

Sampling, bias and sampling distributions

Overview

Review of simple linear regression

Sampling and bias

Sampling distributions

Review of linear regression

Regression

Regression is method of using one variable x to predict the value of a second variable y

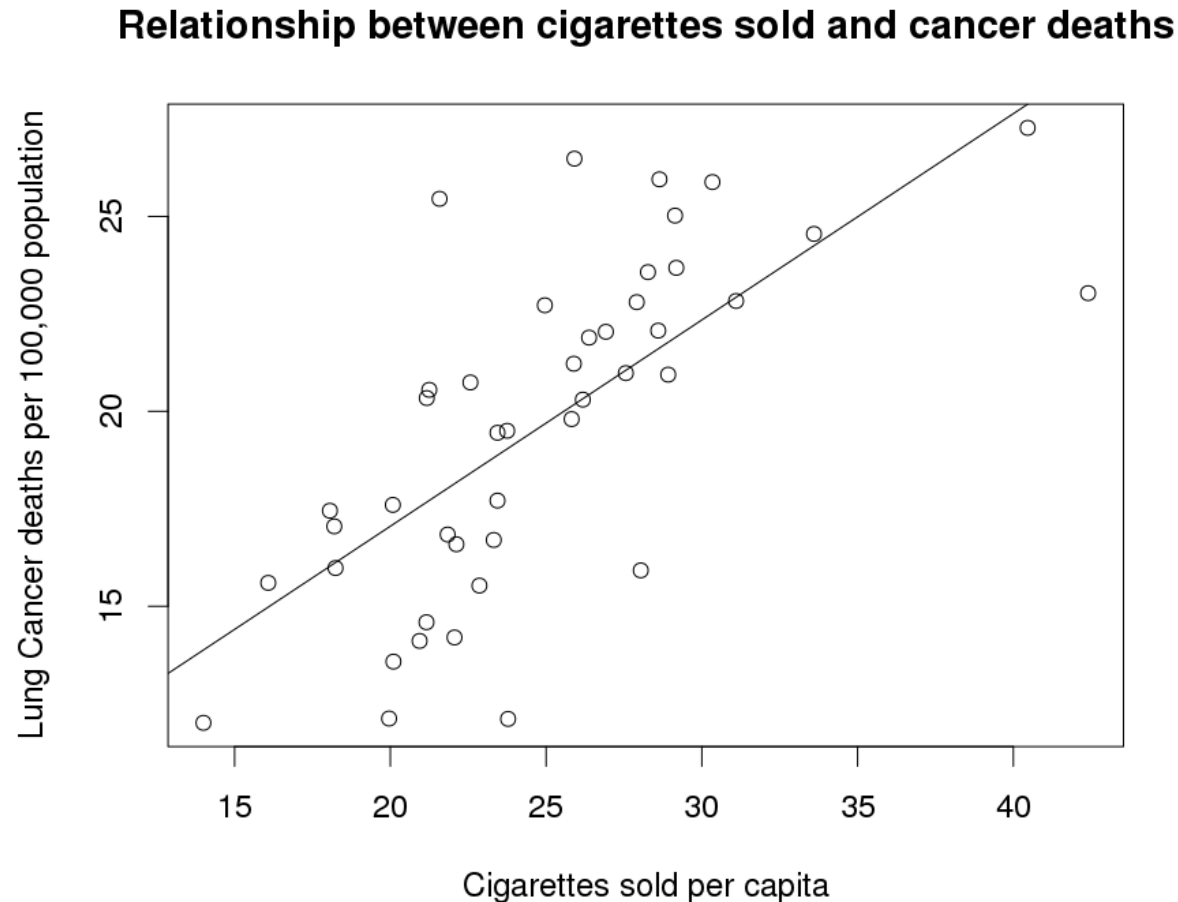
- i.e., $\hat{y} = f(x)$

In **linear regression** we fit a line to the data, called the **regression line**

$$\hat{y} = a + b \cdot x$$

$$\text{Response} = a + b \cdot \text{Explanatory}$$

Cancer smoking regression line



$$\hat{y} = a + b \cdot x$$

R: `my_fit <- lm(y ~ x)`
`coef(my_fit)`

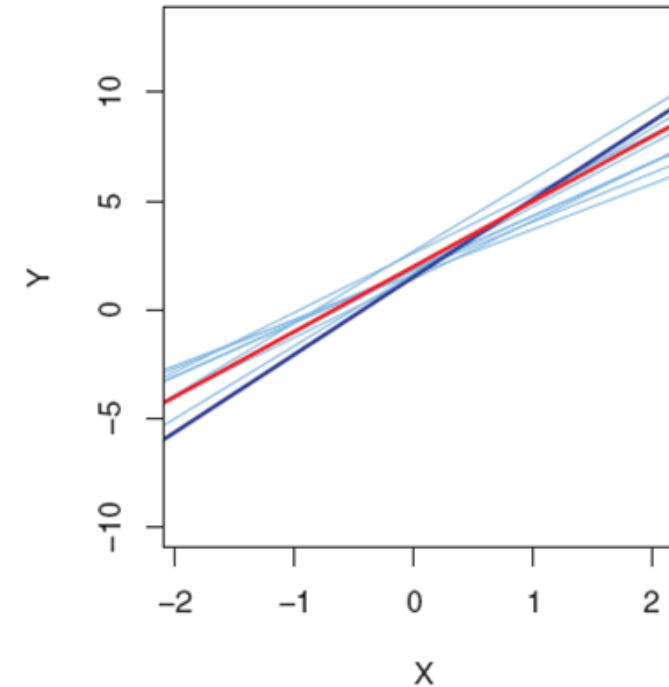
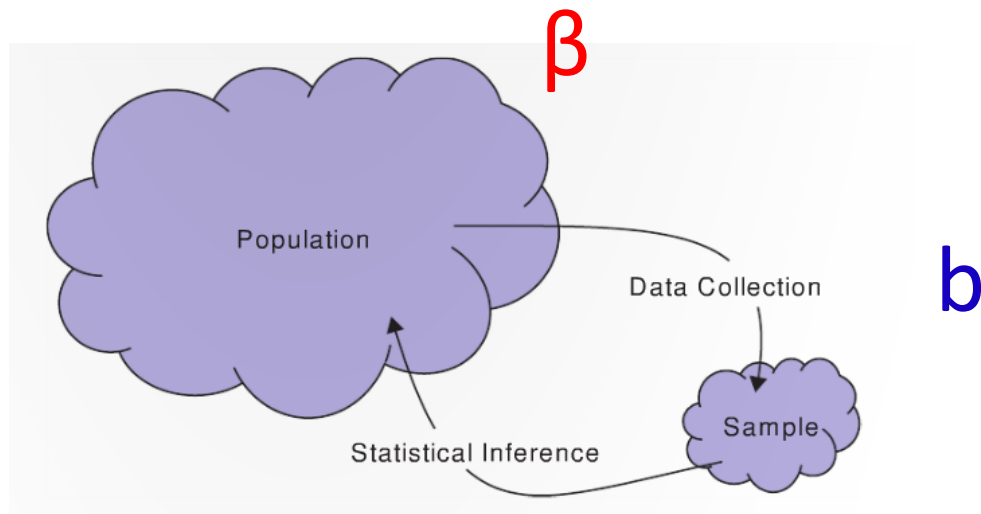
$$a = 6.47 \quad b = 0.53$$

