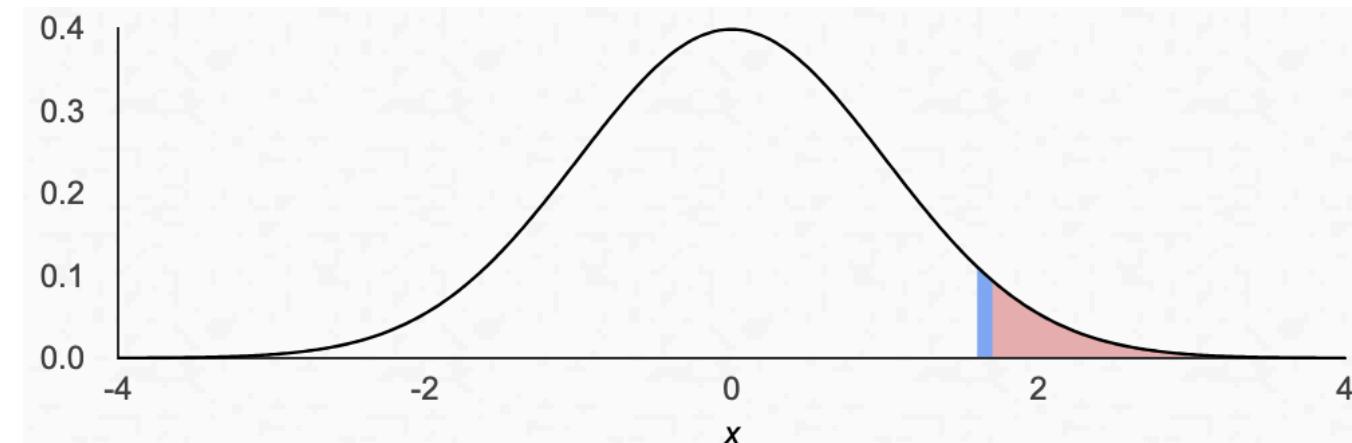
x=2.5



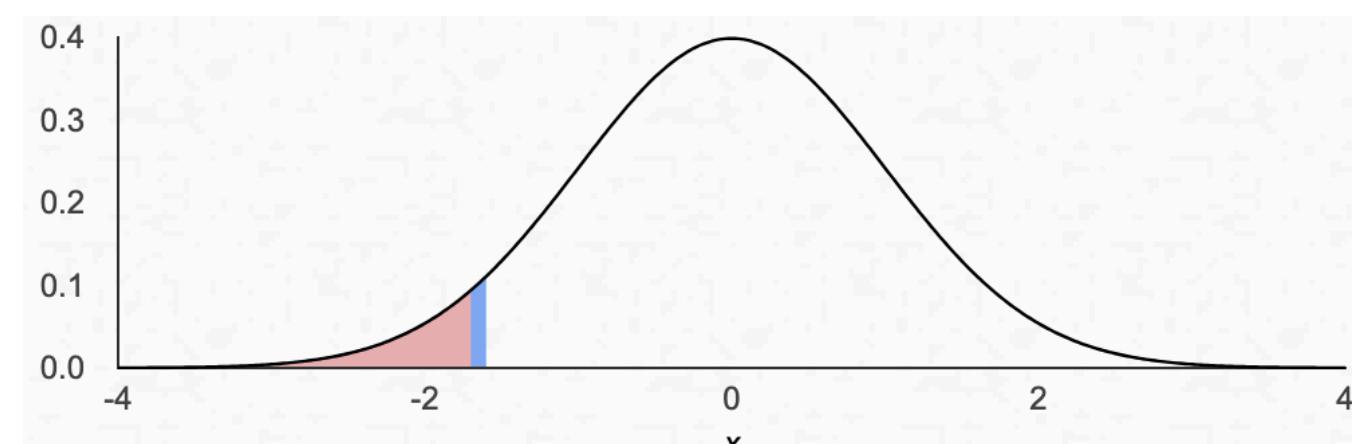
p=0.02

# P VALUES

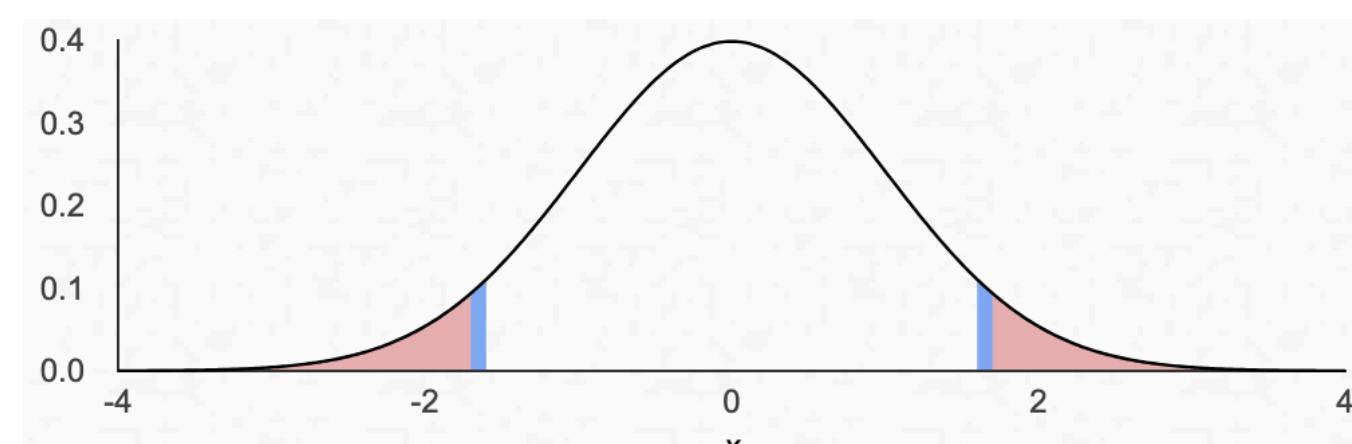
A p value is the probability of observing a statistic at least as extreme as the one we did if the null hypothesis were true.



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



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



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

# RECIPE FOR HYPOTHESIS TESTING

# RECIPE FOR HYPOTHESIS TESTING

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

- In these tests, we can calculate the probabilities by plugging it in an equation

# PARAMETRIC TESTS

- In these tests, we can calculate the probabilities by plugging it in an equation
- For example, if we assume 100 trials of a fair coin, and if we observe 75 heads, then calculate  $P(x=75 | p=0.5, n=100)$  using a binomial distribution's

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

# CENTRAL LIMIT THEOREM

# CENTRAL LIMIT THEOREM

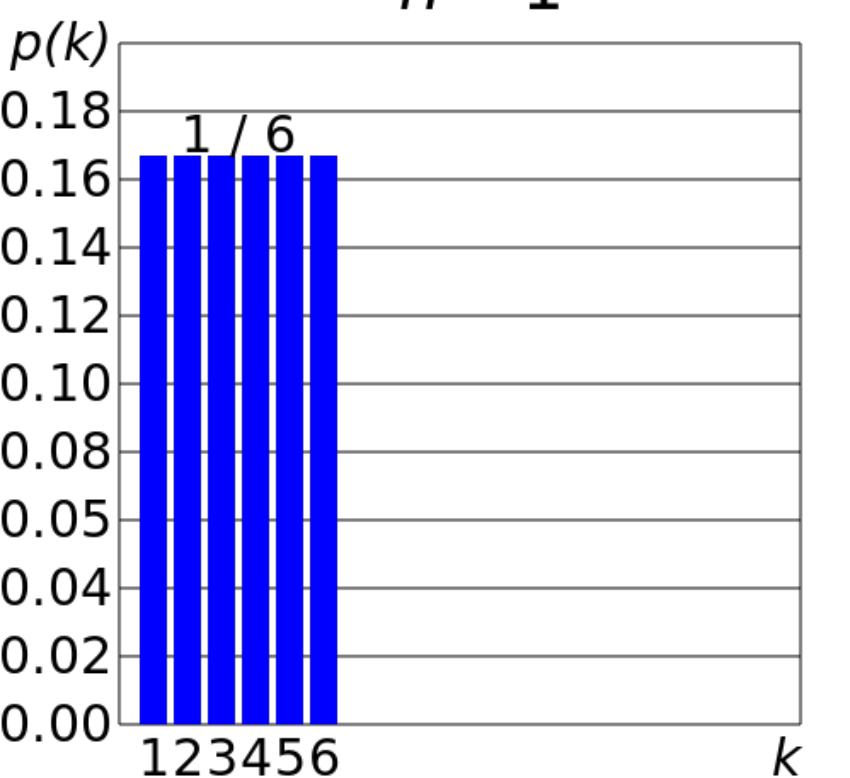
- Often our statistic is calculated by averaging over multiple independent instances (e.g., accuracy)

# CENTRAL LIMIT THEOREM

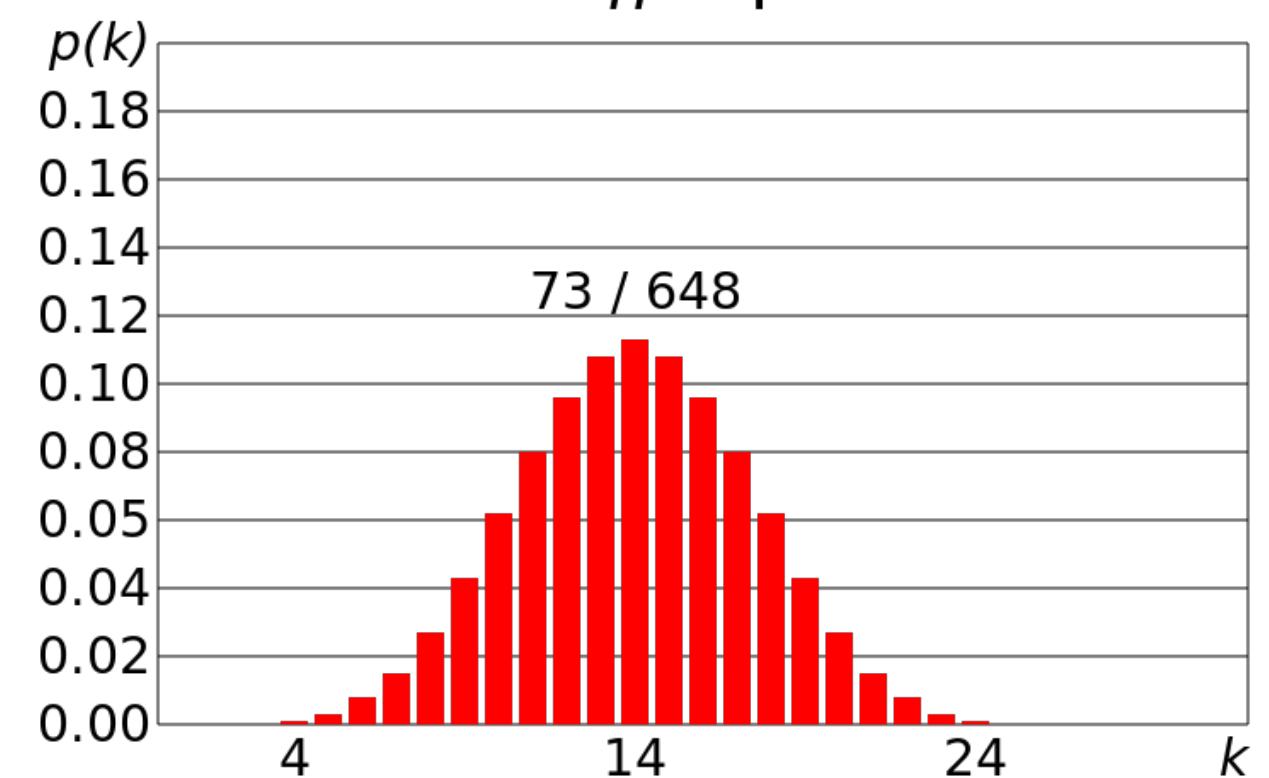
- Often our statistic is calculated by averaging over multiple independent instances (e.g., accuracy)
- According to CLT, the average of independent and identically distributed random variables tends to be a normal distribution, even if the random variables are not normally distributed



