


Managing {Probability, Projects}

September 18, 2018
Statistics Bootcamp



Agenda

- Set notation and probability
- English  Probability notation
- The Sampling: Part II
- Hypotheses and Inferences
- Lunch
- Organizing your Stata/mind

Warm up with hot takes

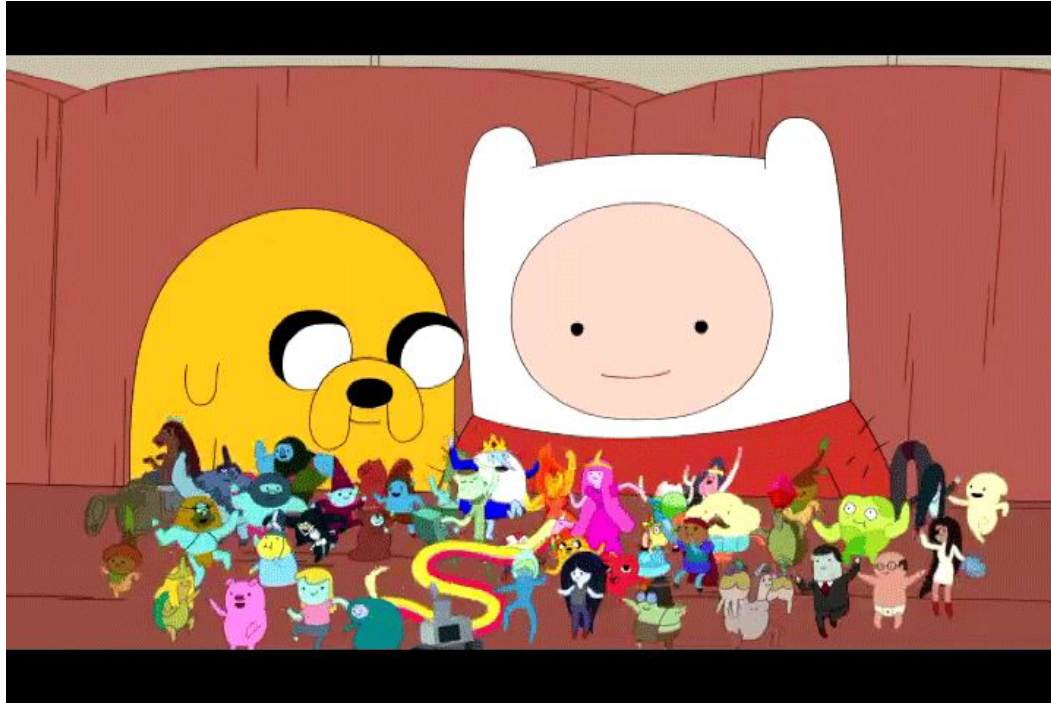
- Write a hot take
 - You don't have to believe it
 - It can be cold
- Describe evidence that would signal your take is correct
- Describe noise that would make people less likely to accept your take

Picture telephone

Goals for Probability and Distributions

- Practical:
 - ...read and use probability notation to describe a situation with bounded uncertainty.
- Conceptual:
 - ...explain some ways that probabilistic model assumptions can undermine our models of society.

Statistical models, probability, and assumptions



“It’s like a finger pointing a way to the moon....

Don’t concentrate on the finger,

Or you will miss all the heavenly glory.”

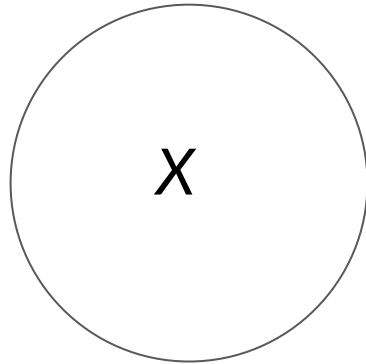
Bruce Lee, Enter the Dragon, 1973

Probability and set notation

The set of all events: $\{\}$

Any given day can be rainy,
sunny, or surreal.

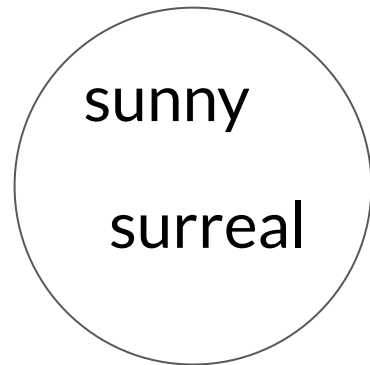
Possible weather events =
 $\{\text{rainy, sunny, surreal}\} = X$ or
some other capital letter



The set of all X such that: {X: some condition}

Imagine the set of all
weather events such that
you could get a tan.

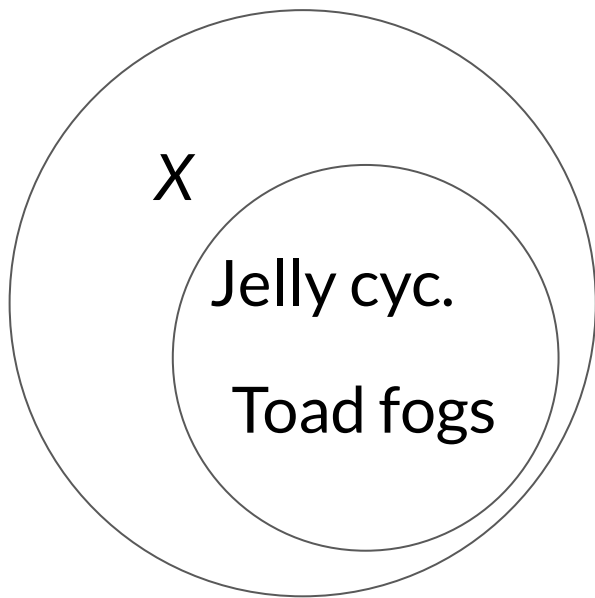
{Weather: you could get a
tan} = {sunny, surreal}



Is an element of: \in

Jelly cyclones and toad fogs
are elements of surreal
weather.