$$\hat{y} = 6.47 + .53 \cdot x$$

Notation

The Greek letter β is used to denote the slope of the **population**

The letter b is typically used to denote the slope of the **sample**



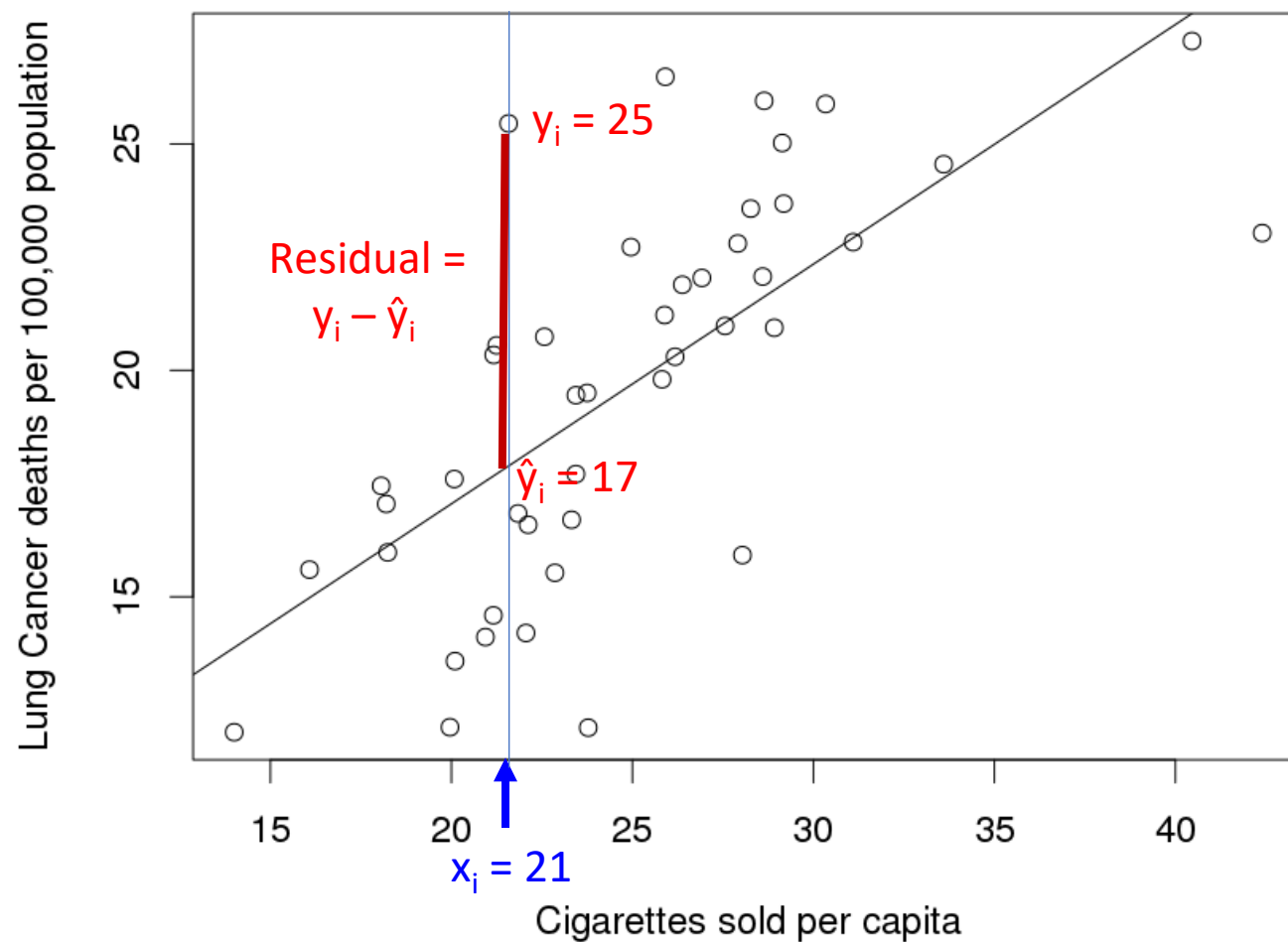
Residuals

The **residual** is the difference between an observed (y_i) and a predicted value (\hat{y}_i) of the response variable

$$Residual_i = Observed_i - Predicted_i = y_i - \hat{y}_i$$

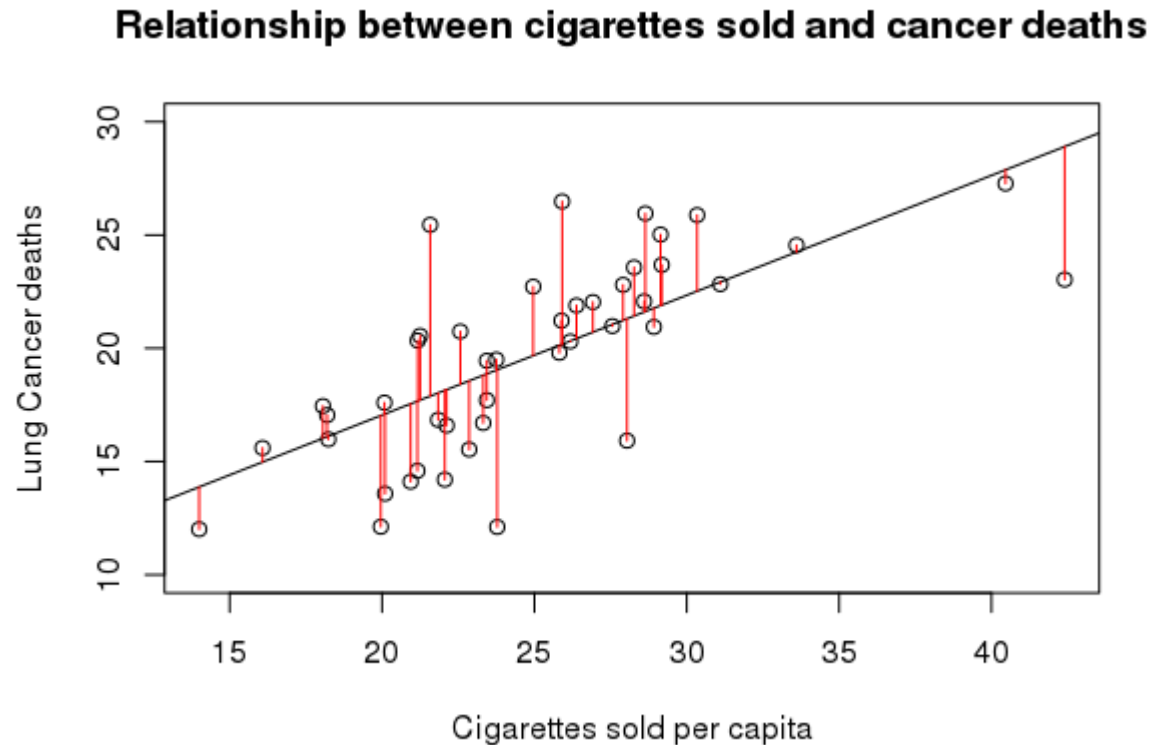
Cancer smoking residuals

Relationship between cigarettes sold and cancer deaths



Line of 'best fit'