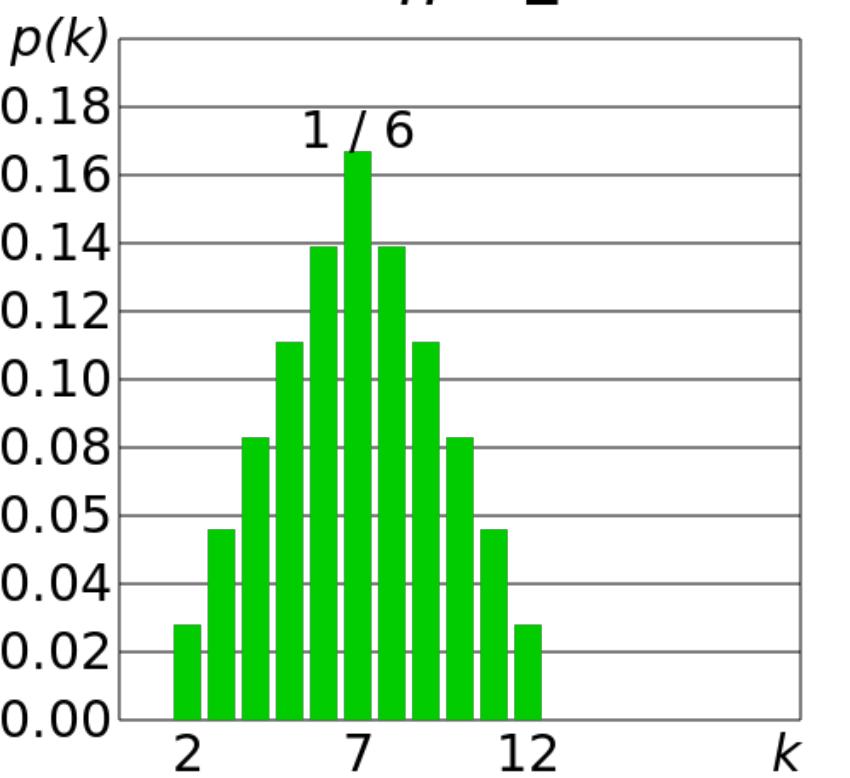
$n = 1$



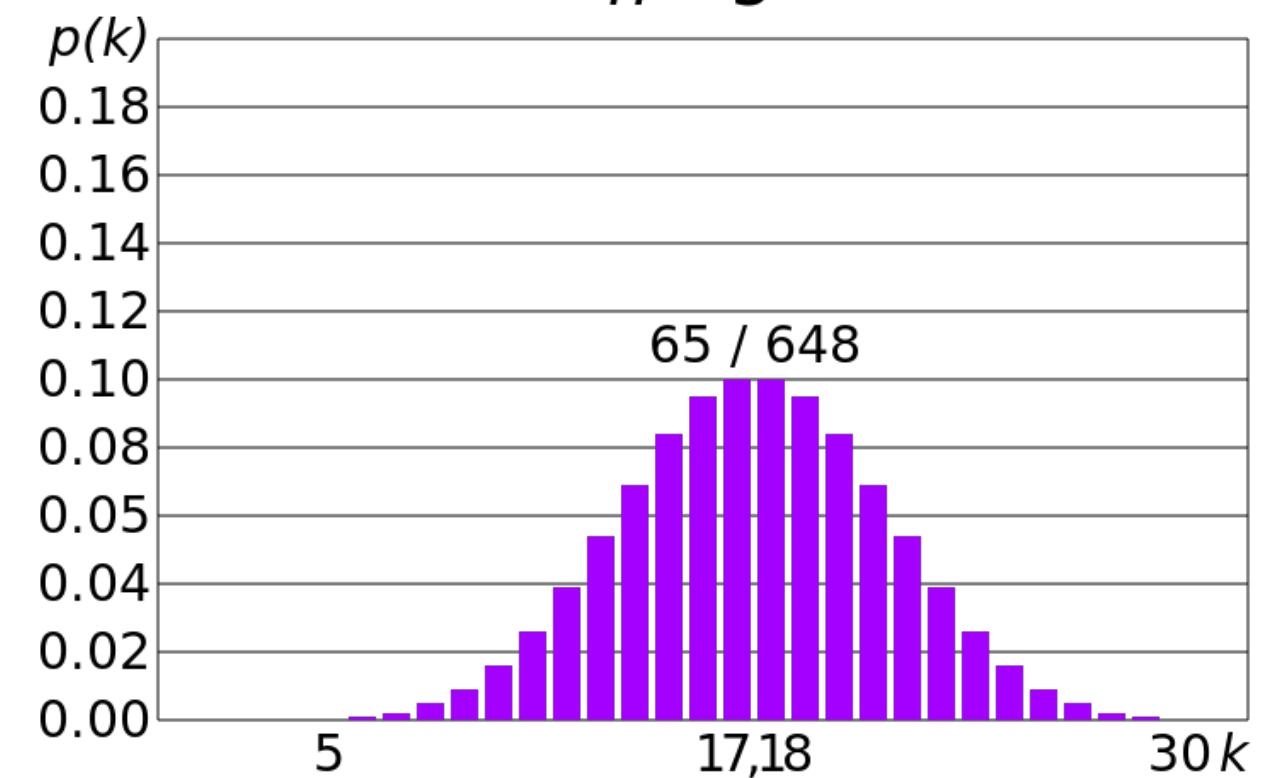
$n = 4$



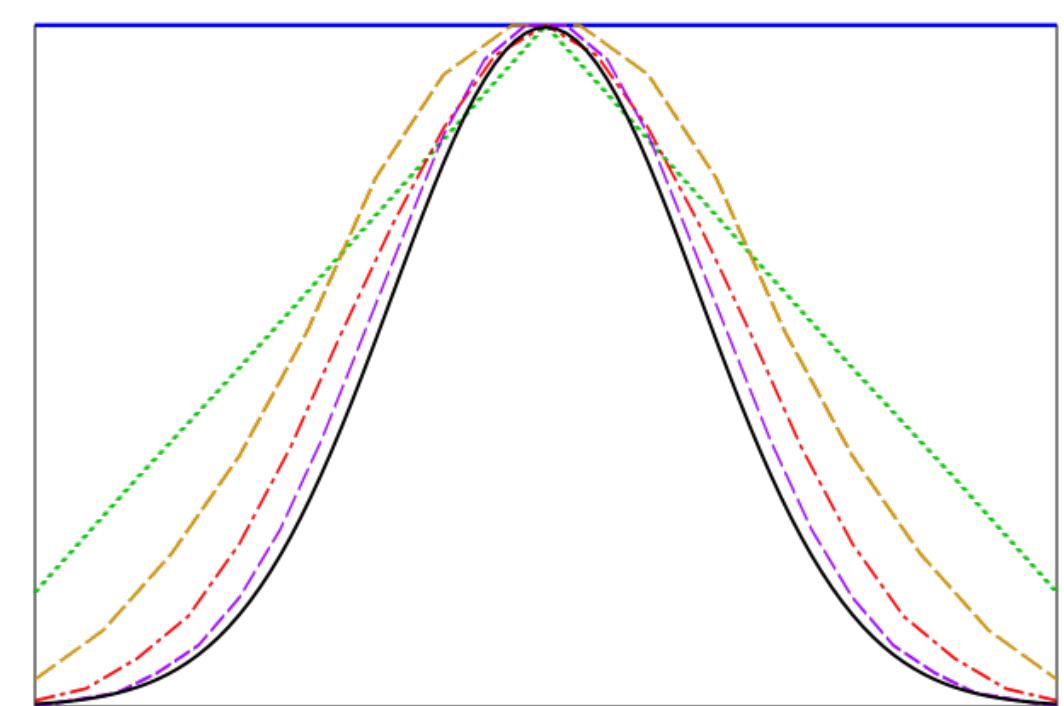
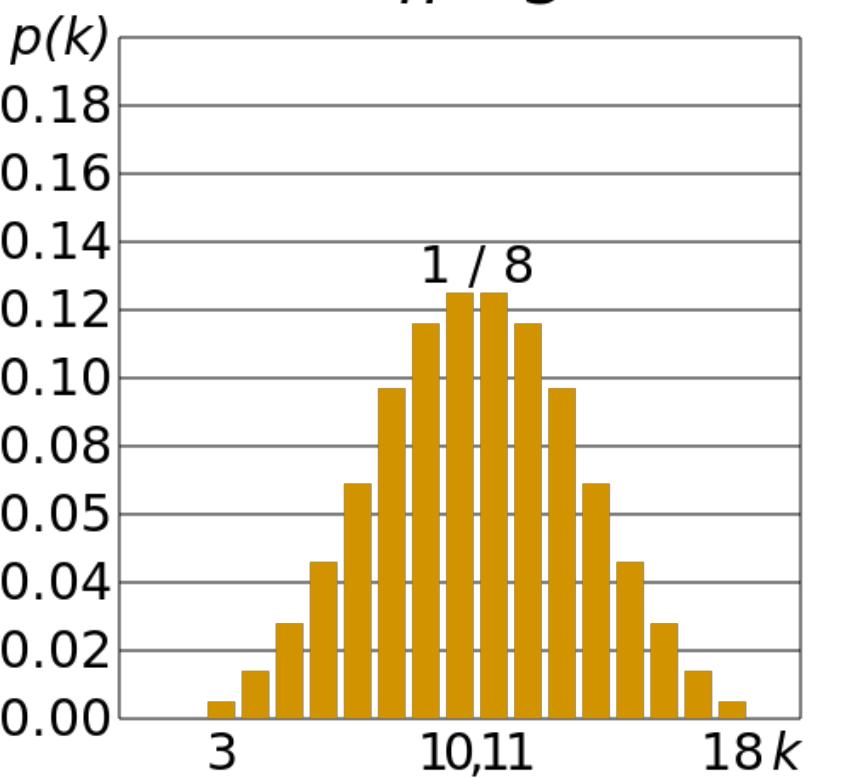
$n = 2$