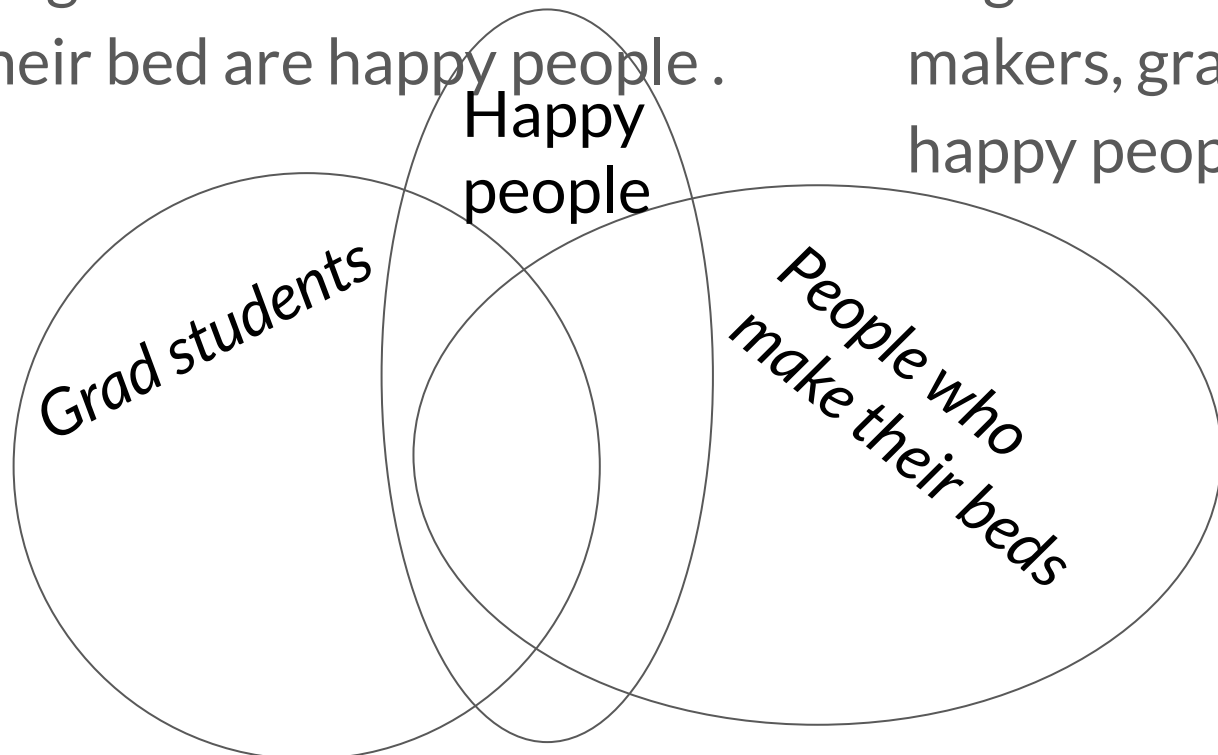
$\{\text{jelly cyclones, toad fogs}\} \in$
Surreal weather \in Weather



For all: \forall

All grad students who make
their bed are happy people.

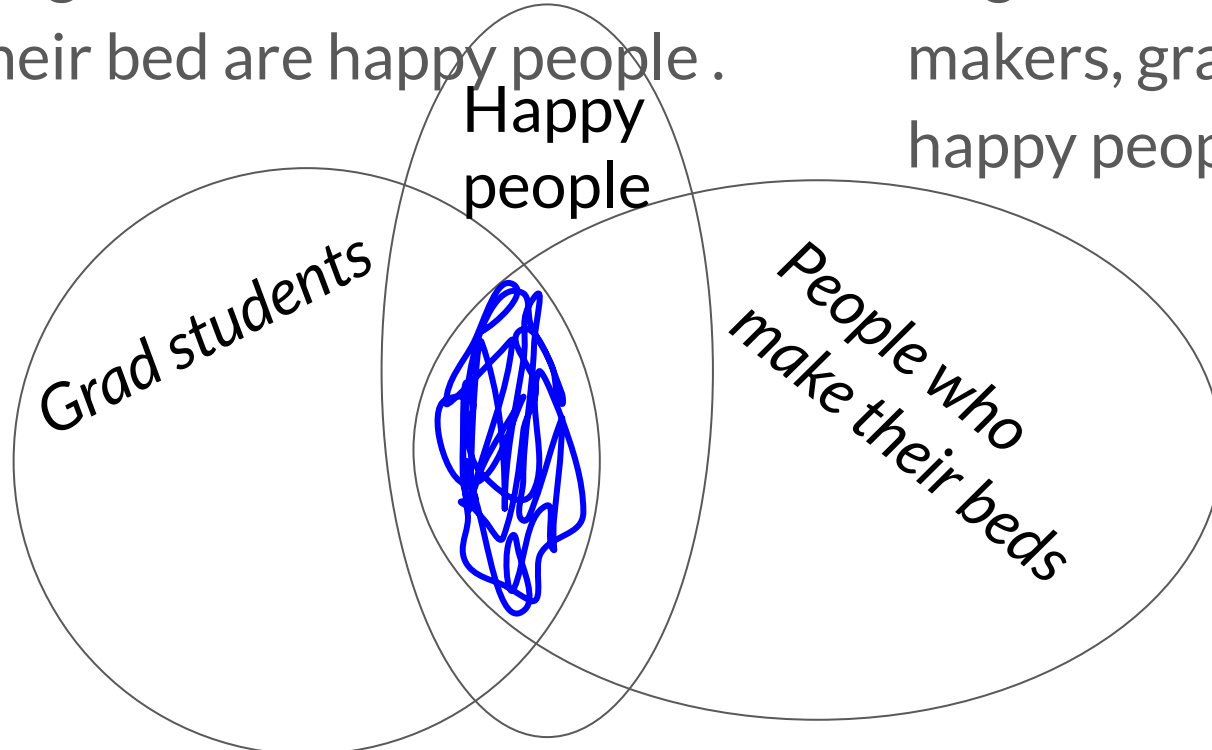
\forall grad students \in bed
makers, grad students \in
happy people



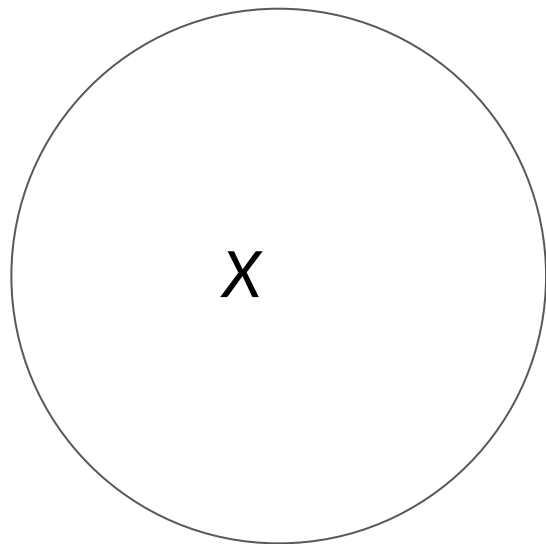
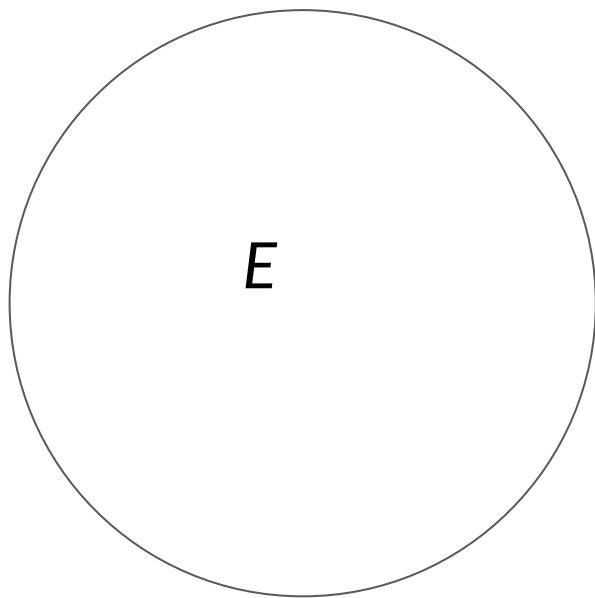
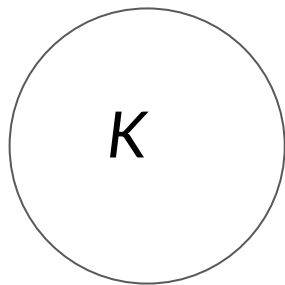
For all: \forall

All grad students who make
their bed are happy people.

\forall grad students \in bed
makers, grad students \in
happy people



Exercise



Write a set notation phrase to the neighbor on your left
Translate!

Flip index card, write a plain-english phrase to right
neighbor
Translate!