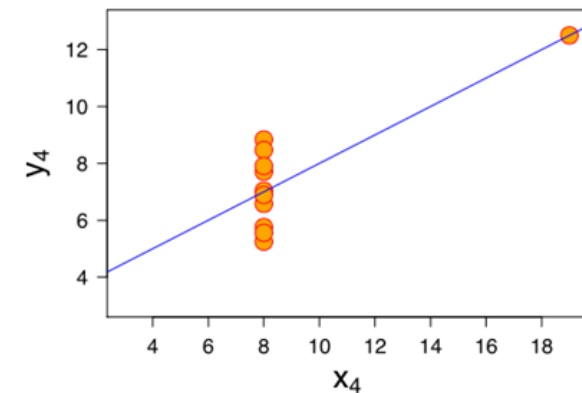
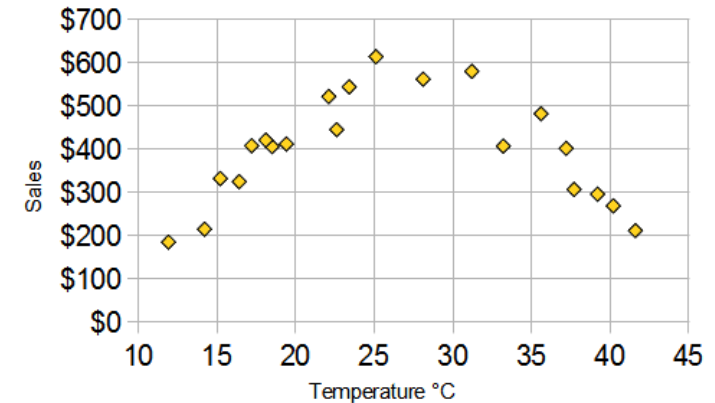
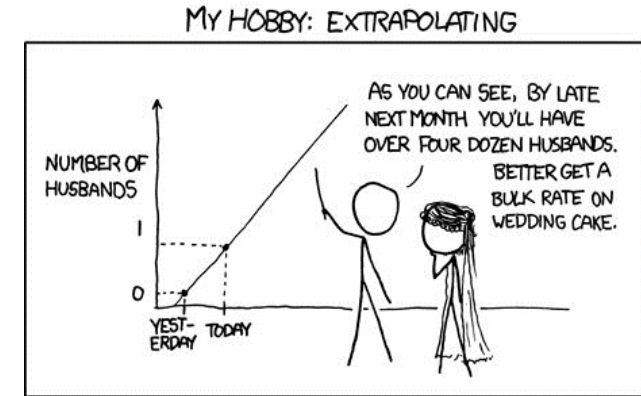
The **least squares line**, also called '**the line of best fit**', is the line which minimizes the sum of squared residuals



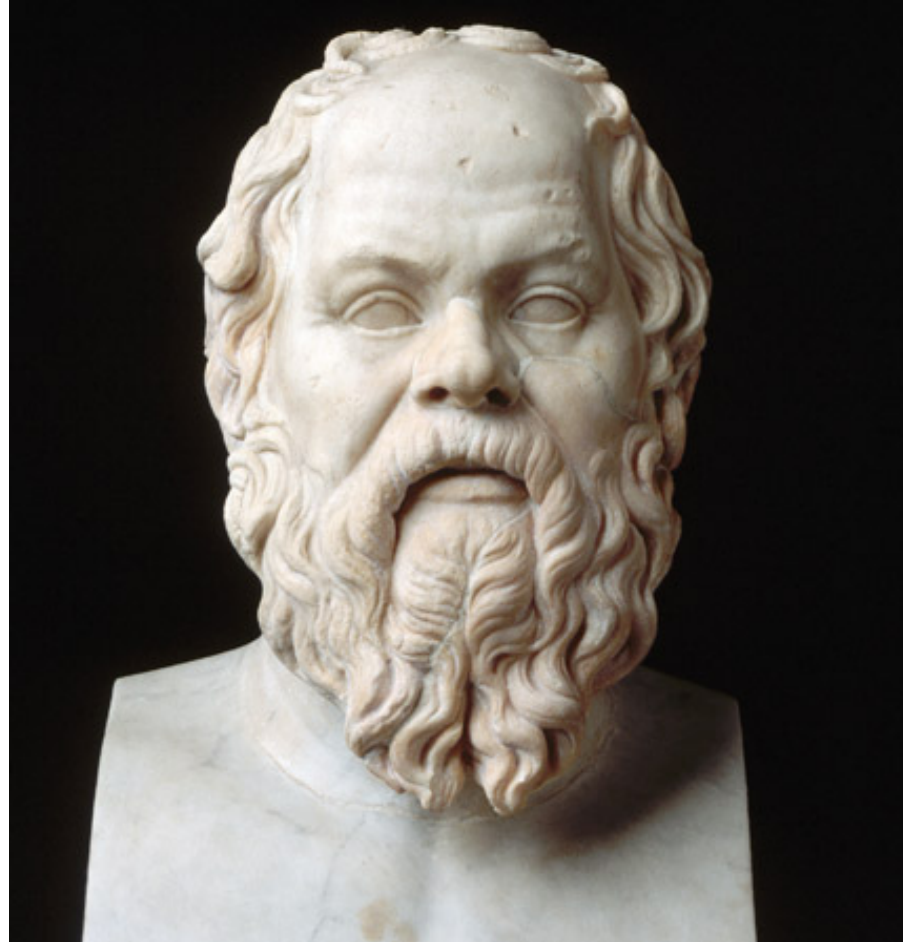
[Find the line of best fit](#)

Regression cautions

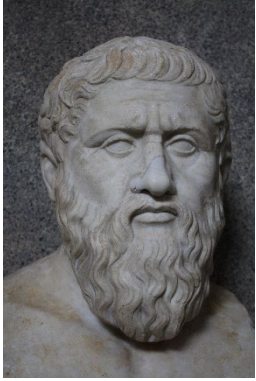
1. Avoid trying to apply the regression line to predict values far from those that were used to create the line.
2. Plot the data! Regression lines are only appropriate when there is a linear trend in the data.
3. Be aware of outliers – they can have an huge effect on the regression line.



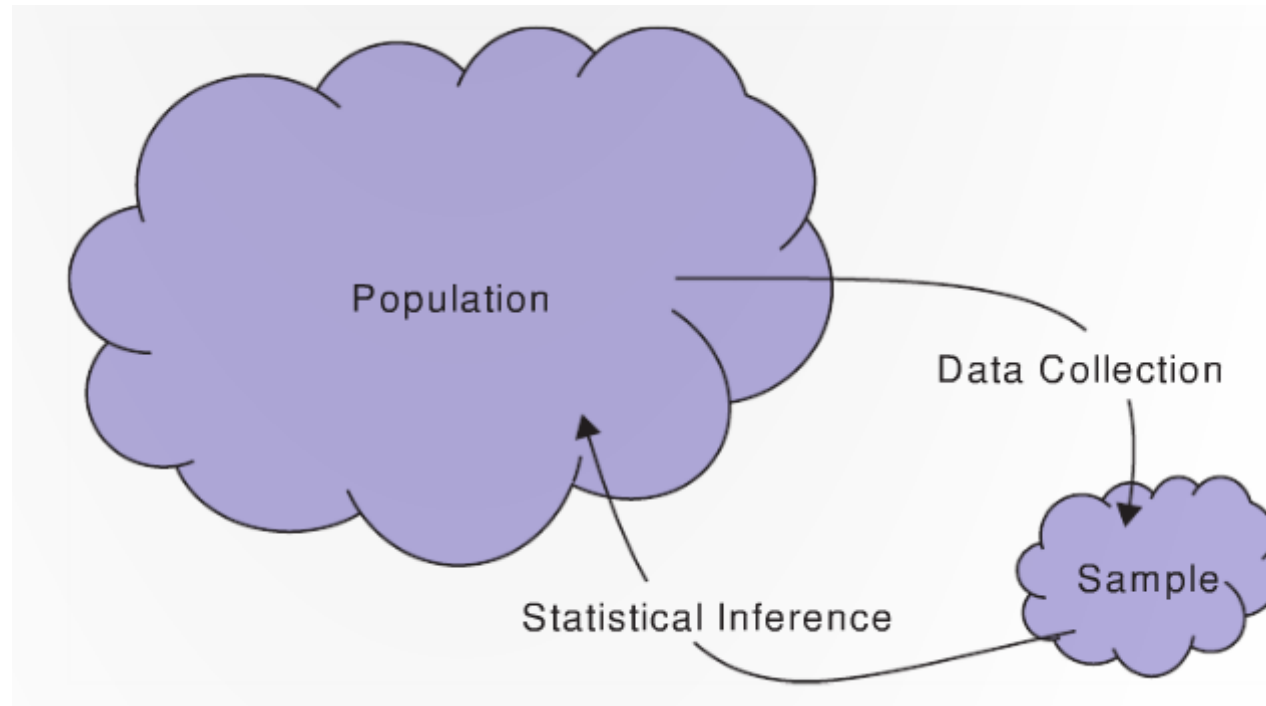
Socratic method to review descriptive statistics