$n = 5$



$n = 3$



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

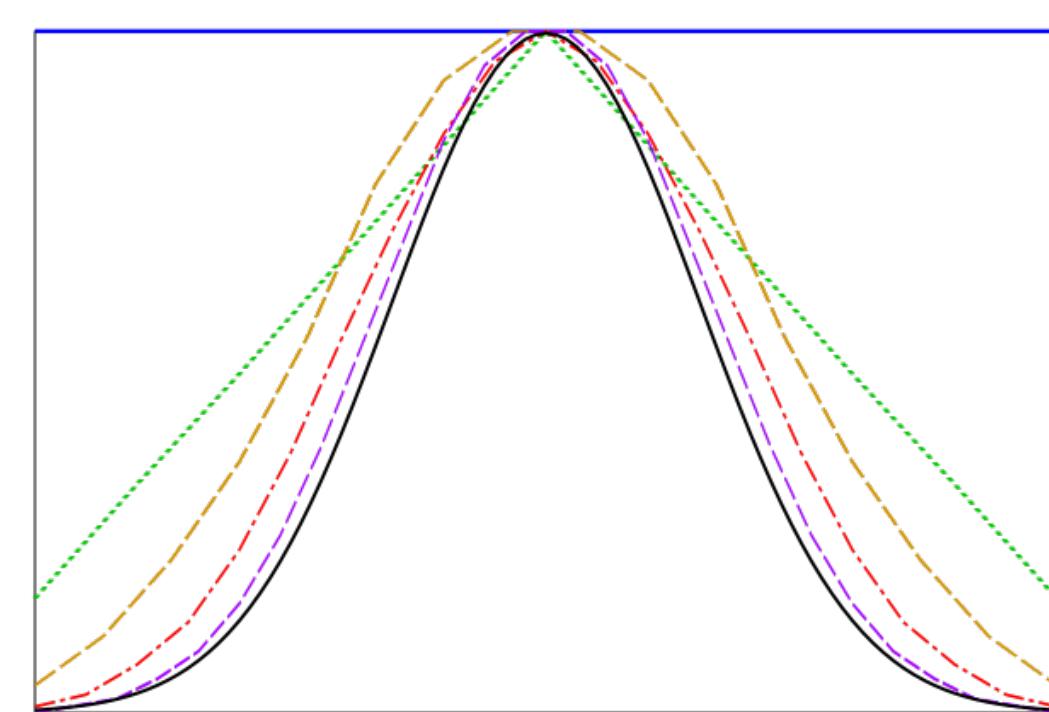
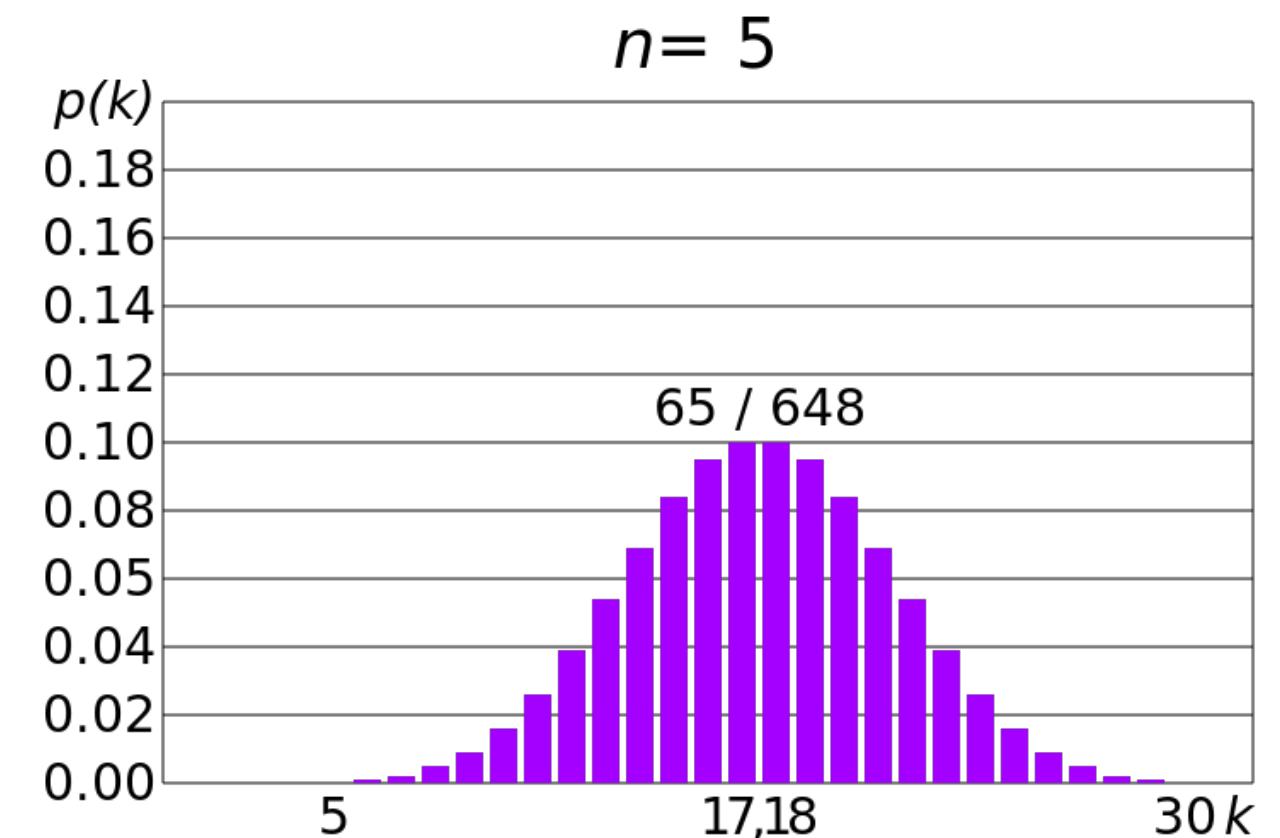
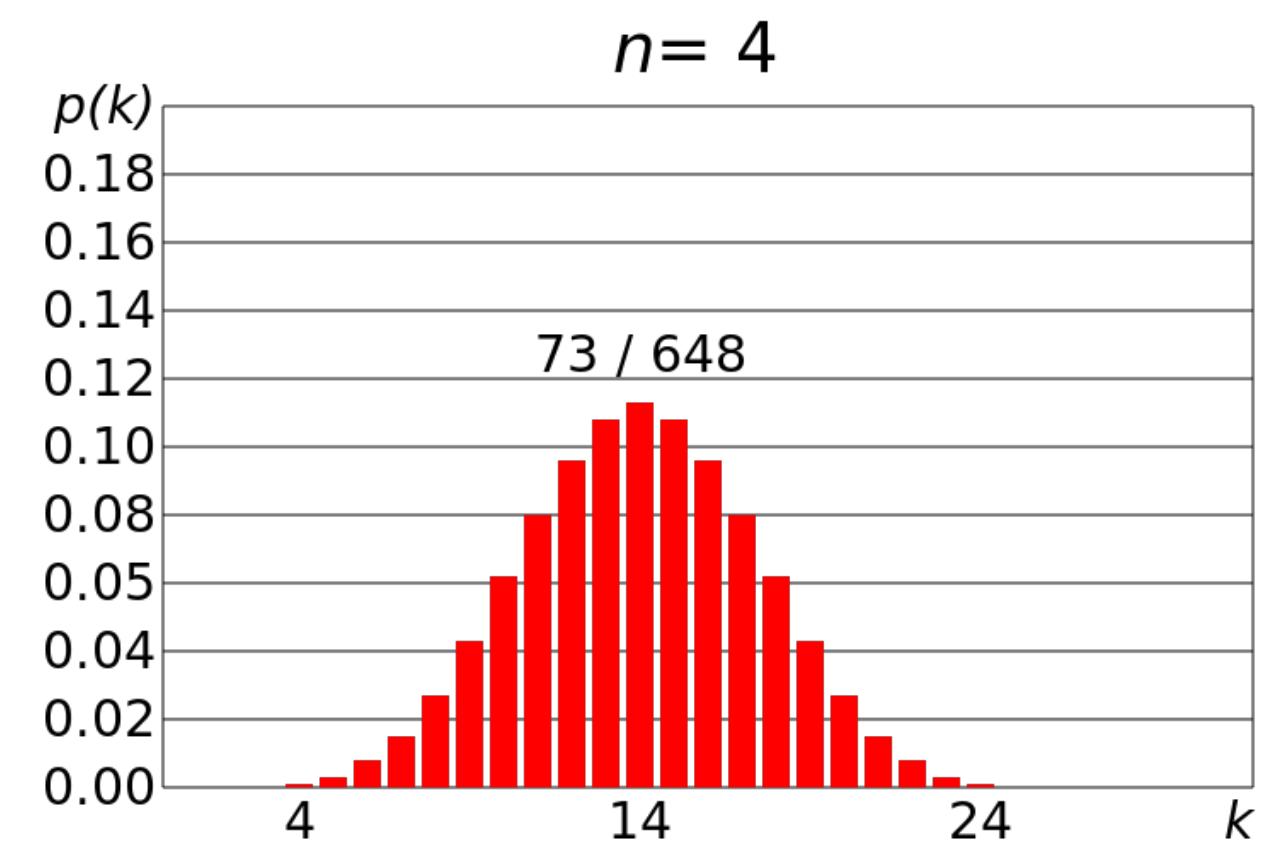
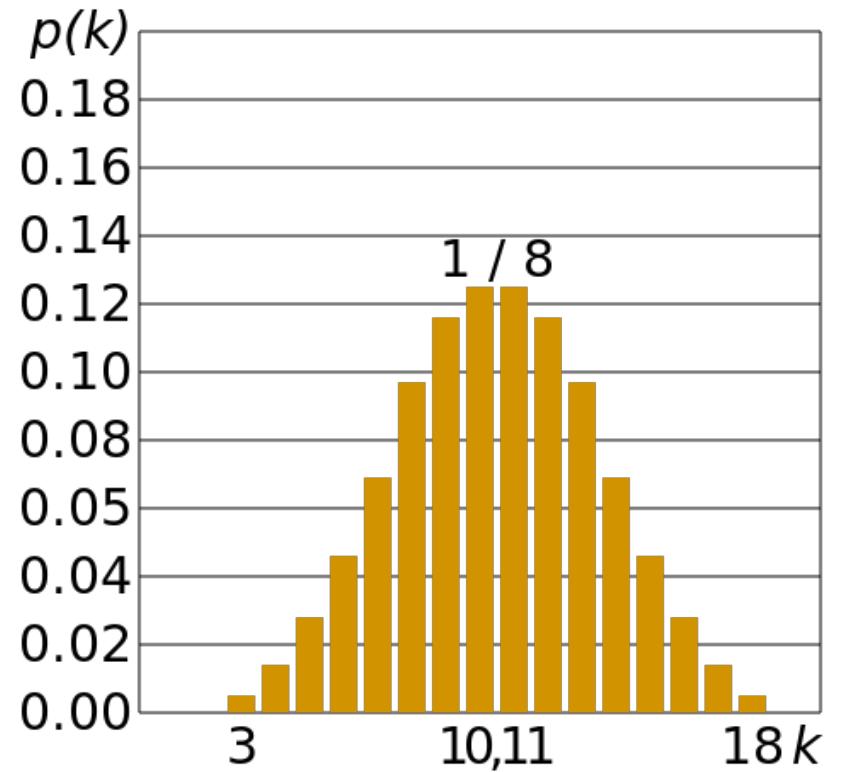
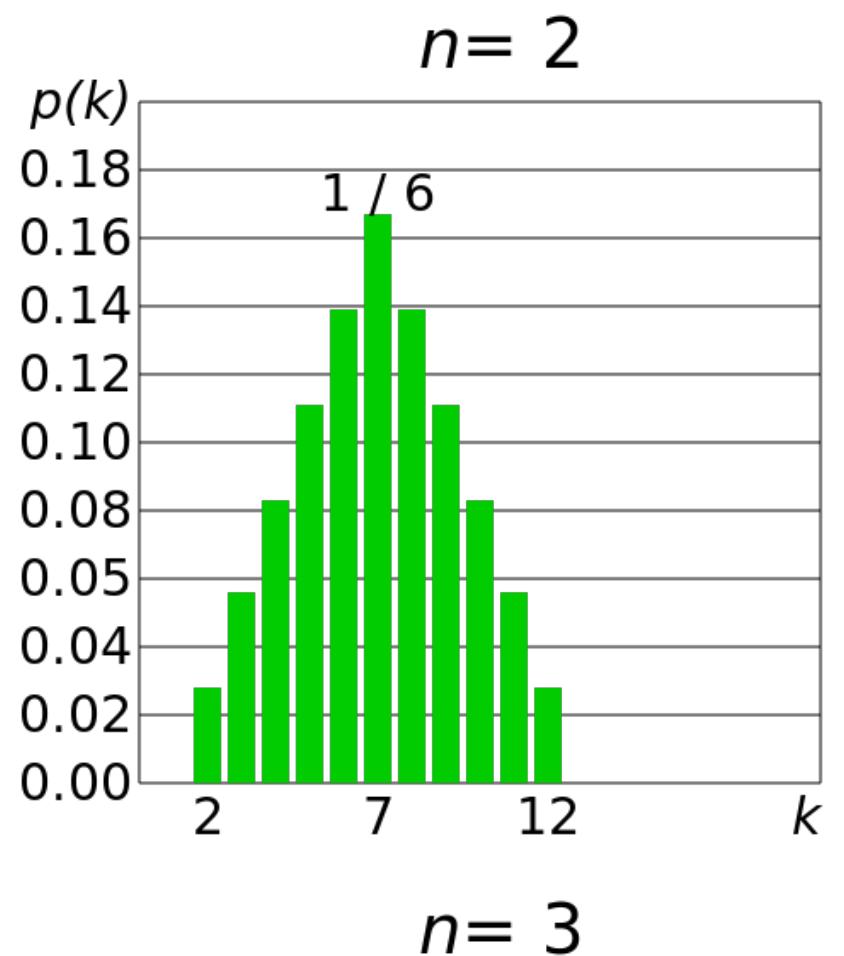
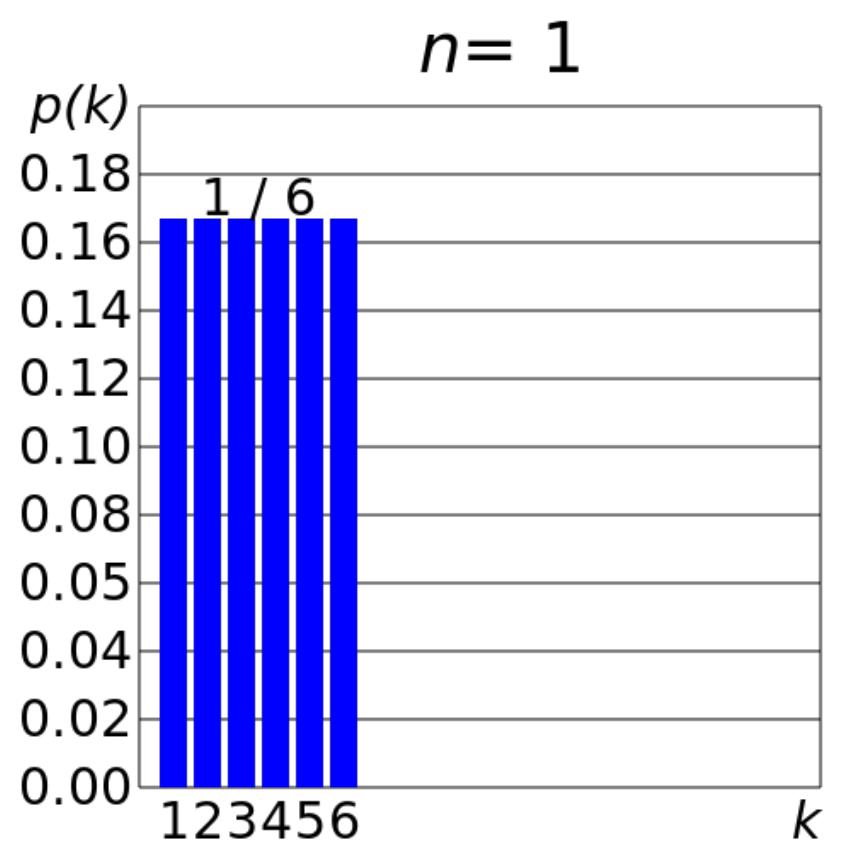


Image credit: Wikipedia

# IS CLT ALWAYS APPLICABLE?

# IS CLT ALWAYS APPLICABLE?

- For some statistics (e.g., median, reciprocal rank, f-score), we're not doing simple averaging

# IS CLT ALWAYS APPLICABLE?

- For some statistics (e.g., median, reciprocal rank, f-score), we're not doing simple averaging
- 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 some difference between two categories
- Null hypothesis will be that there is no difference (i.e labels should not matter)
- If labels don't matter then repeatedly randomizing the label assignments to examples should yield us the null distribution