The probability that an event occurs: $\Pr(x)$

There is a 70% chance of
rain.

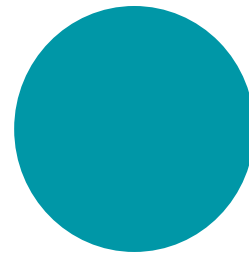
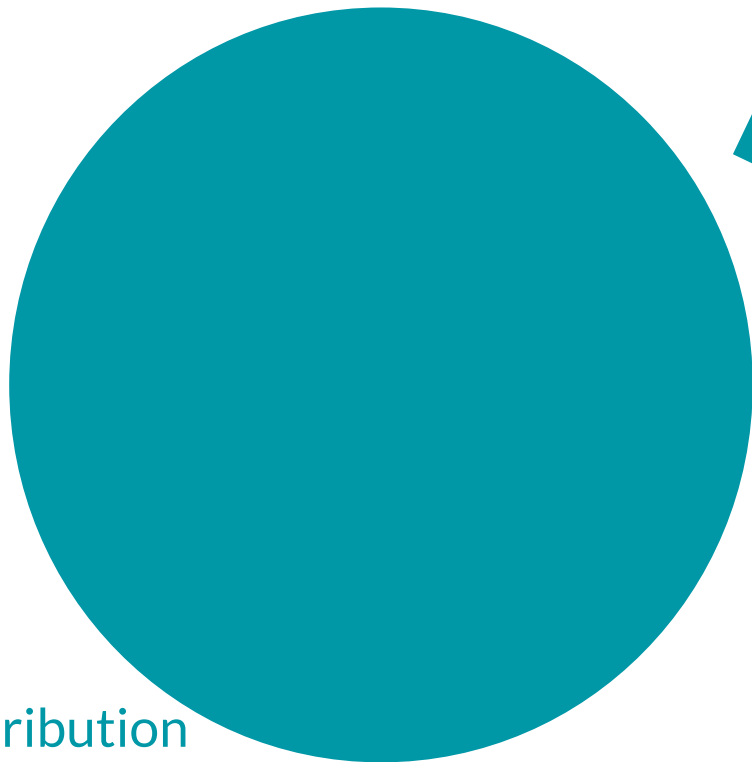
$$\Pr(\text{rain}) = 0.7$$

$$\Pr(\text{no rain}) = 0.3$$

$$\Pr(\text{weather no matter how surreal}) = 1$$

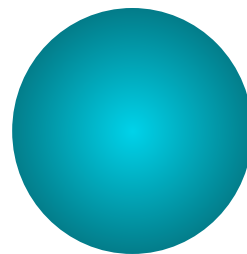
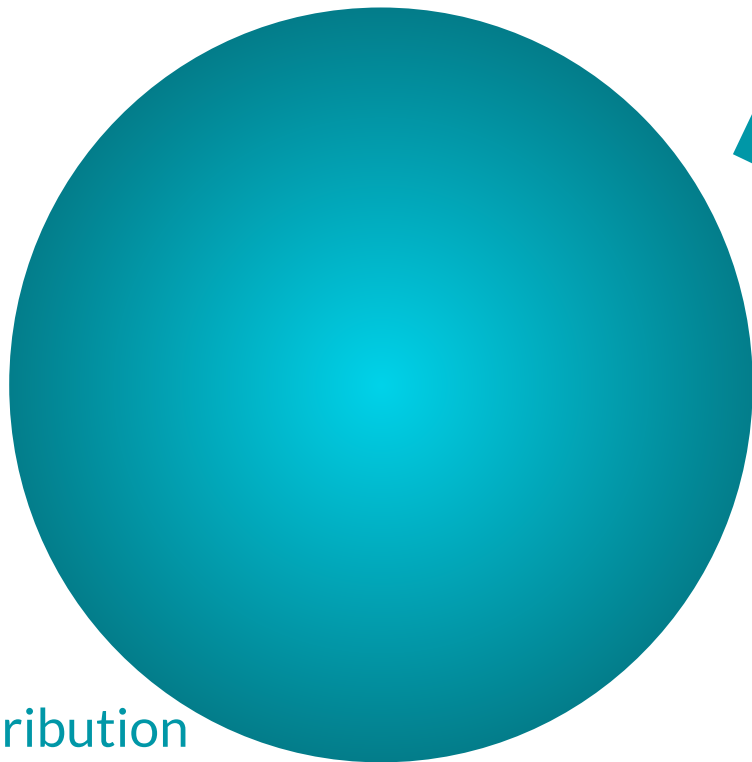
REVIEW: Samples, distributions, and the Sampling

Population distribution

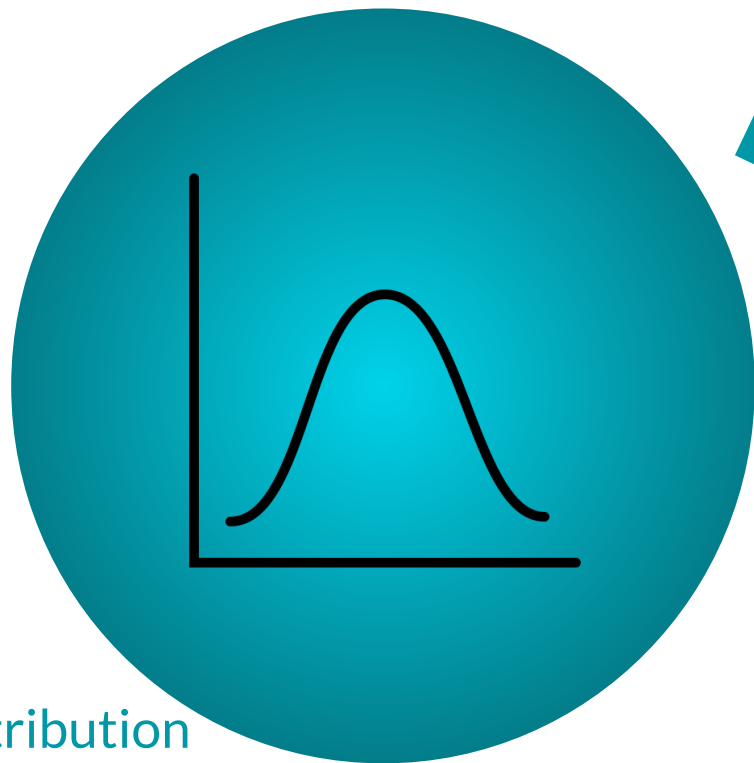


Sample distribution

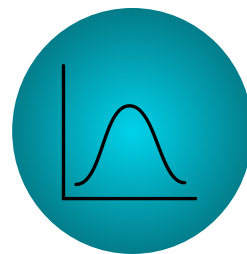
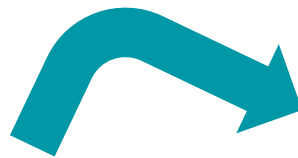
Population distribution