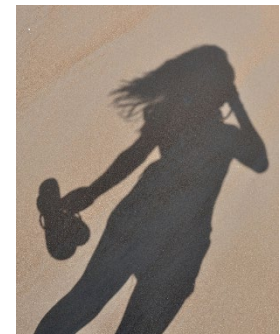
Parameters and statistics



$\pi, \mu, \sigma, \rho, \beta$



$\hat{p}, \bar{x}, s, r, b$



THE TRUTH IS OUT THERE



Can you handle the Truth[®]?

Sampling

Where do samples/data come from?



Example: sampling 100 sprinkles



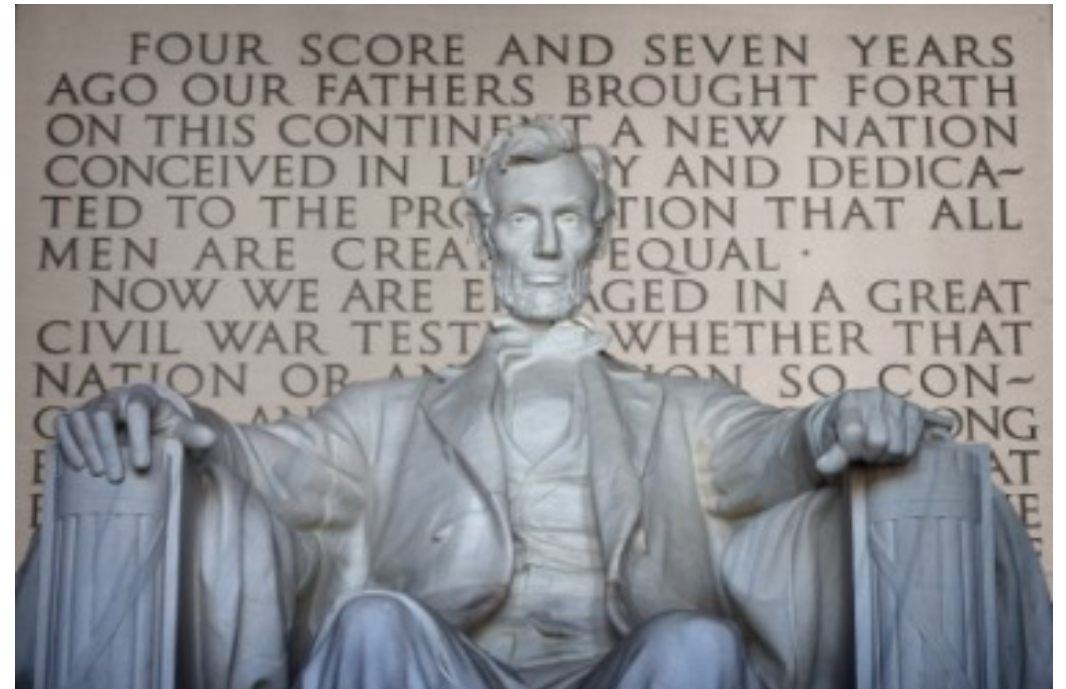
1	orange
2	red
3	green
4	white
5	white
6	white
7	white
8	white
9	red

The **sample size** (n) is the number of items in the sample
What is **n** in the sprinkle example?

Let's try some sampling ourselves...

Fill out the worksheet where you need to randomly sample 10 words from the Gettysburg address

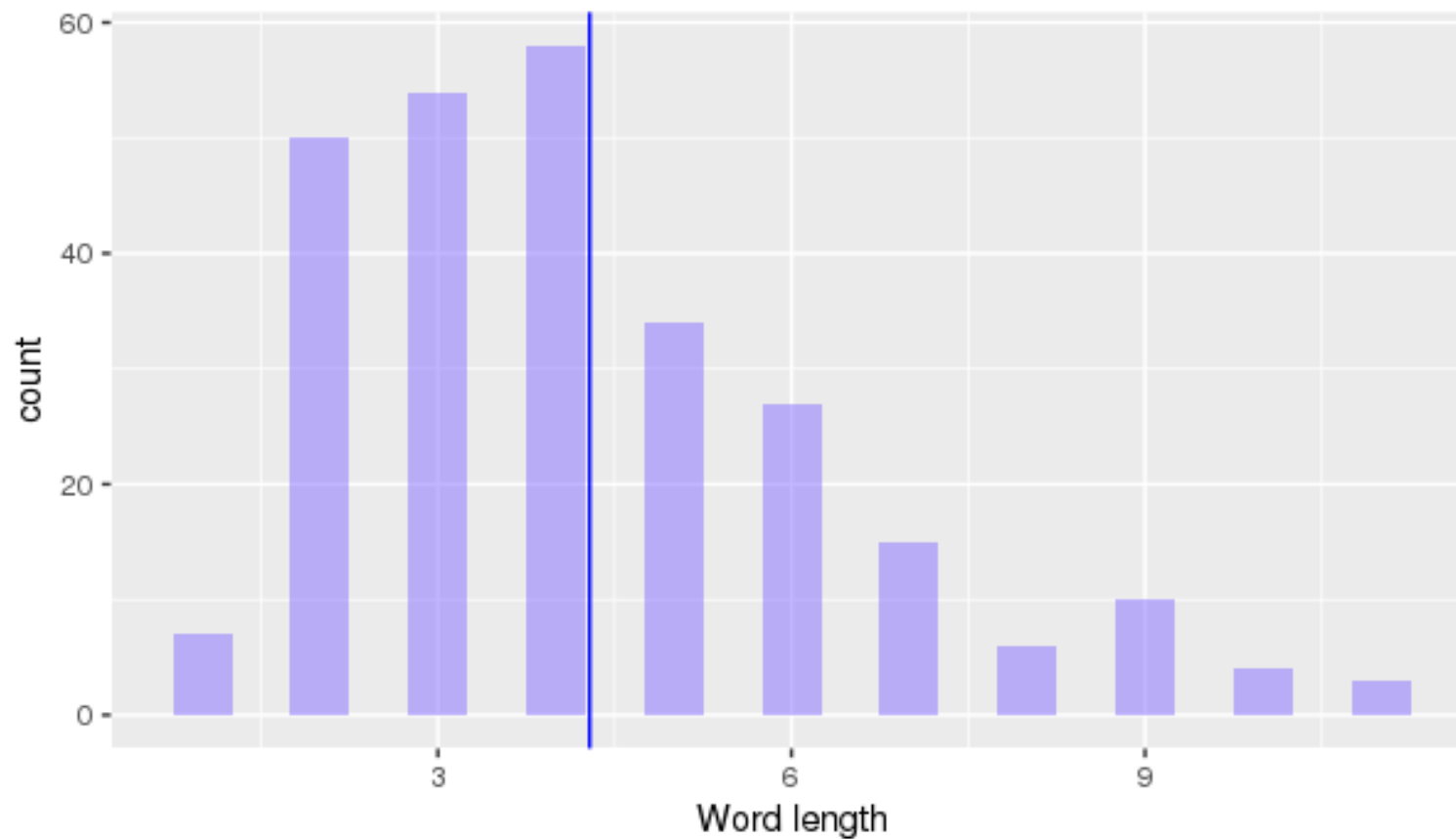
Please report the mean of the 10 words in the Canvas survey/quiz.