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

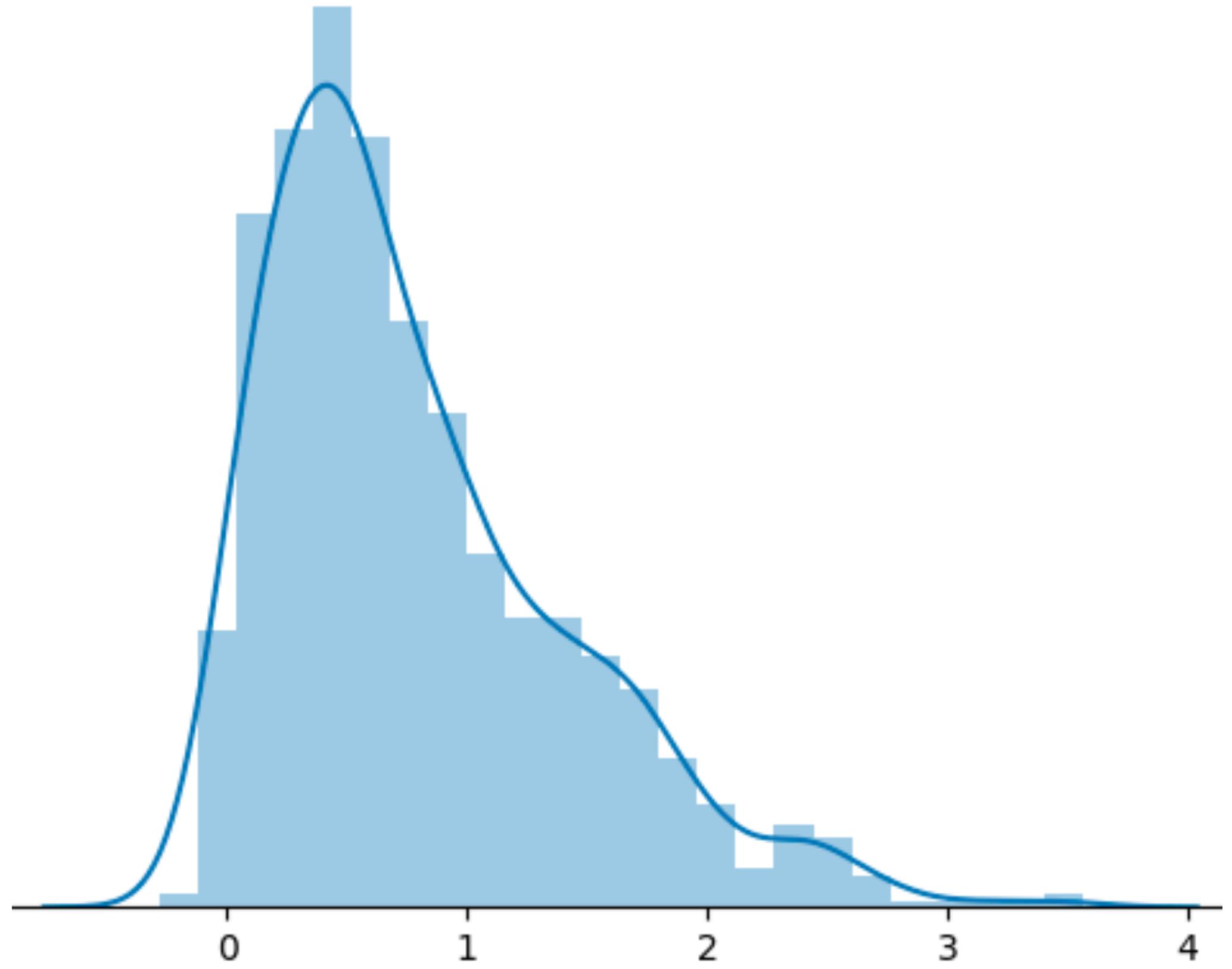
-1.4

2.1

...

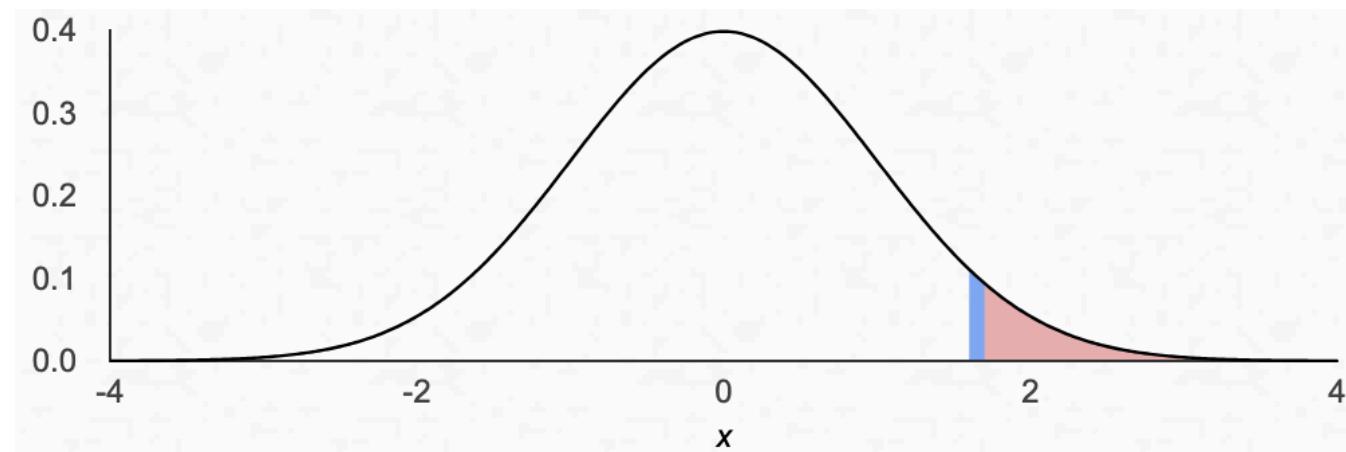
1.6

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

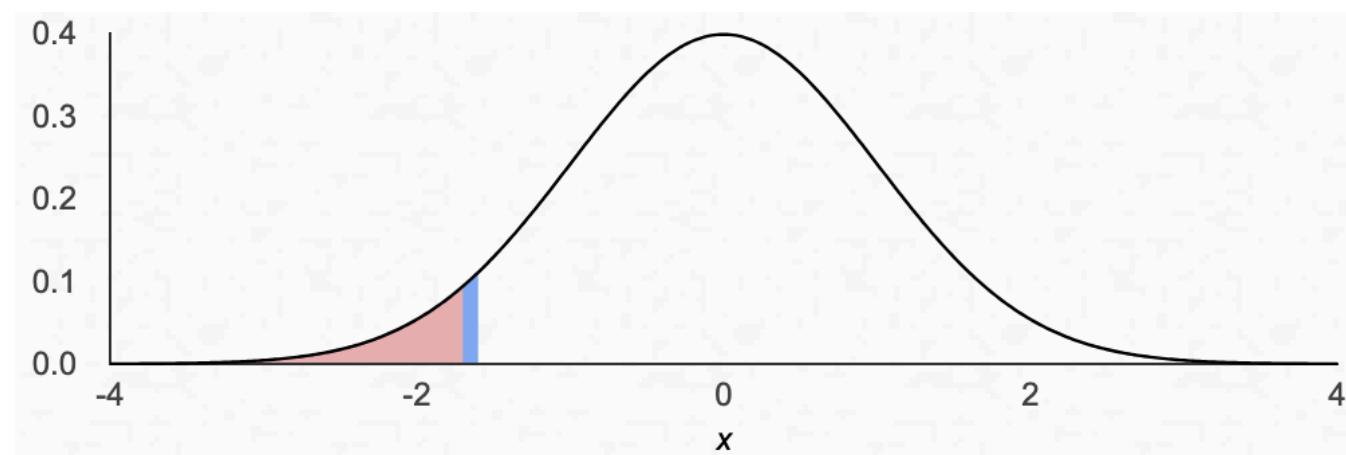


Is 3.9 unusual for this constructed null distribution?

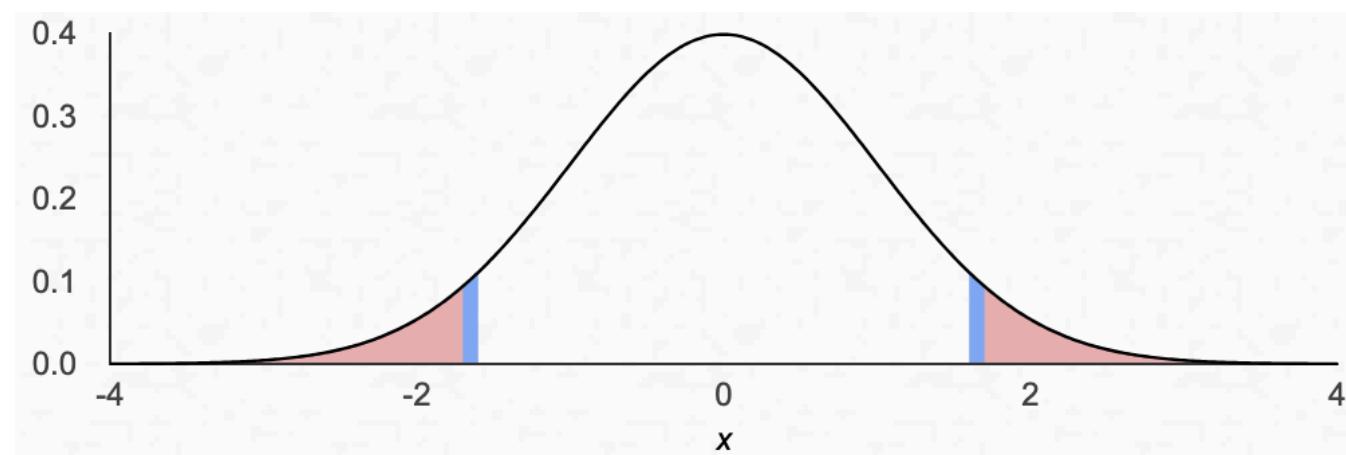
# P VALUE



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



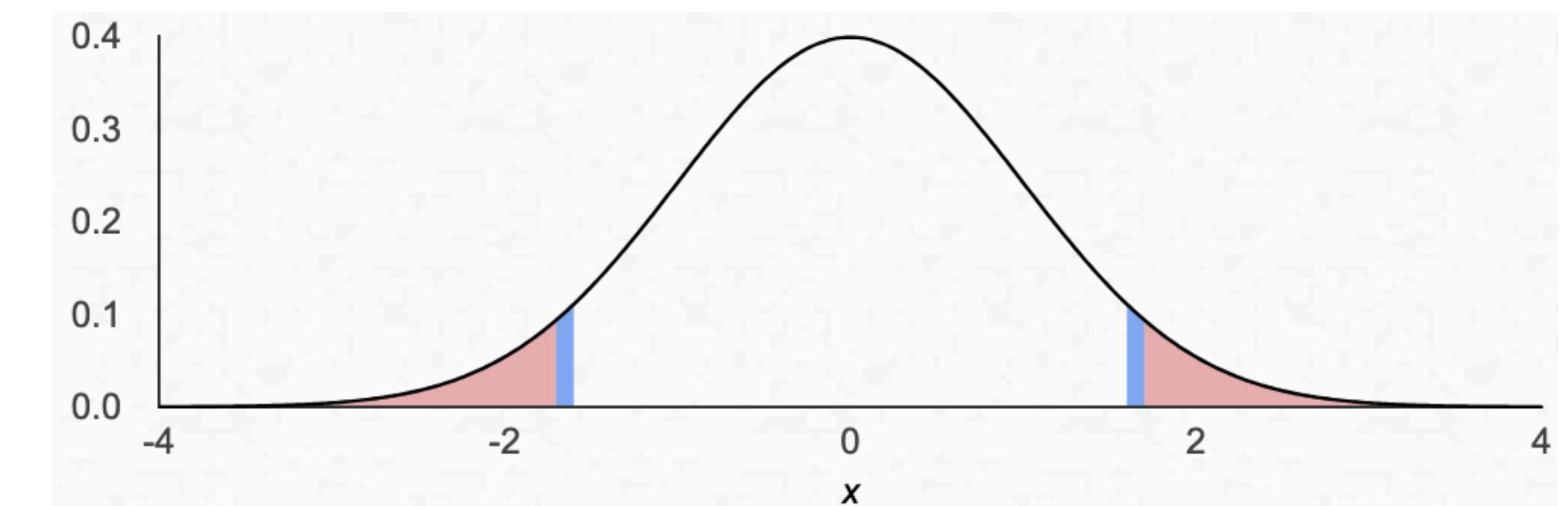
$$\text{p-value}(x) = 1 - P(X \leq x | H_0)$$