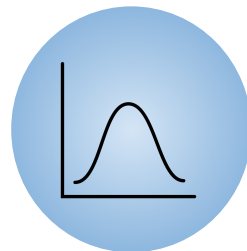
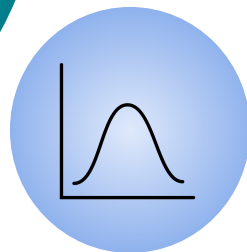
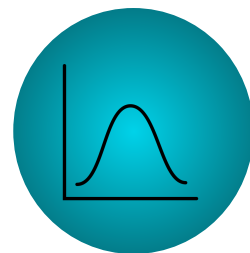
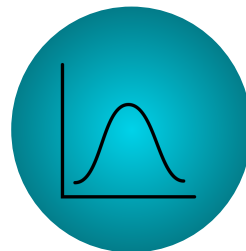
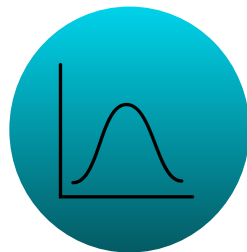
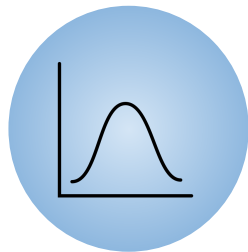
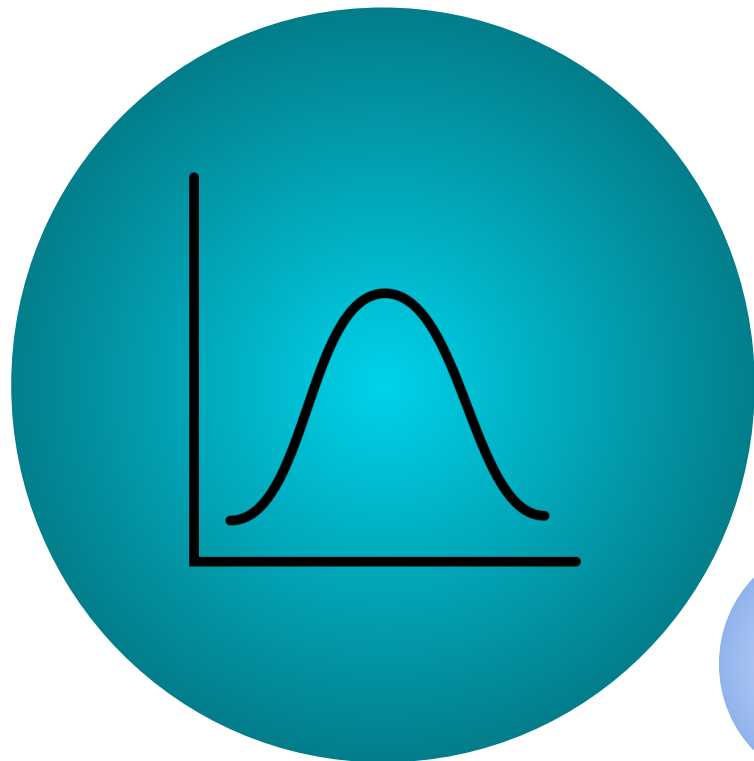
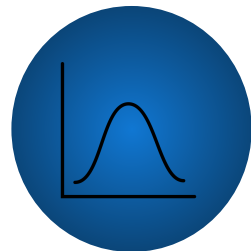
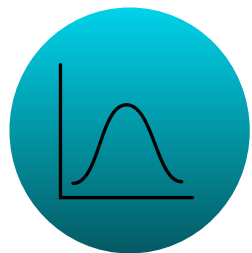
Sample distribution

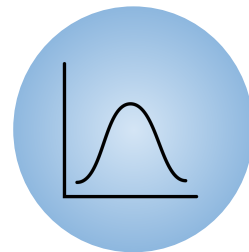
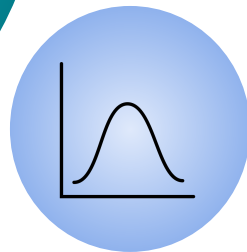
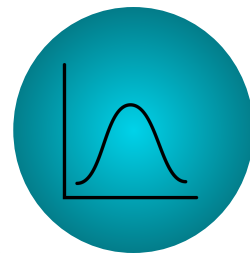
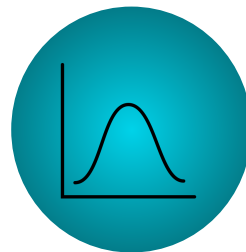
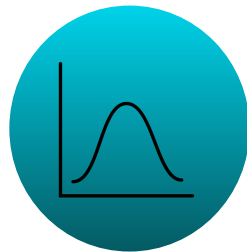
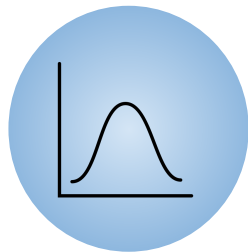
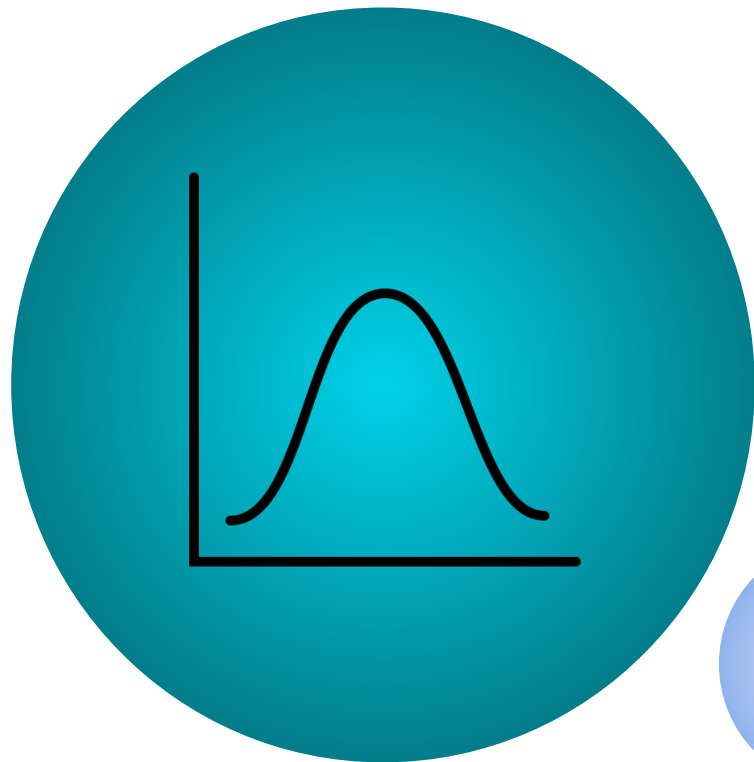
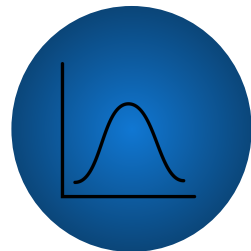
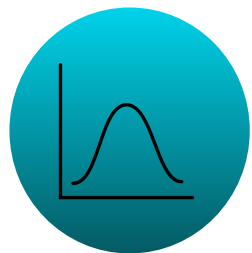