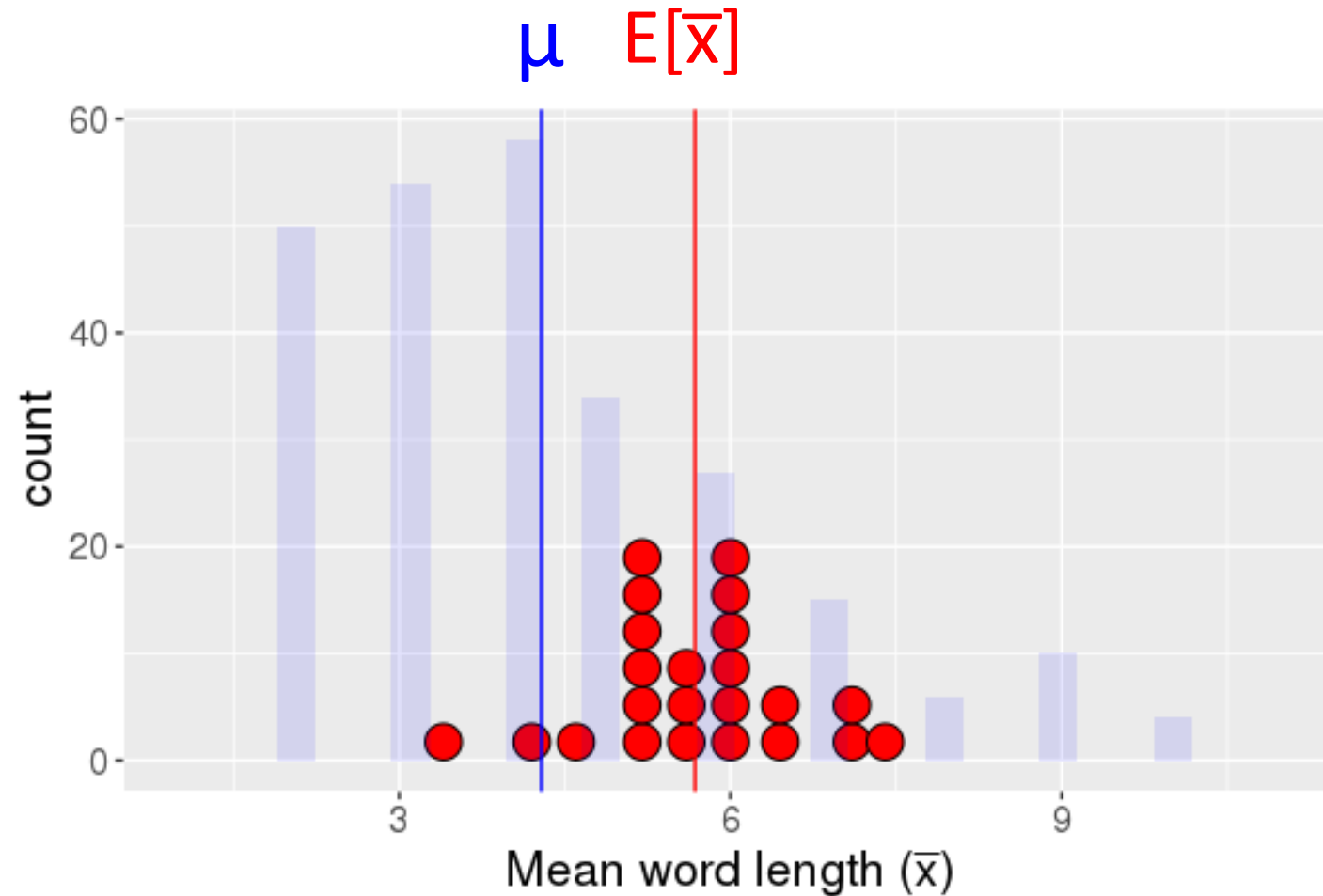
Sampling and bias

Gettysburg address: lengths of 268 words in the population

$$\mu = 4.287$$

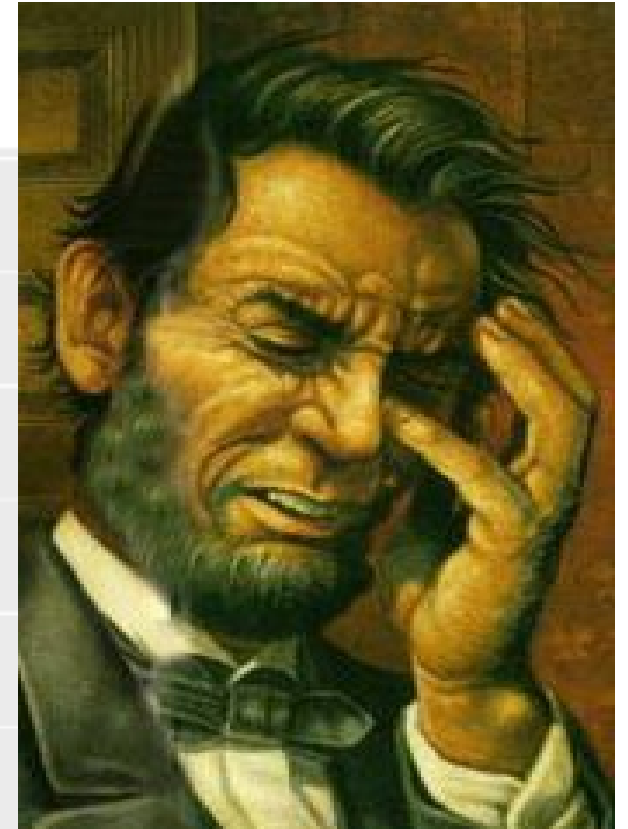
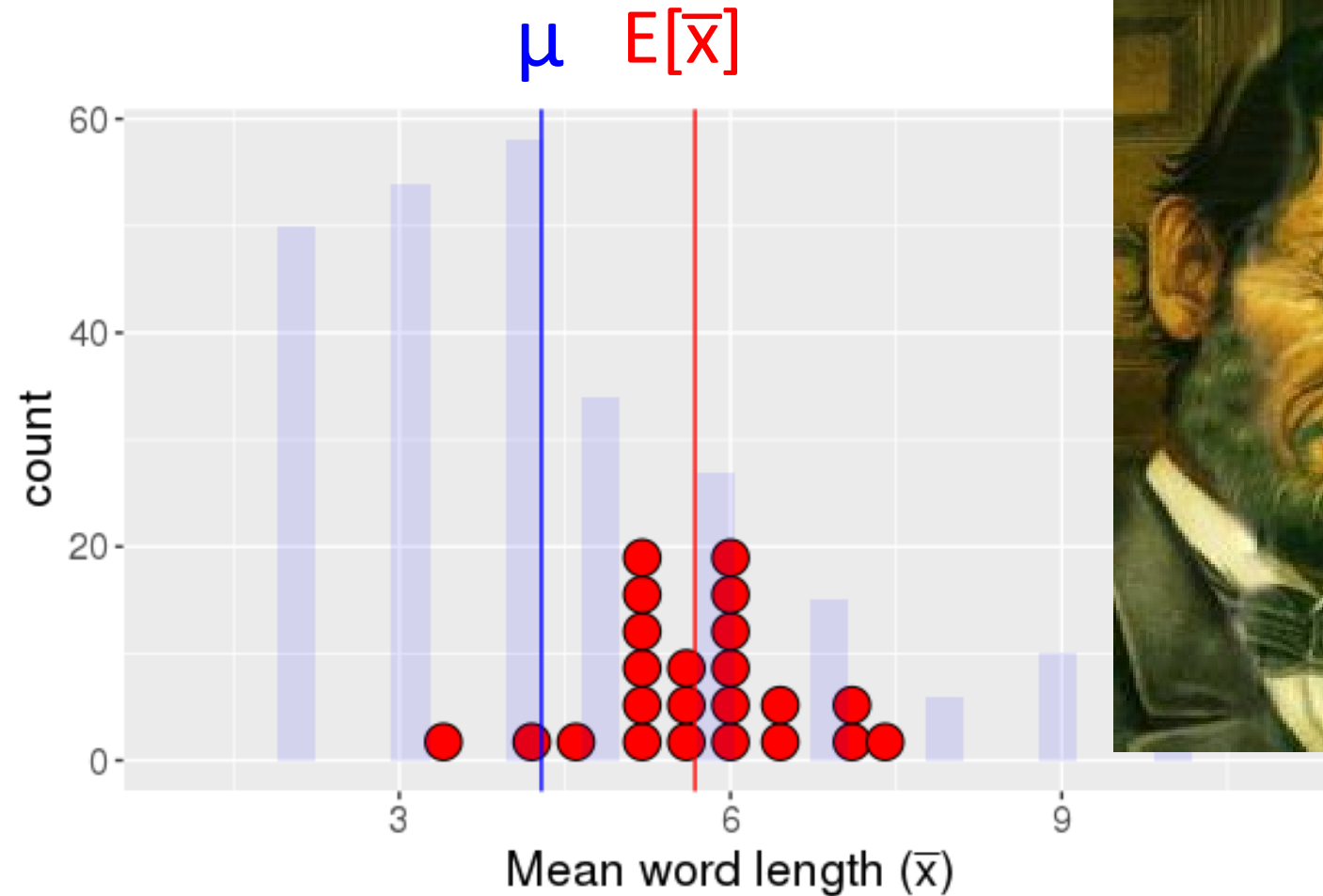


Gettysburg address, mean word length



Gettysburg address, mean word length