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

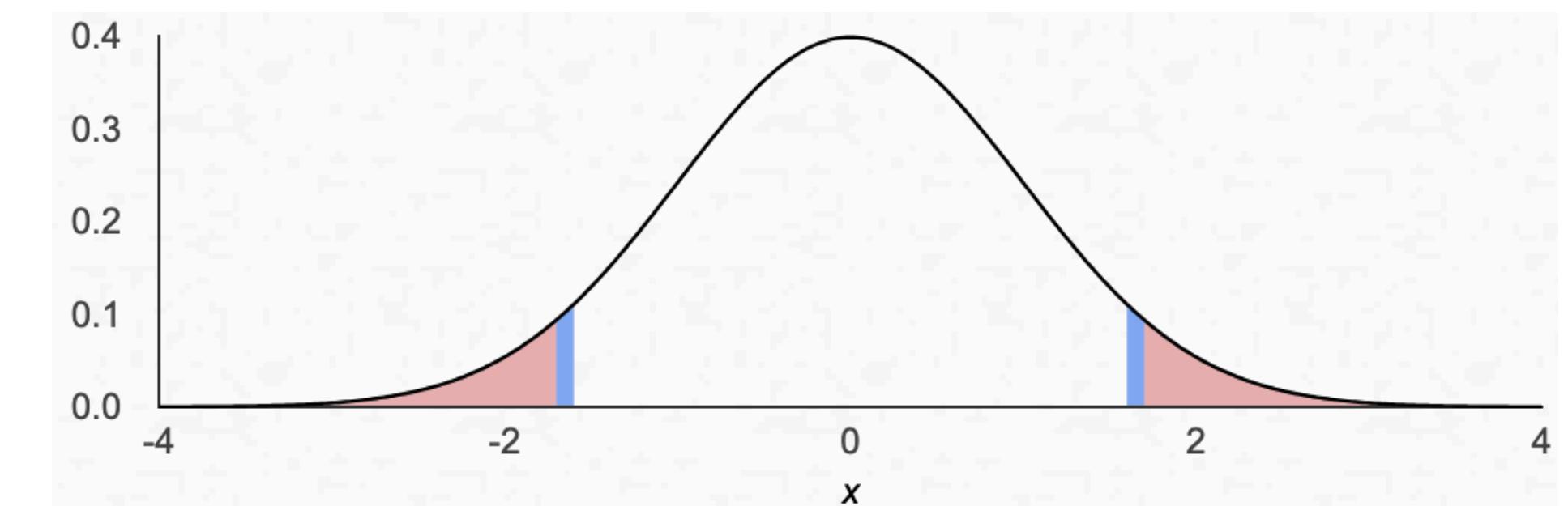
# EMPIRICAL P VALUE

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



# 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

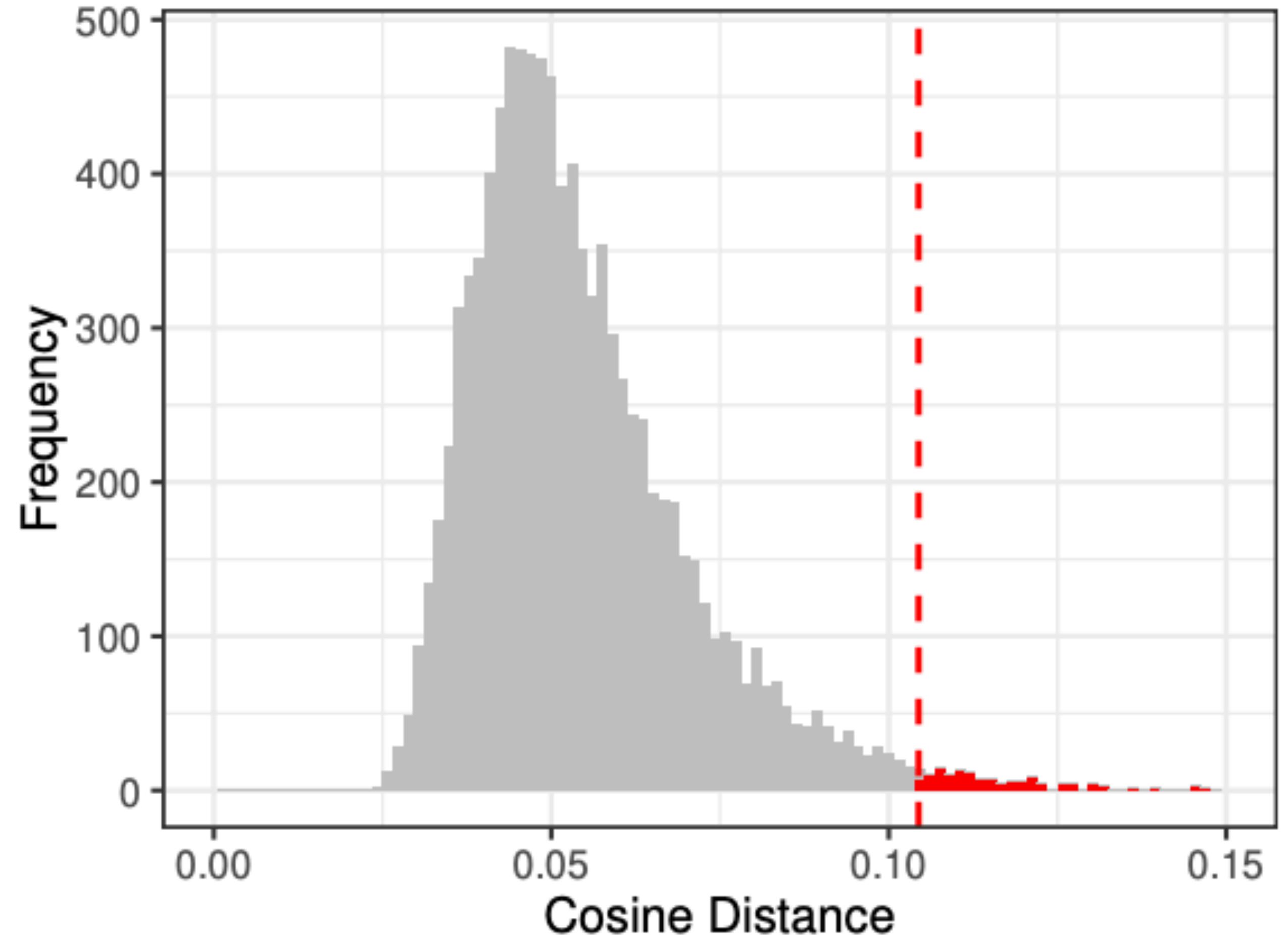
# PERMUTATION TESTS

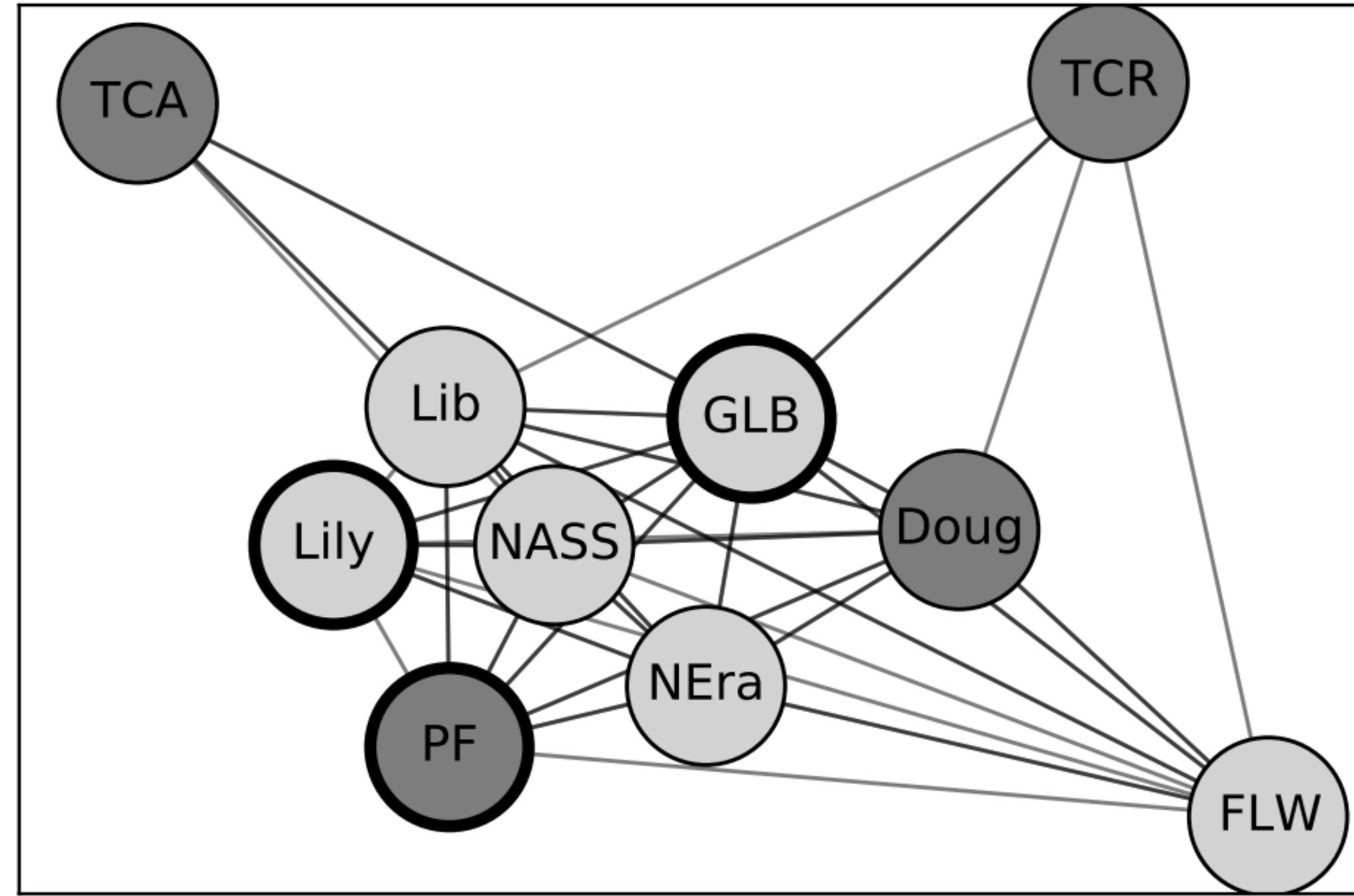
- Permutation tests have broad application
- You don't make any assumptions about the shape and form of the null distribution for the statistic