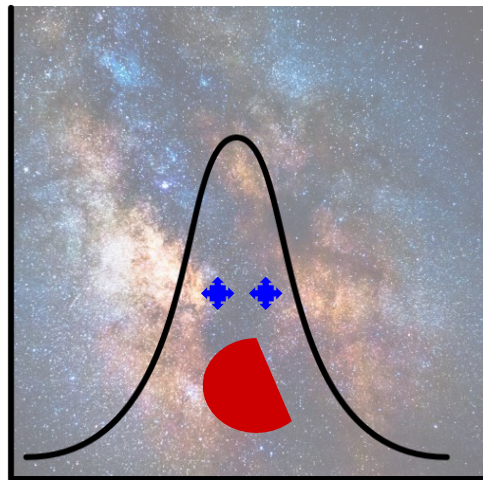
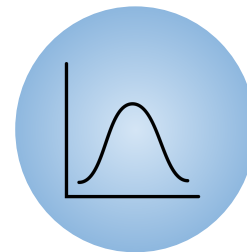
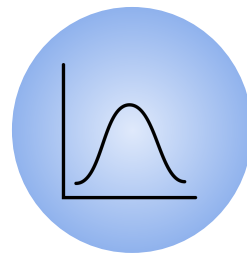
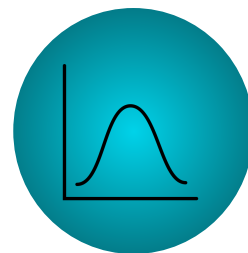
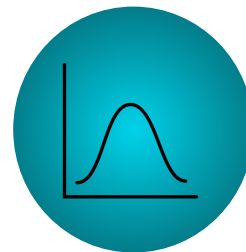
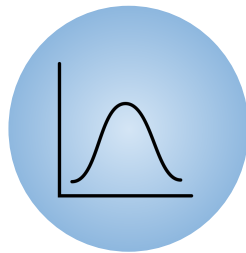
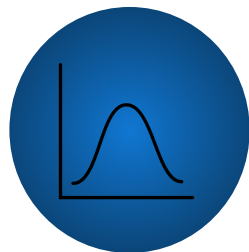
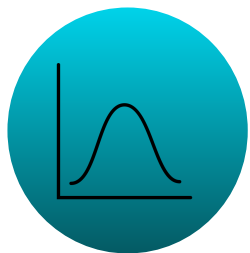
Population distribution



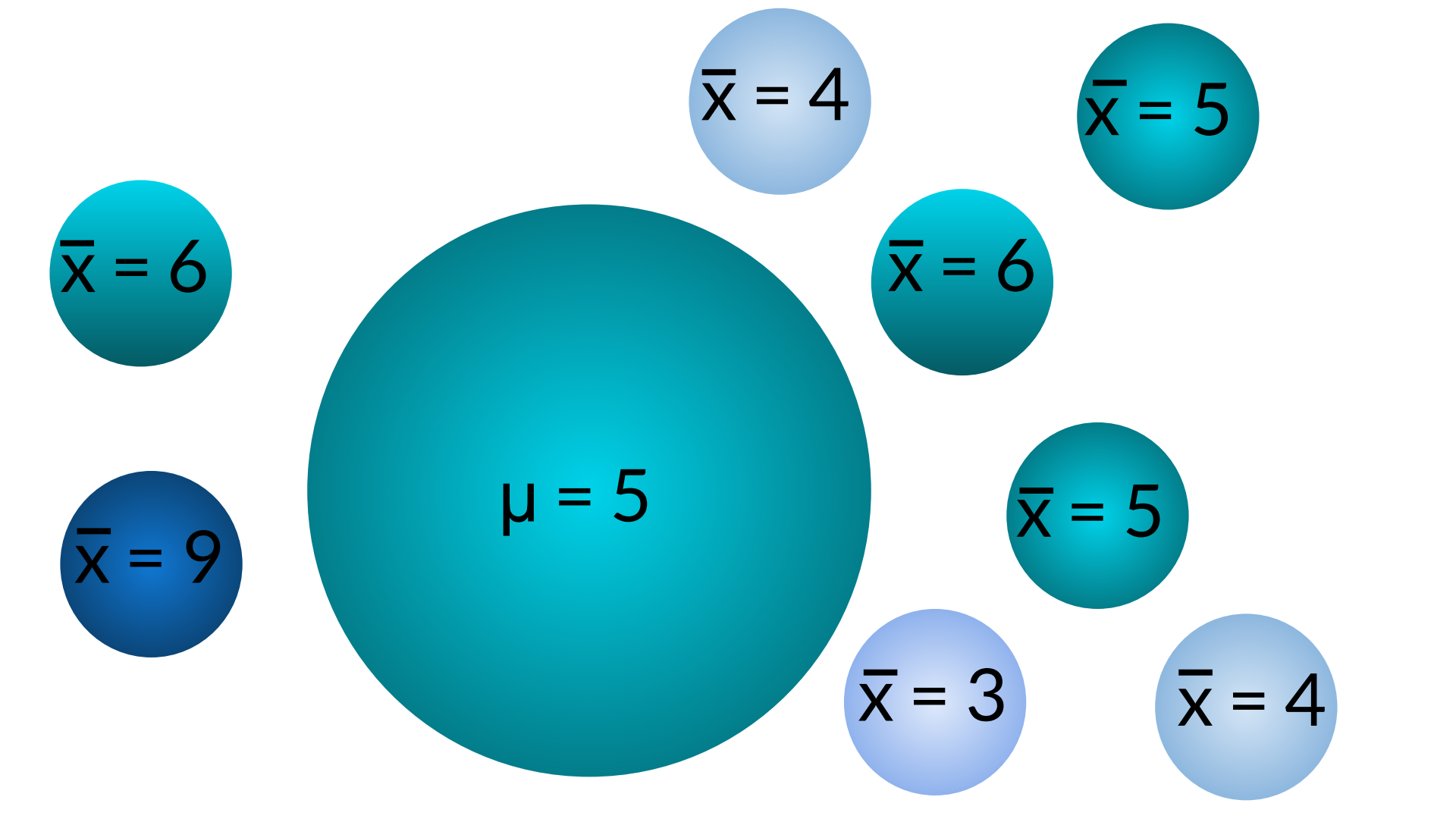
Sample distribution







Sampling
distribution



A diagram illustrating a central value $\mu = 5$ surrounded by seven other values \bar{x} . The central circle is the largest and is a dark teal color. The surrounding circles are smaller and have varying shades of teal and blue. The values of \bar{x} are: 6 (top-left), 4 (top), 5 (top-right), 6 (middle-right), 5 (bottom-right), 4 (bottom), 3 (bottom-left), and 9 (far-left).