Observations?



Bias

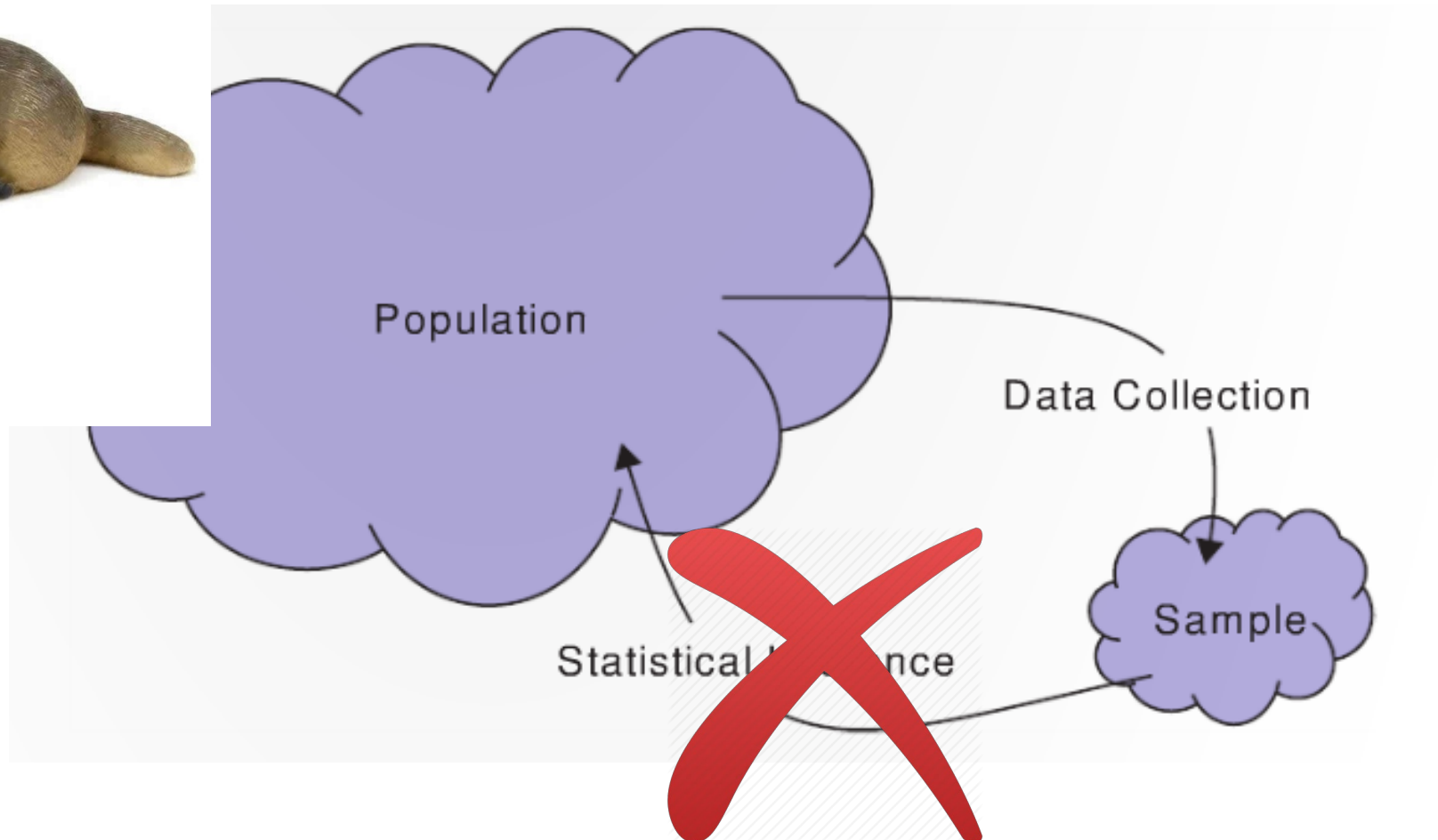
Sampling bias exists when the method of collecting the data causes the sample to inaccurately reflect the population.

This leads to ***biased statistics*** where our average statistic value does not equal the parameter value.