# 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

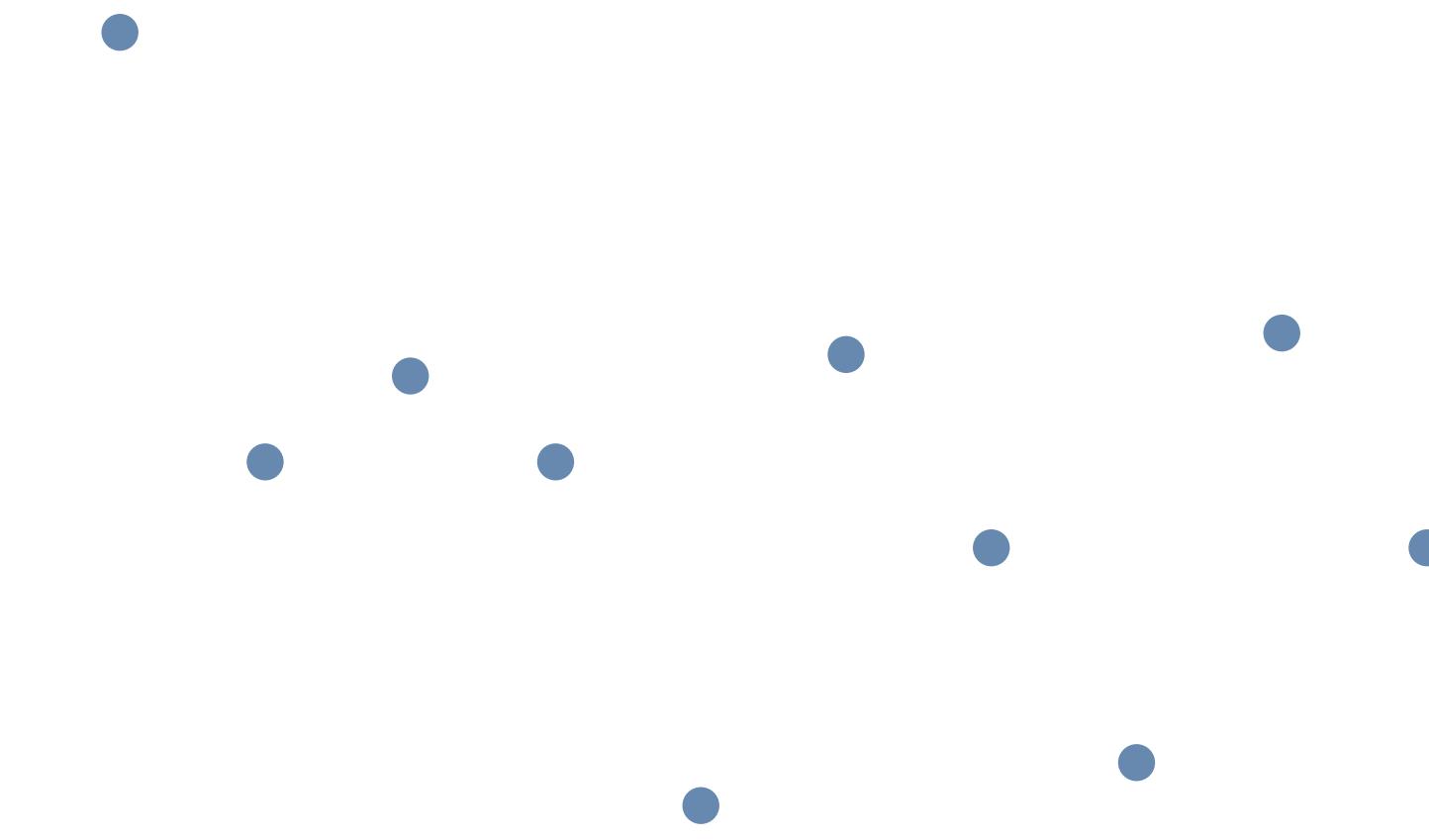
# BOOTSTRAP

- 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



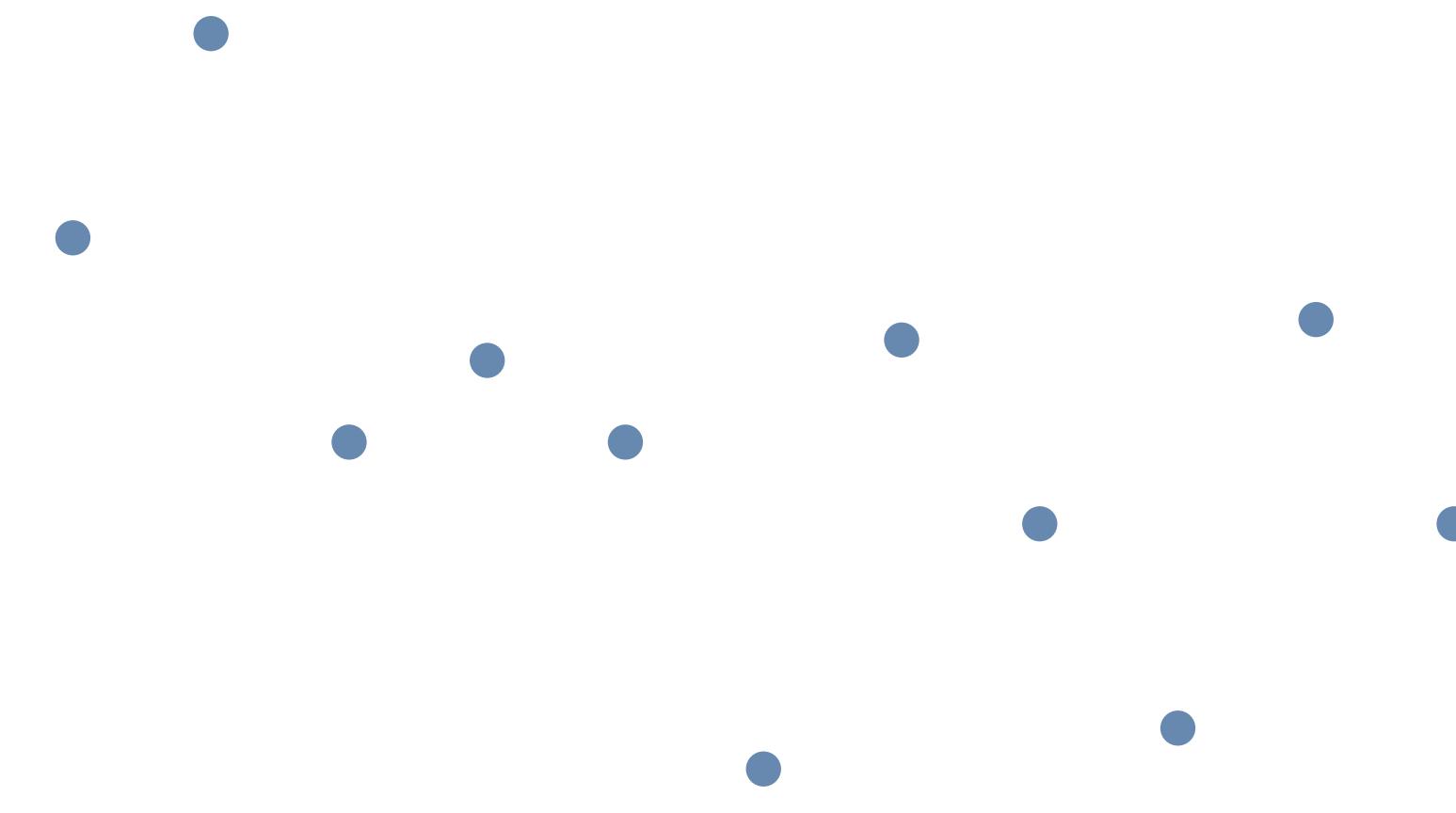
# BOOTSTRAP

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

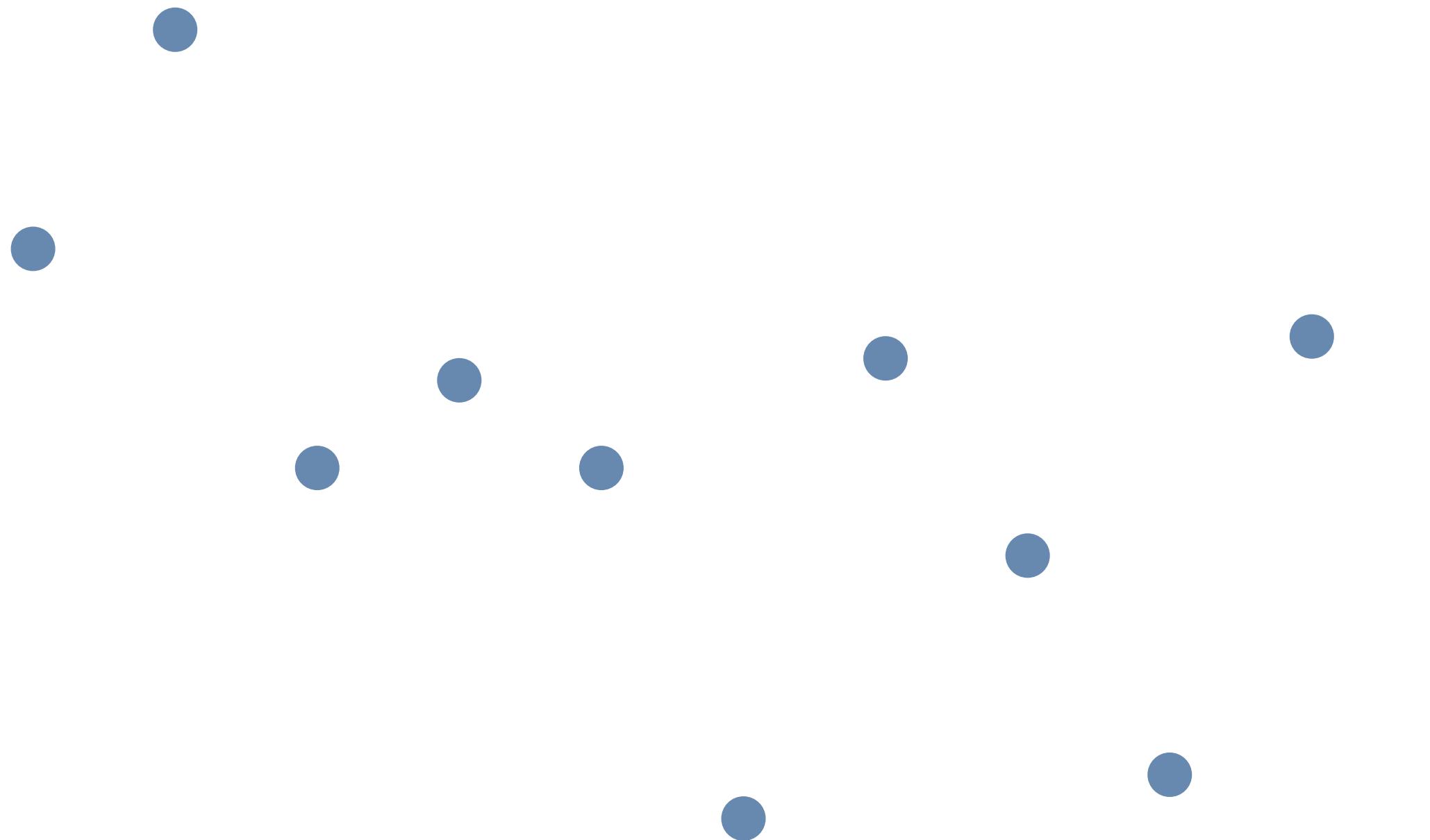


# BOOTSTRAP

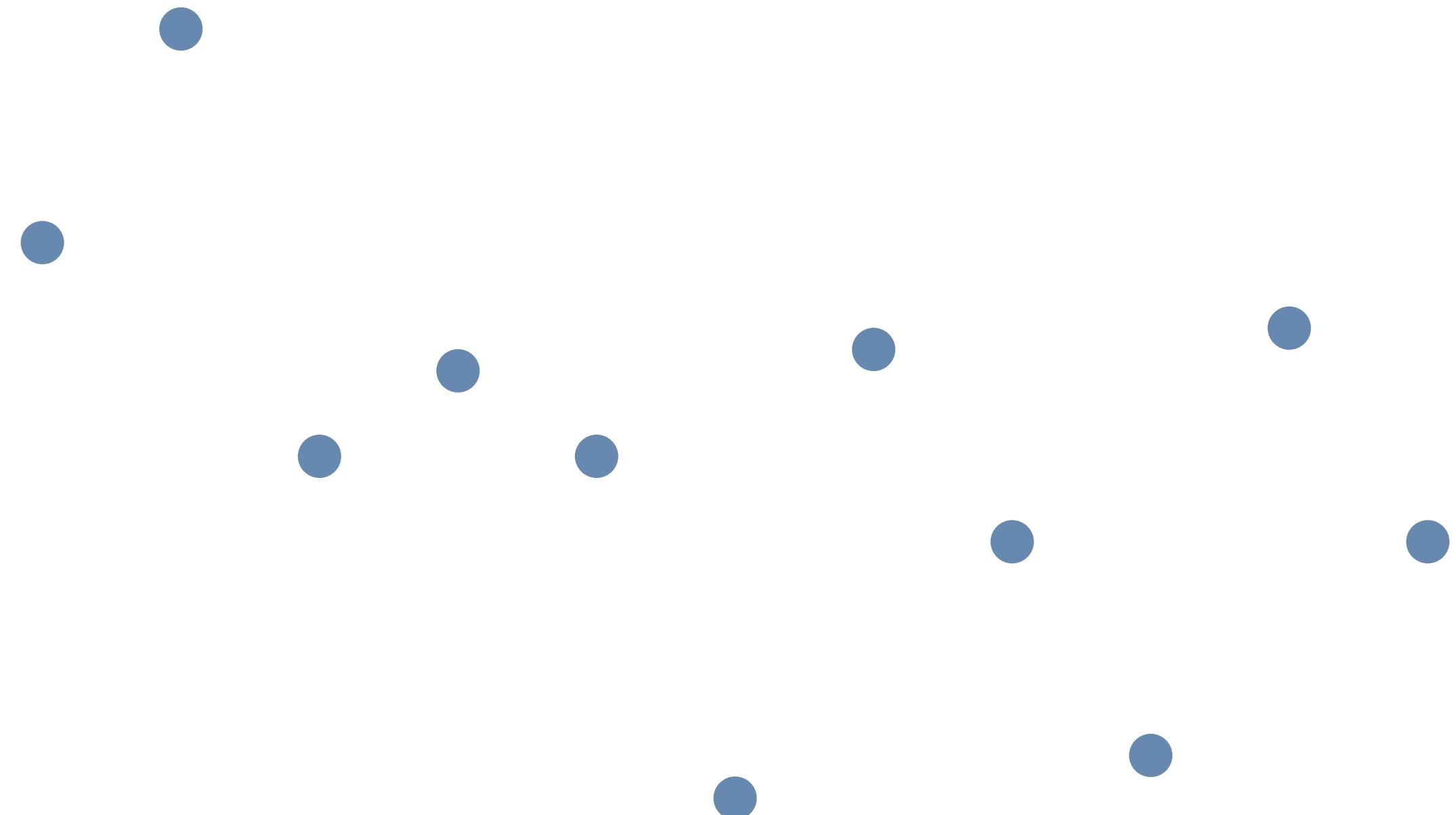
- Our sample is assumed to be iid, and represents the population distribution
- Repeatedly **sample** from our **sample** to generate many alternative datasets