$$\bar{x} = 6$$

$$\bar{x} = 4$$

$$\bar{x} = 5$$

$$\bar{x} = 6$$

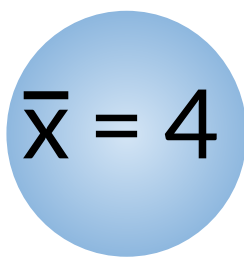
$$\bar{x} = 9$$

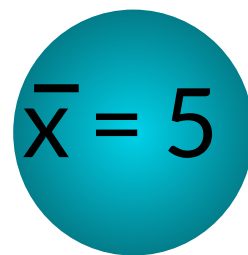
$$\mu = 5$$

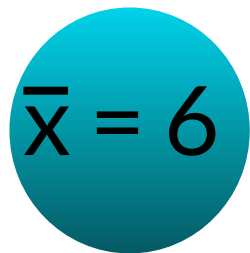
$$\bar{x} = 5$$

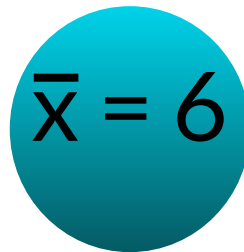
$$\bar{x} = 3$$

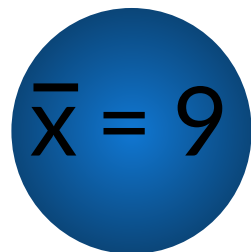
$$\bar{x} = 4$$

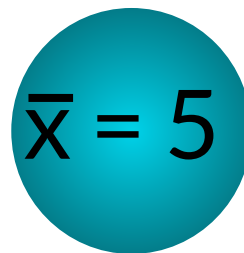

$$\bar{x} = 4$$

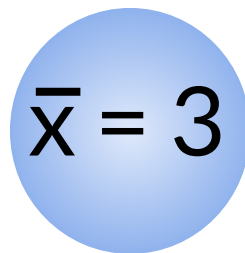

$$\bar{x} = 5$$

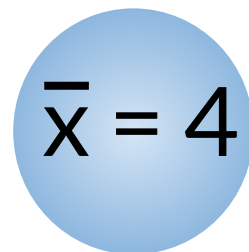

$$\bar{x} = 6$$

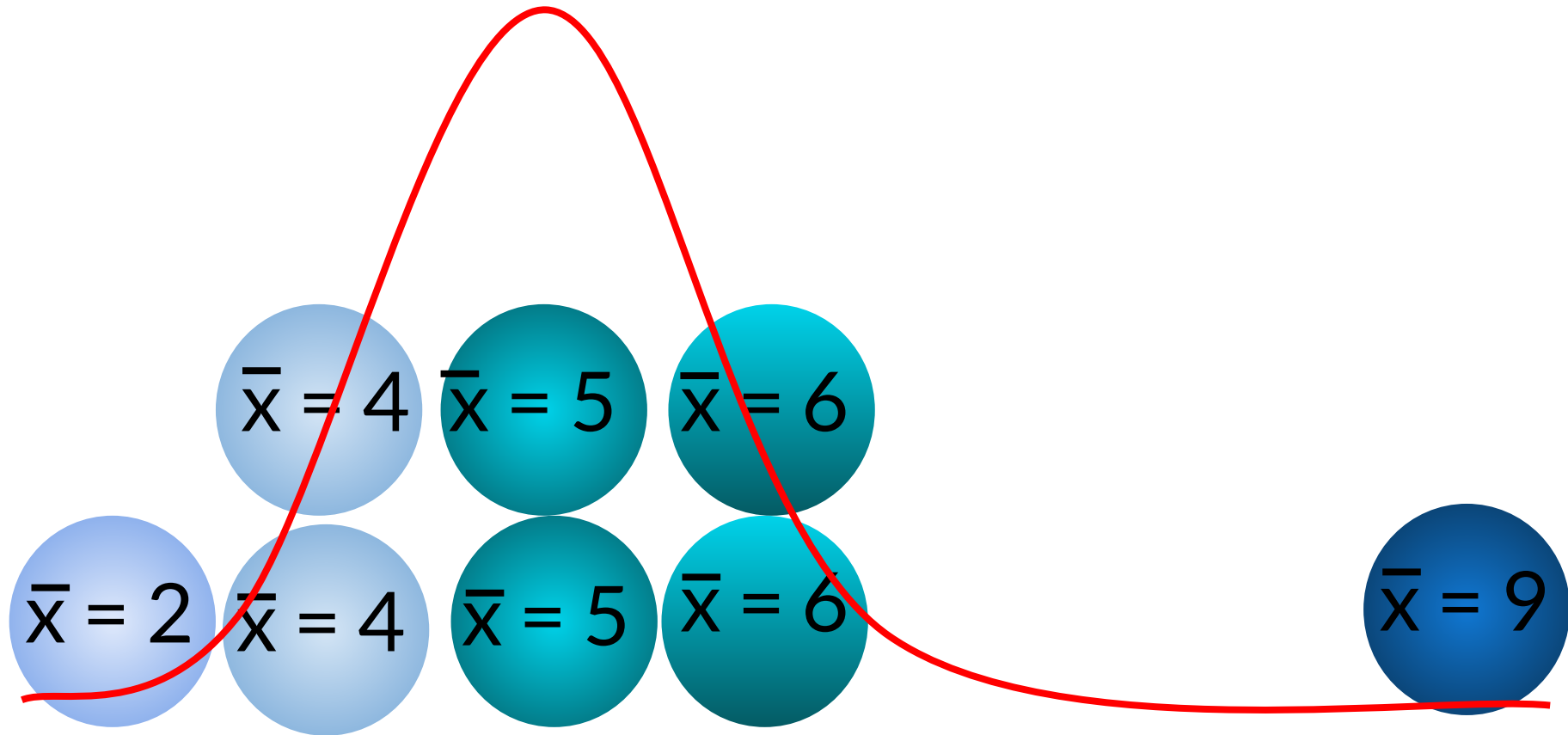

$$\bar{x} = 6$$

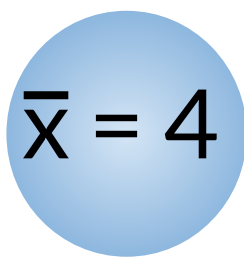

$$\bar{x} = 9$$

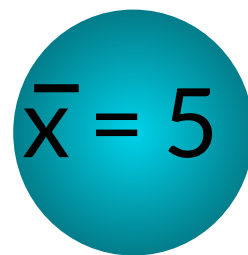

$$\bar{x} = 5$$

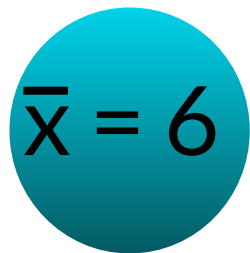

$$\bar{x} = 3$$

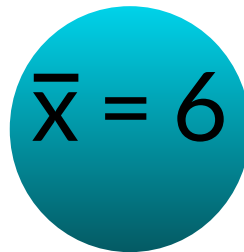

$$\bar{x} = 4$$

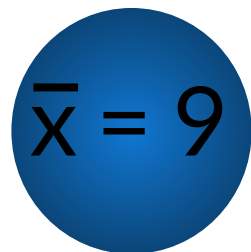


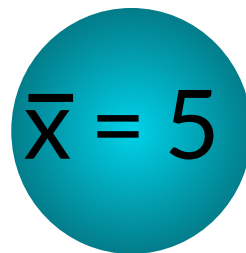

$$\bar{x} = 4$$

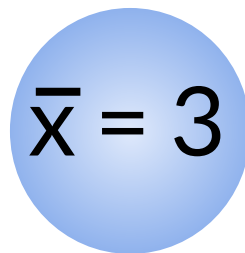

$$\bar{x} = 5$$

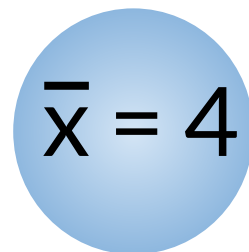

$$\bar{x} = 6$$

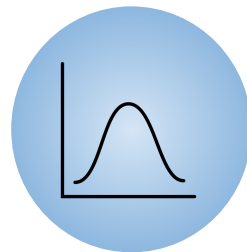
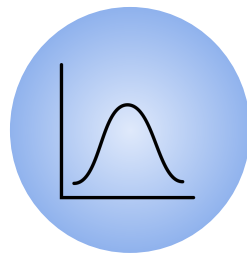
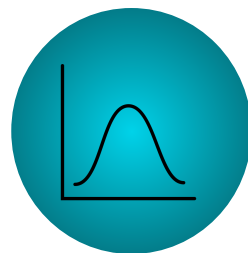
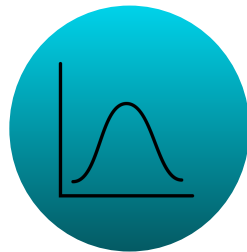
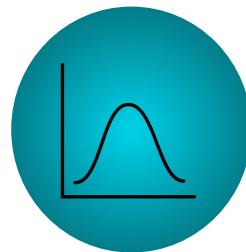
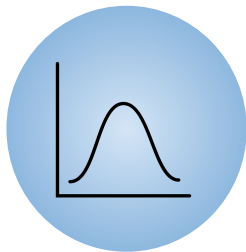
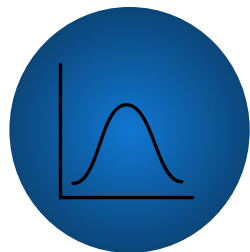
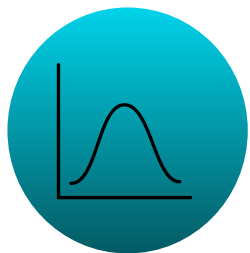

$$\bar{x} = 6$$


$$\bar{x} = 9$$


$$\bar{x} = 5$$


$$\bar{x} = 3$$


$$\bar{x} = 4$$



http://onlinestatbook.com/stat_sim/sampling_dist/

Or search: onlinestatbook sampling

5m Questions:

1. How does the sampling distribution compare to the population distribution (normal/uniform/skewed)?
2. What happens when we change the sample size used to generate the sampling distribution?
3. What custom population distribution produces a non-normal sampling distribution?

Now we want to know if dogs are cuter than cats.



What could we do?

Now we want to know if dogs are cuter than cats.



We could:

1. Get two samples: one of dogs and one of cats.
2. Find the mean for dogs and the mean for cats.
3. Take the difference. (**Signal** = $\text{mean}_{\text{dogs}} - \text{mean}_{\text{cats}}$)

Now we want to know if dogs are cuter than cats.



We could:

1. Get two samples: one of dogs and one of cats.
2. Find the mean for dogs and the mean for cats.
3. Take the difference. (**Signal** = $\text{mean}_{\text{dogs}} - \text{mean}_{\text{cats}}$)



Testing the null hypothesis

Null hypothesis (H_0): There is **no difference** in the cuteness of dogs and cats.

Testing the null hypothesis

Null hypothesis (H_0): There is **no difference** in the cuteness of dogs and cats.

Alternative hypothesis (H_A): There **is** a difference in the cuteness of dogs and cats.

Testing the null hypothesis

Null hypothesis (H_0): There is **no difference** in the cuteness of dogs and cats.

Alternative hypothesis (H_A): There **is** a difference in the cuteness of dogs and cats.

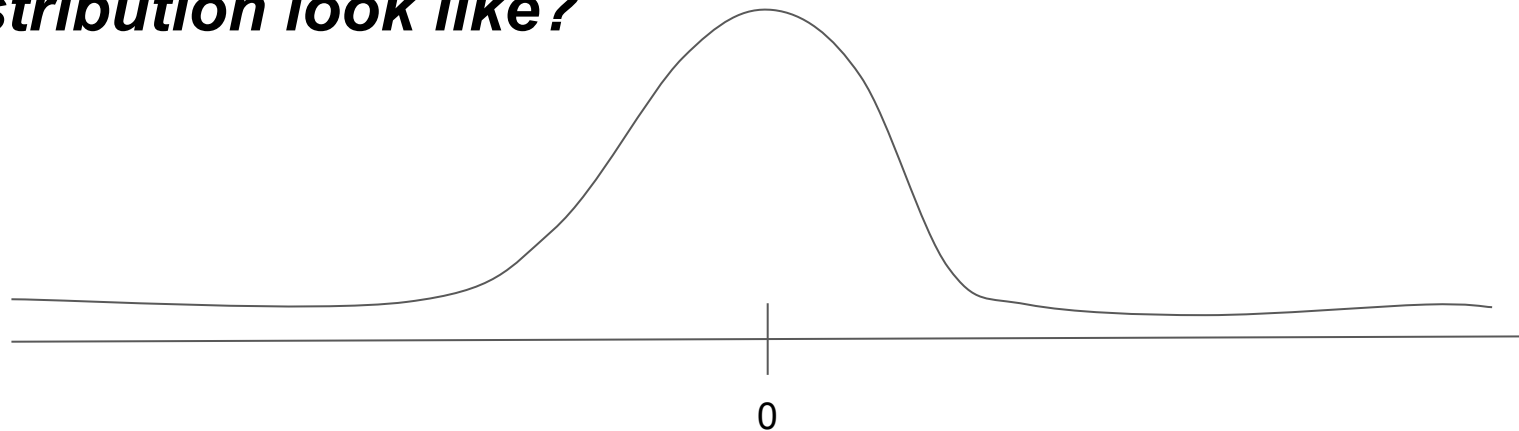
IF the null hypothesis were true, what would the sampling distribution look like?

Testing the null hypothesis

Null hypothesis (H_0): There is **no difference** in the cuteness of dogs and cats.

Alternative hypothesis (H_A): There **is** a difference in the cuteness of dogs and cats.

IF the null hypothesis were true, what would the sampling distribution look like?