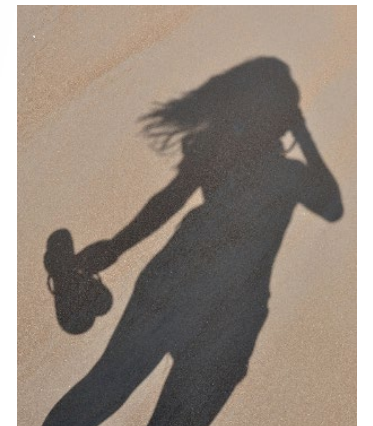
- E.g., $E[\bar{x}] \neq \mu$

Statistical bias

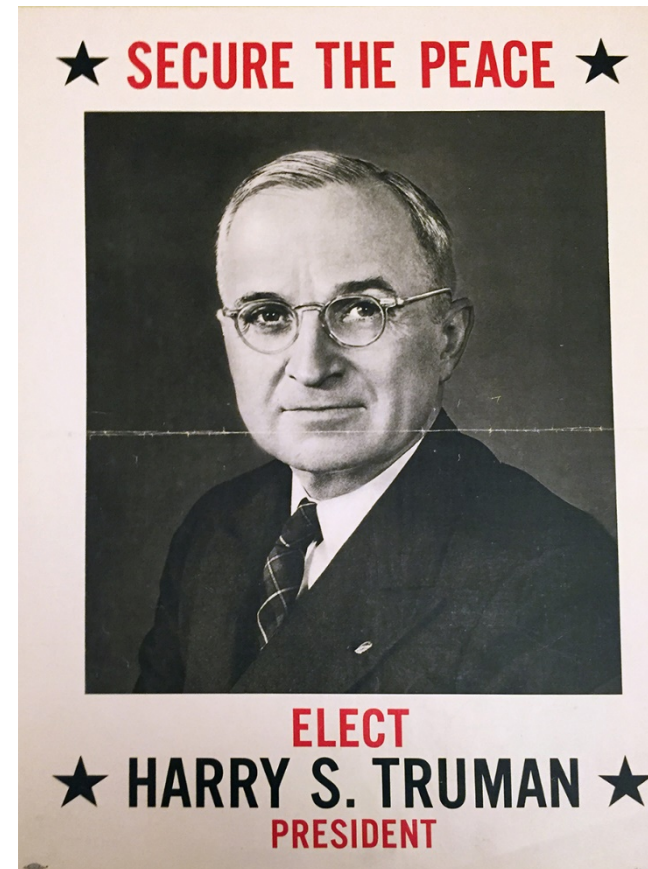
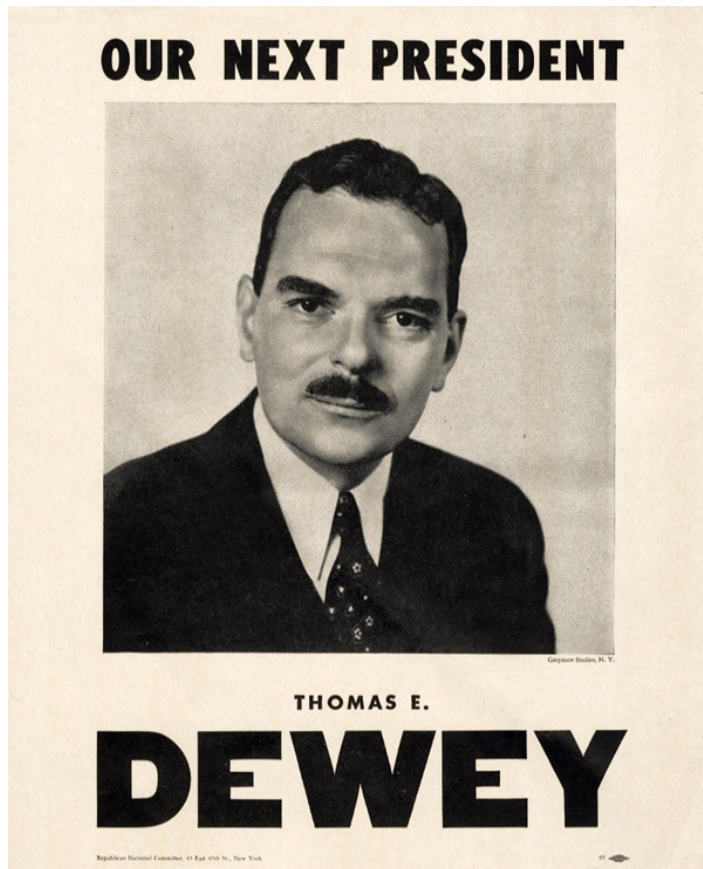
μ



\bar{x}



1948 US election



Newspaper title: Dewey Defeats Truman (1948)

The newspaper was published before the conclusion of the 1948 presidential election

The results were based on a large telephone poll which showed Dewey sweeping Truman

However, Harry S. Truman won the election

Q: What went wrong?



Basic questions for sampling

What is the population?

What is the sample?

Do they differ in a meaningful way?

Basic questions for sampling

What is the population?

What is the sample?

Do they differ in a meaningful way?

To prevent bias: use simple random sample!

Simple random sample: each member in the population is equally likely to be in the sample.

Allows for generalizations to the population!

Soup analogy



How do we select a random sample?

Mechanically:

- Flip coins

- Pull balls from well mixed bins

- Deal out shuffled cards, etc.

Use computer programs



Bias or no bias?

1948 US election: Dewey vs. Truman

Suppose there was a poll for the Truman/Dewey election that randomly chose 6,000 people from all voters in the USA and calculated who they voted for.





As part of a strategic-planning process, in spring 2013 Hampshire College launched a survey of alums.

Via email, the college invited 8,160 alums to fill out an online questionnaire administered by the campus's offices.

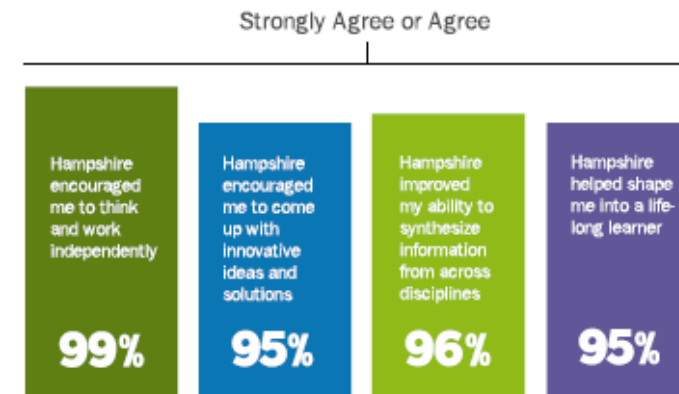
A total of 1,920 surveys were completed, yielding a response rate of 24%.

Alumni Survey Results

As part of a strategic-planning process, in spring 2013 Hampshire College launched a survey of alums. Via email, the College invited 8,160 alums to fill out an online questionnaire administered by the campus's Alumni and Family Relations and Institutional Research offices. A total of 1,920 surveys were completed, yielding a response rate of 24%.

Note: The percentages in the data (below) are based on the number of responses received for each question.

To what extent do you agree with the following statements?



Please rate your student experience at Hampshire.



65% of our alumni earn advanced degrees within ten years of graduating.

1 in 7 alumni holds a Ph.D. or other terminal degree.

Hampshire ranks in the **top 1%** of colleges nationwide in the % of grads that go on to earn doctorates.