# DISTRIBUTION OF OUR SAMPLE

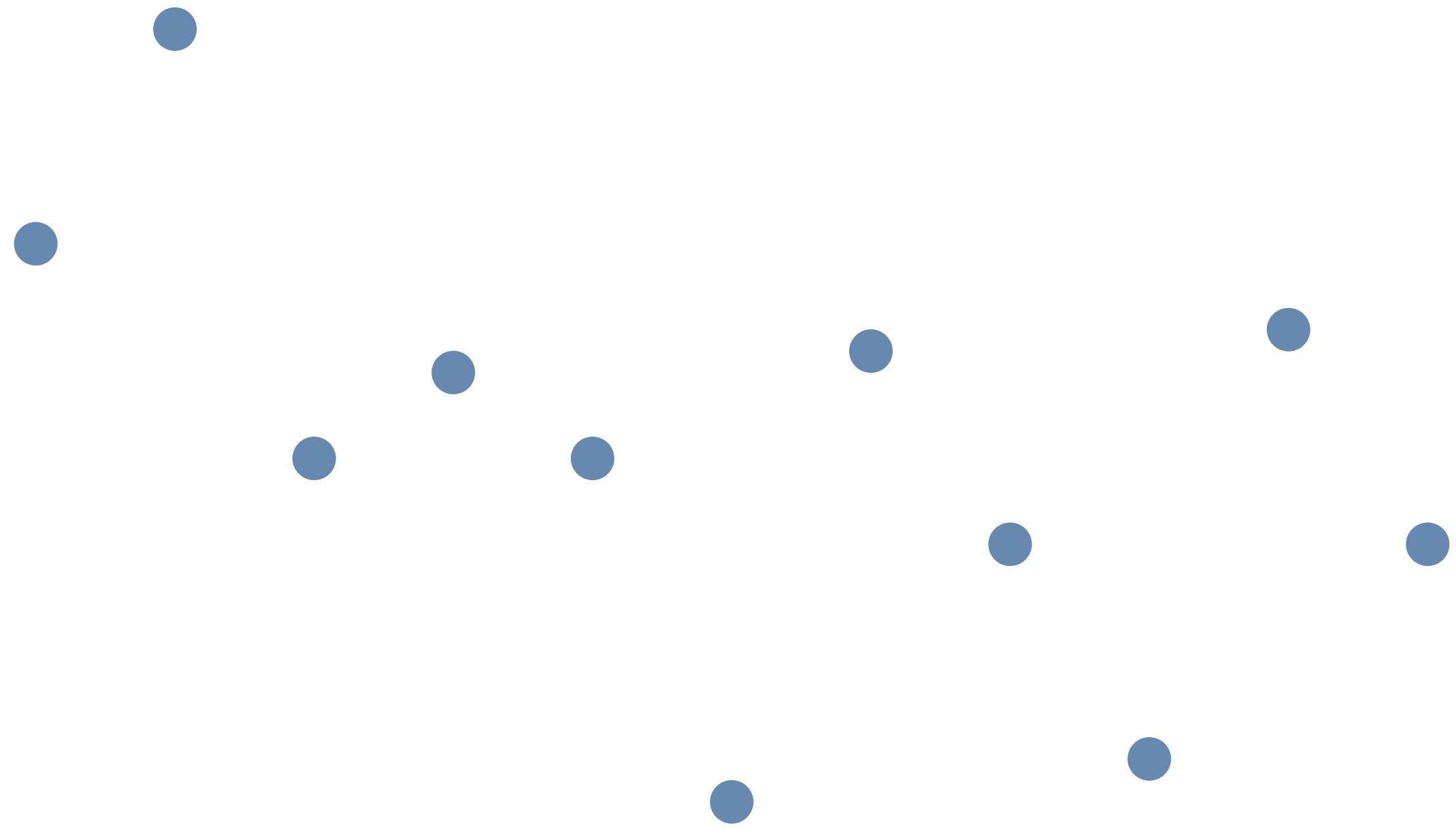


# DISTRIBUTION OF OUR SAMPLE



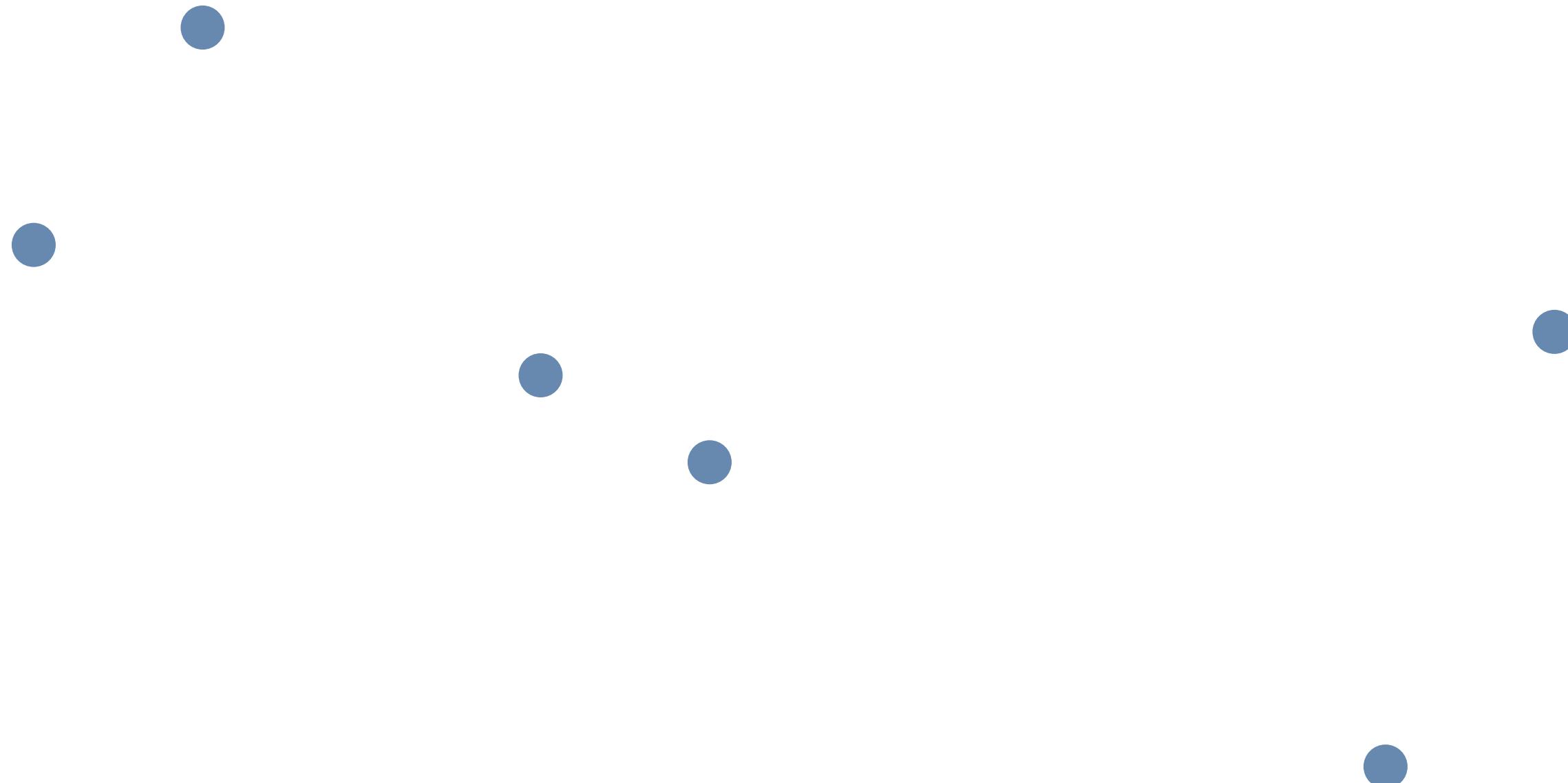
- We can calculate some metric on this sample such as accuracy

# DISTRIBUTION OF OUR SAMPLE

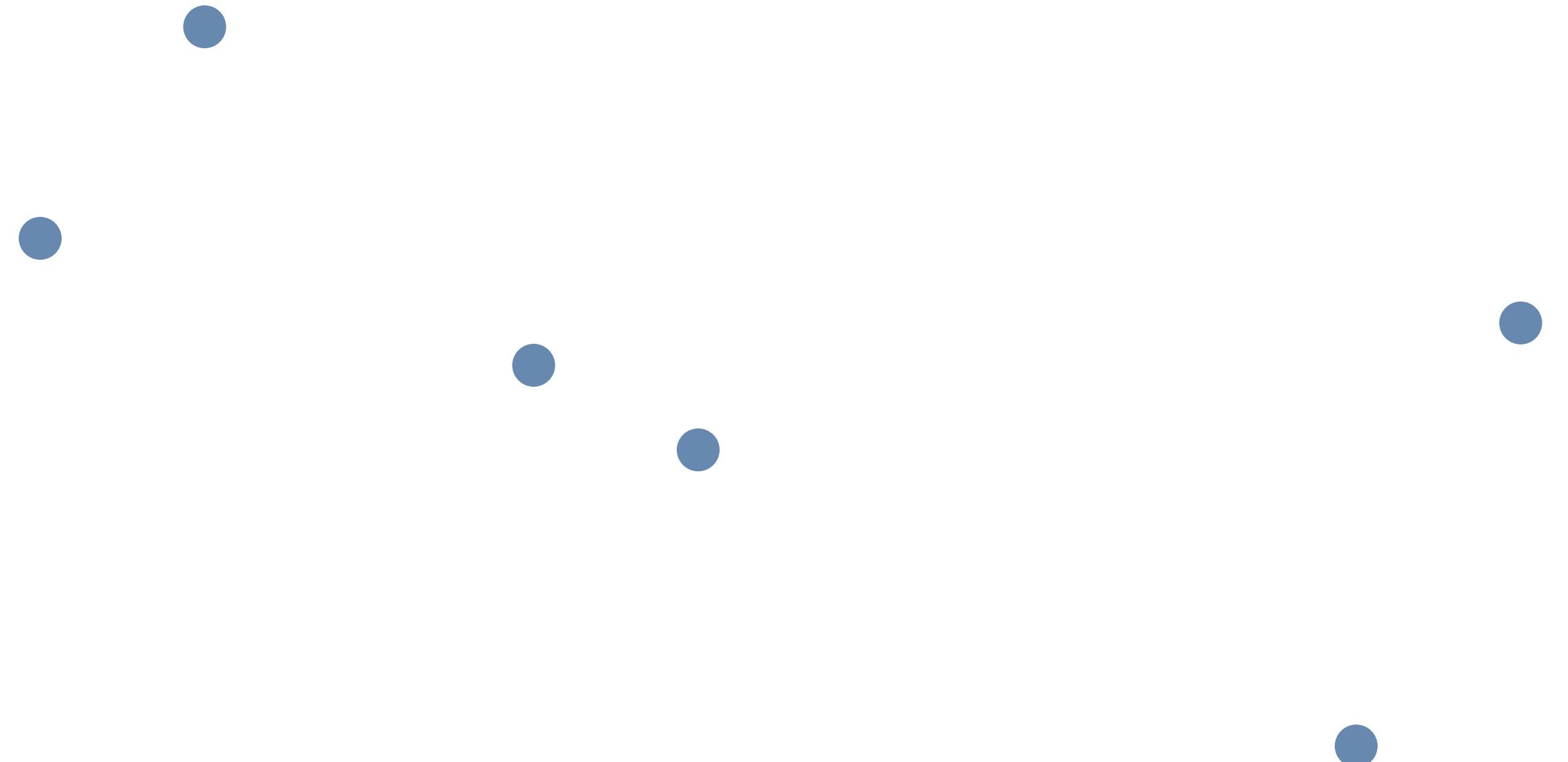


- We can calculate some metric on this sample such as accuracy
- $\text{acc}_{\text{obs}}$