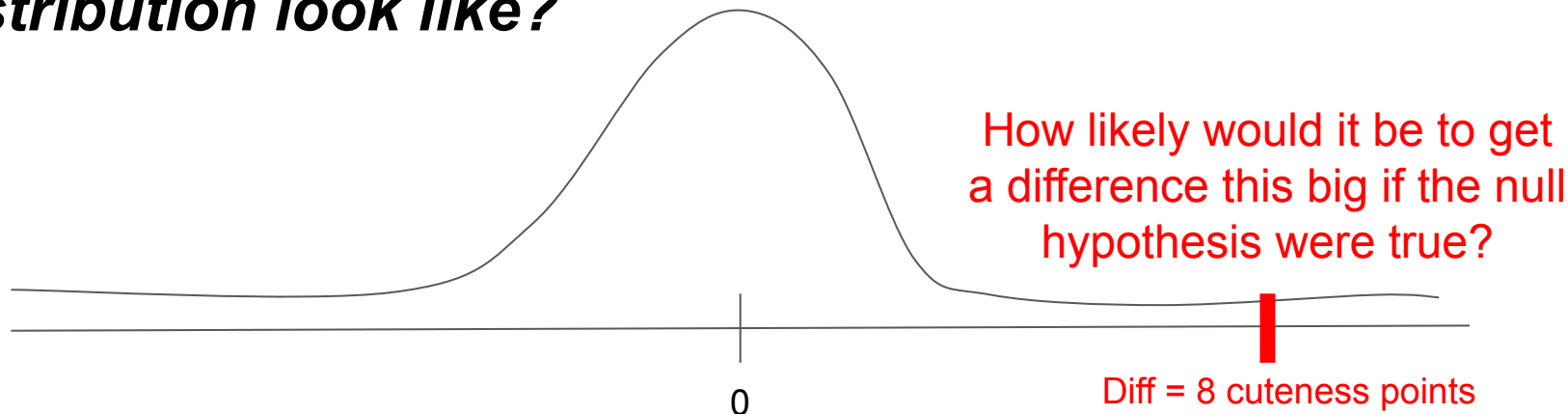


Testing the null hypothesis

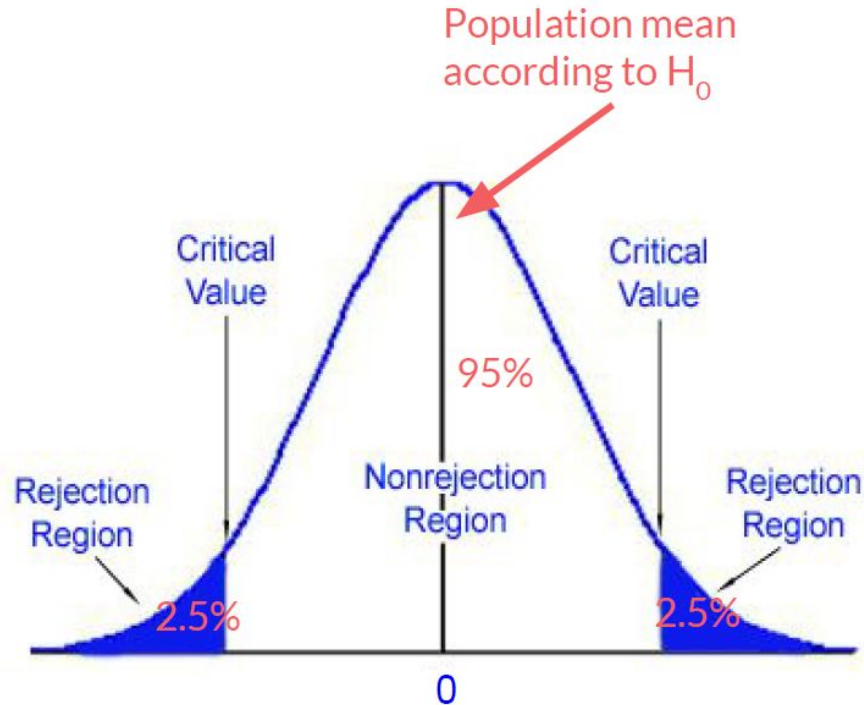
Null hypothesis (H_0): There is **no difference** in the cuteness of dogs and cats.

Alternative hypothesis (H_A): There **is** a difference in the cuteness of dogs and cats.

IF the null hypothesis were true, what would the sampling distribution look like?



Sampling distribution, with a standardized scale



(t-distribution, or Normal distribution)

Testing the null hypothesis

How likely would it be to get a difference this big if the null hypothesis were true?

If the probability < 0.05 , we say it would be UNLIKELY.

→ We reject the null hypothesis.

Otherwise...

→ We FAIL to reject the null hypothesis.

Testing the null hypothesis

How likely would it be to get a difference this big if the null hypothesis were true?

If the probability < 0.05 , we say it would be UNLIKELY.

- We reject the null hypothesis.

- We conclude there is a ***statistically significant*** difference in the cuteness of dogs and cats.

Otherwise...

- We FAIL to reject the null hypothesis.

- We cannot conclude there is any significant difference in the cuteness of dogs and cats.

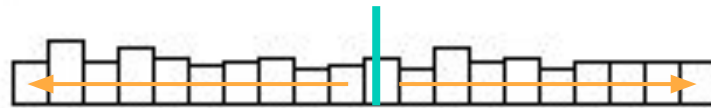
Pros and cons of using $p < 0.05$ as a cutoff:

Pros and cons of using $p < 0.05$ as a cutoff:

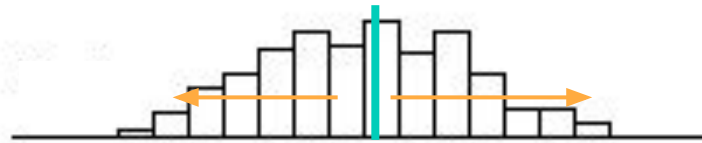
Some recent challenges to the null hypothesis significance test:

- Arbitrariness of 95% confidence intervals
- The difference between significant and not significant is itself not statistically significant
- The null hypothesis that $\mu_1 = 0$ is frequently literally incredible. This becomes more and more apparent as datasets get larger and larger.

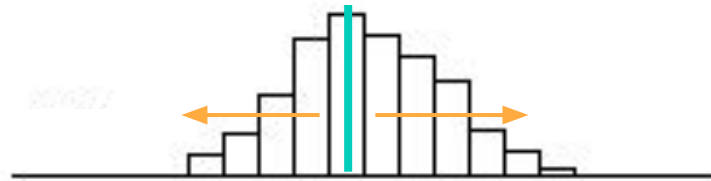
75% 50



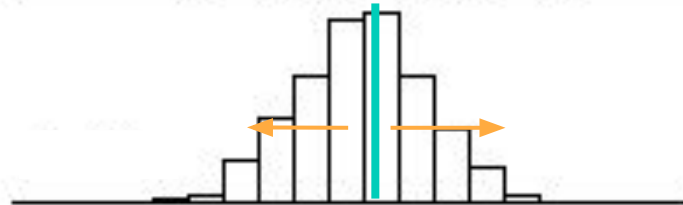
75% 10



75% 20



75% 40



Lunch



Managing projects






Goals for Project Management

- Practical:
 - ...organize programming files to help future you and others.
 - ...organize yourself for 381
- Conceptual:
 - ...intentionally develop your style.
 - ...imagine stats projects as rivers, trees, or webs.

Get some data!

1. Go to www.kaggle.com/datasets
2. Make an account (sorry)
3. Use the search filter to download a .csv file < 2MB



hotaling_cocktails - Cocktails.csv (252.87 KB)						
9 of 9 columns ▾ Views     						
	▲ Cocktail Name ▼ Cocktail name	▲ Bartender ▼ Bartender who created this cocktail (Optional)	▲ Bar/Company ▼ Bar or company the bartender is associated with (Optional)	▲ Location ▼ Location of the company (Optional)	▲ Ingredients ▼ Ingredients and quantities, comma-separated	▲ Garnish Garnishes separated
	684 unique values	[null] 34% Francesco Lafran... 6% Other (249) 60%	[null] 61% Dirty Habit 4% Other (167) 35%	[null] 50% San Francisco 23% Other (40) 27%	686 unique values	[null] Luxardo (Other (31
1	Flor de Amaras	Kelly McCarthy		Boston	1.5 oz Mezcal, 1 oz Hibiscus Simple Syrup*, .5 oz Lime Juice. too Soda	Marigold

Make a folder structure and project

1. First make a project folder with the name of your project
2. Inside your project folder include folders for:
 - a. Raw data, Clean data
 - b. Cleaning, Analysis
 - c. Graphs, Logs
 - d. Archive
3. Create a new Project in Stata in your project folder
4. Drag your folder structure into the project manager!

Convert you raw data to .dta

1. Move your .csv data to your Raw Data folder
2. Create a new .do file in your Cleaning folder
3. Add starter code and comment block
4.

```
import delimited "Raw Data/YOUR_CSV.csv",  
varnames(1)
```
5.

```
save "Clean Data/NICE_DATA.dta", replace
```
6. Drag your nice new data from the finder into your project manager folder tree

Analyze!

1. Generate summary statistics for a couple variables
2. Make 2 histograms for two interesting continuous variables
 - a. `histogram varname`
3. Output each figure to your Plots folder
 - a. `graph save "Graphs/filename", replace`
4. Make a scatter plot of the two variables and save to your graphs folder.
5. Play with the advanced stata graphics in the Day 2 folder!

TWIST

1. Add two different mistakes to your analysis .do file
2. Add a to-do item [be merciful and don't make it too hard]
 - e.g., summarize a var, generate a new var, plot a graph
3. Compress your project folder, share it to a partner you haven't worked with yet.
4. Open the file, find the errors, complete the to-do!

Reflect/Radiate

<https://forms.gle/UwyPtY6kFEPgm47JA>