# BOOTSTRAP 1

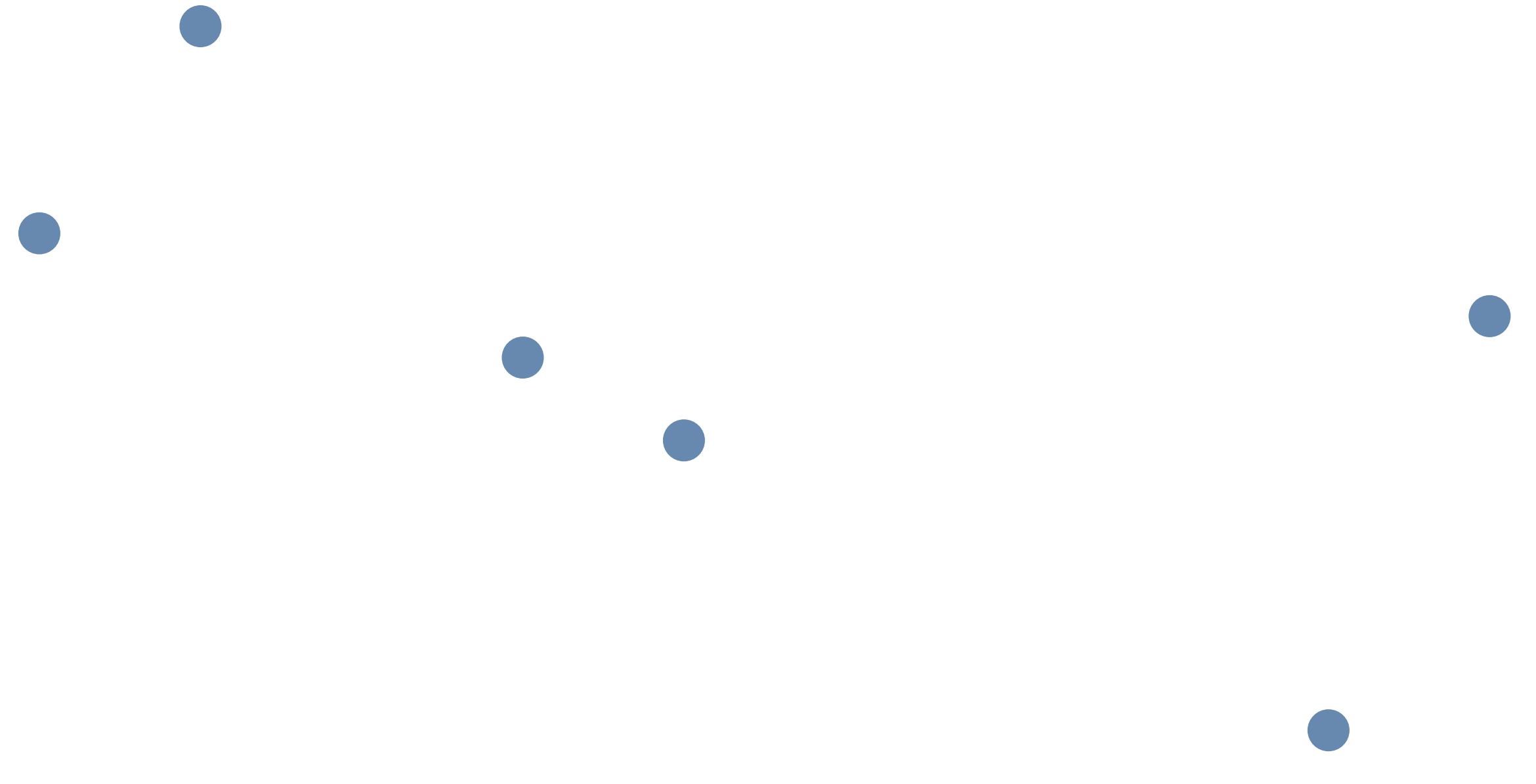


# BOOTSTRAP 1



- Calculate accuracy on this bootstrapped sample

# BOOTSTRAP 1

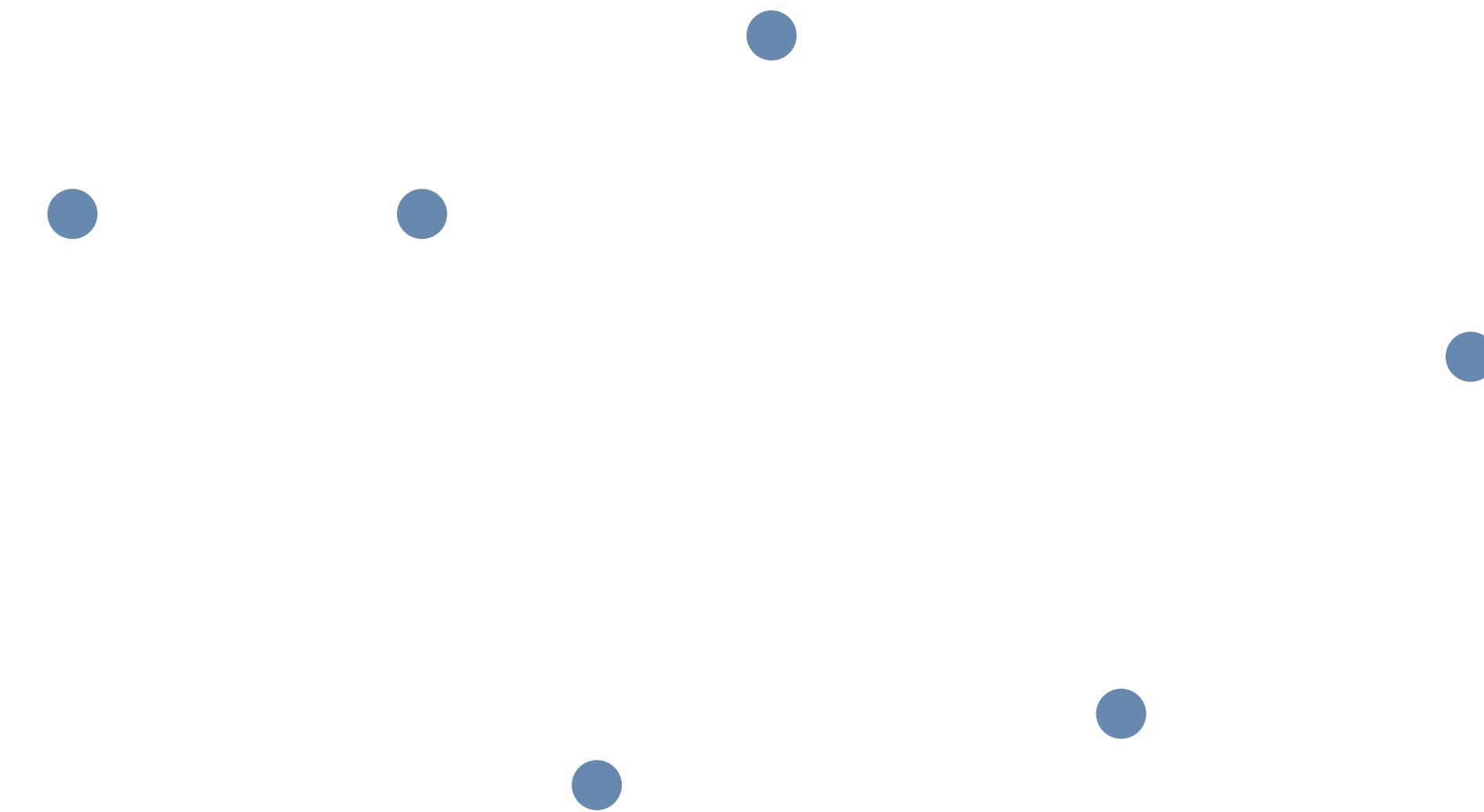


- Calculate accuracy on this bootstrapped sample
- $\text{acc}_{\text{boot}}$

# 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

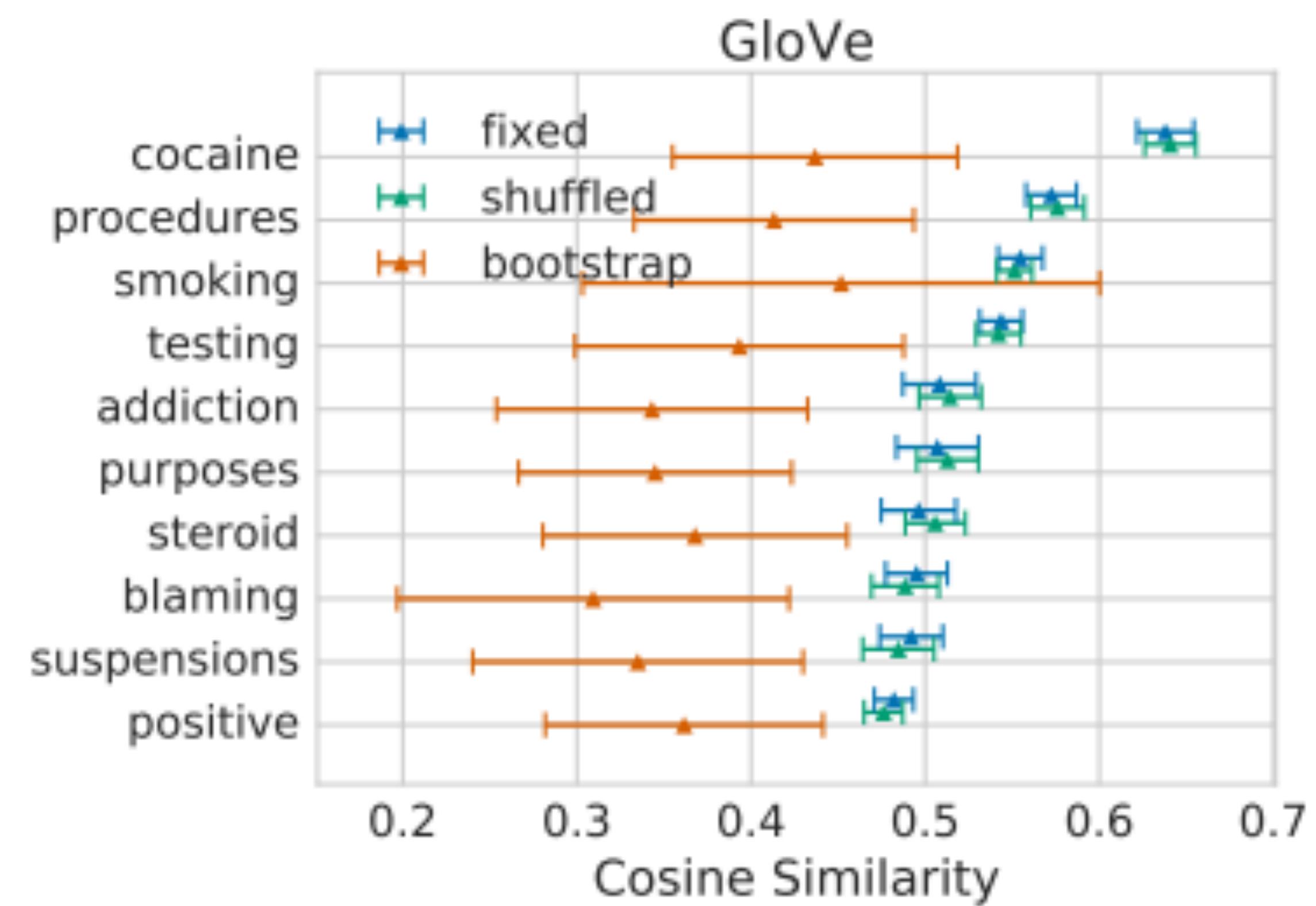
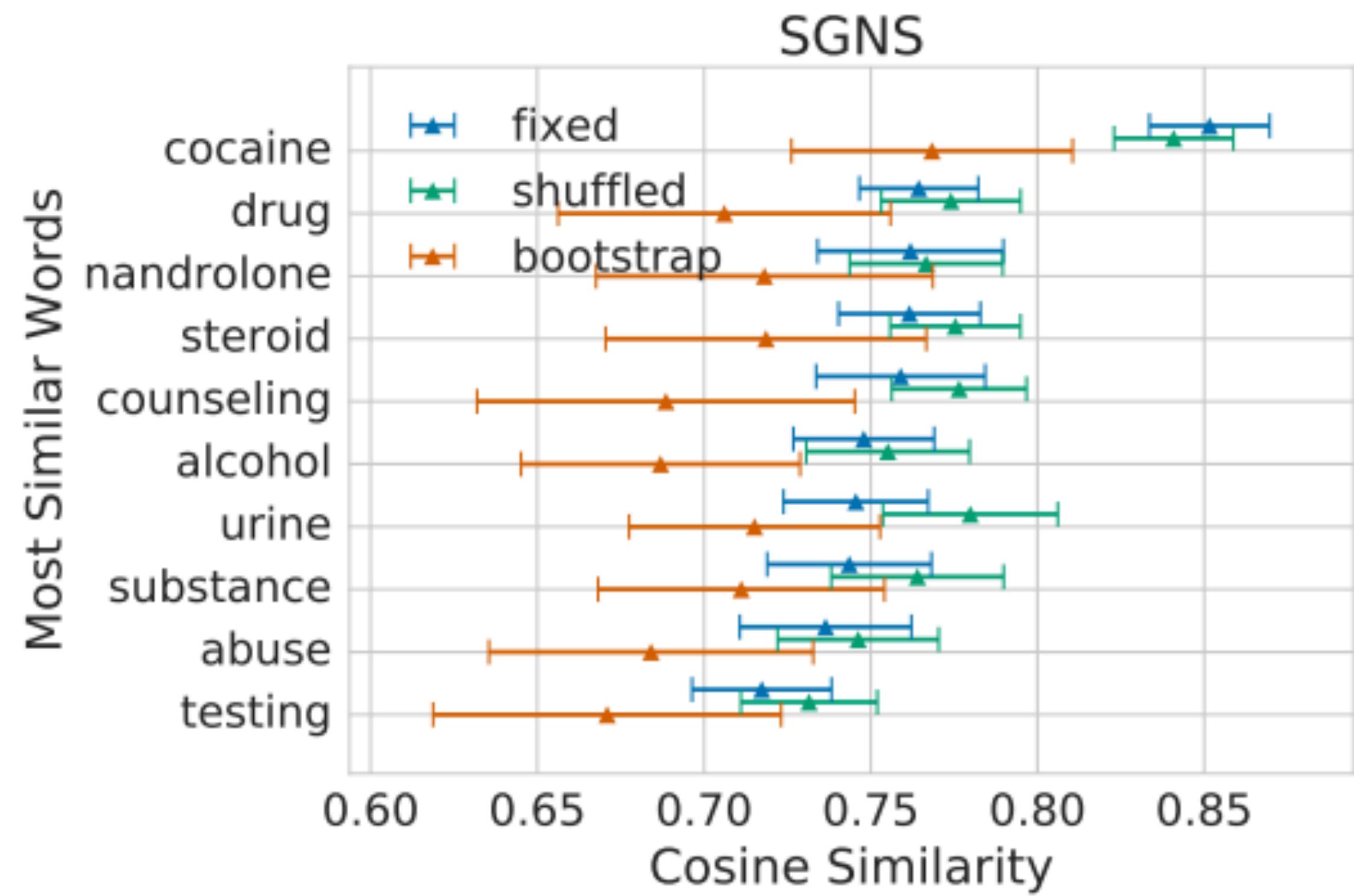
- From  $b$  bootstrap samples, we can get a size  $b$  array as estimates for accuracy
- We can rank these estimates and find the middle 95% to define a range in which the accuracy estimate falls

# 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?

# IN CLASS

- Non parametric test