26% of our graduates have started their own business or organization.

“

Hampshire does a great job fostering the ability to ask good questions and to look at ideas with a critical lens.

Hampshire has encouraged me to be more engaged, socially aware and more of a critical thinker than my peers.

I feel more able to adapt to a range of environments because Hampshire taught me skills and ideas rather than just knowledge.

”



Yelp reviews of restaurants?

An anonymous survey randomly select 6,000 people and asked them if have they used an illicit drug in the past month?

<https://www.billoreilly.com/poll-center>

The way you frame the question matters!

Quinnipiac University conducted two polls on November 5, 2015

First poll they asked: do you support “stricter gun control laws”?

- Yes = 46% No = 51% Difference = -5%

Second poll asked: do you support “stricter gun laws”?

- Yes = 52% No = 45% Difference = 7%

How could this affect the newspaper headlines?

- “Majority of Americans **oppose** stricter gun control laws” vs.
- “Majority of Americans **support** stricter gun laws”

Also see textbook section 1.2:

- “If you had to do it over again, would you have children?”

Practicalities...

It might not be feasible to randomly select equally from all members of a population.

This might not be a problem as long as the sample is representative of the population.

Example: If we wanted to know proportion of people left-handed in the US, randomly sampling Yale students might be good enough.

Need to think carefully to avoid bias!

Statistics requires thought!

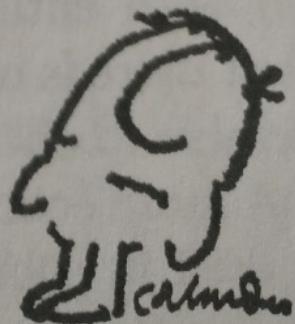
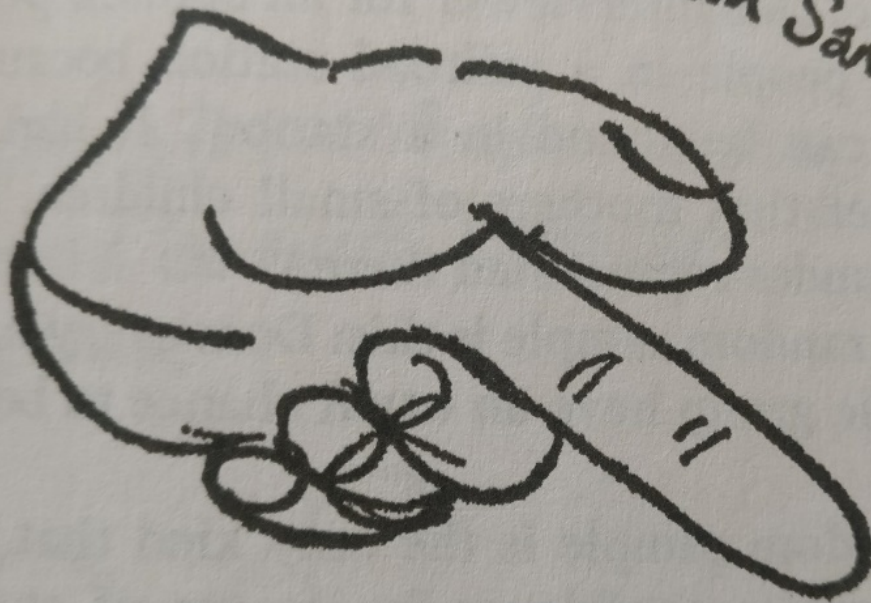
Use your own reasoning:

- What is the population I am interested in?

- Does the sample reflect the population of interest?

- Be your own worst critic!

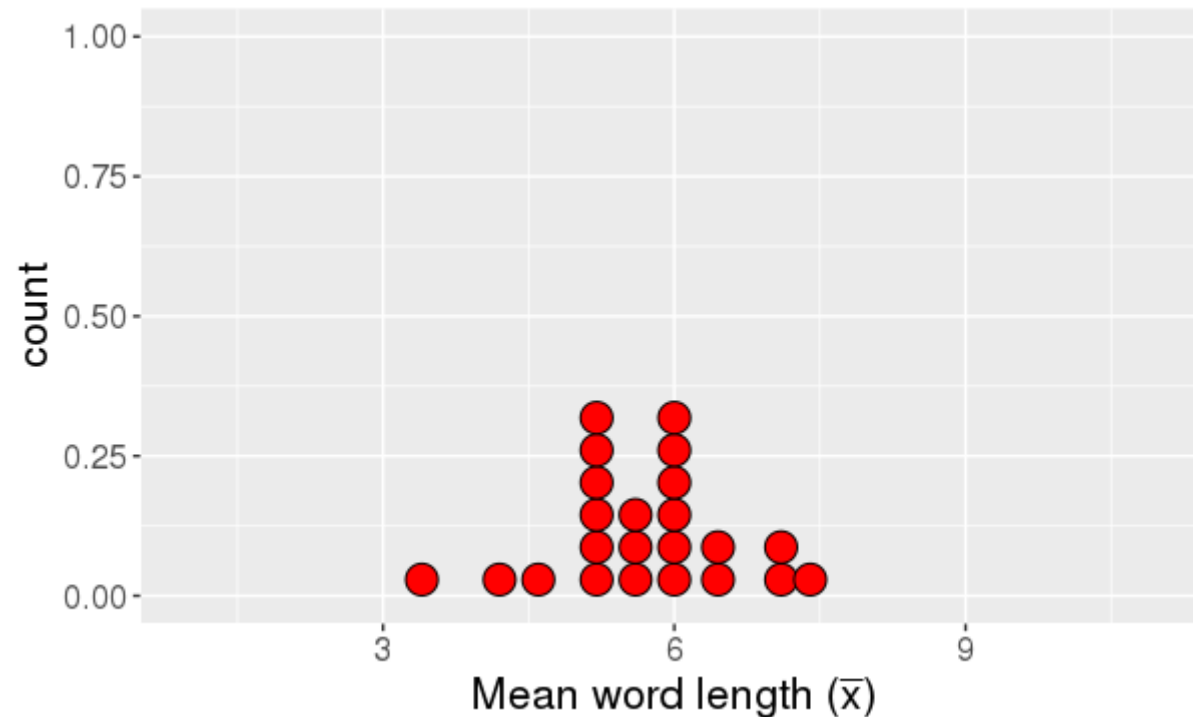
/// You've been CHOSEN
to be a Random Sample!



Sampling distributions

Recall for our distribution of Gettysburg word lengths...

Q: What does each case that is plotted correspond to?



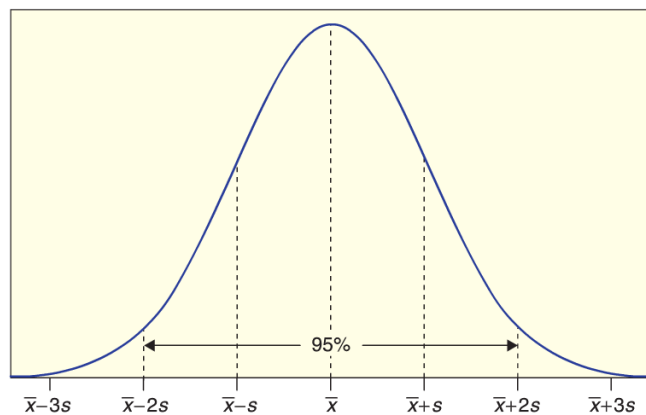
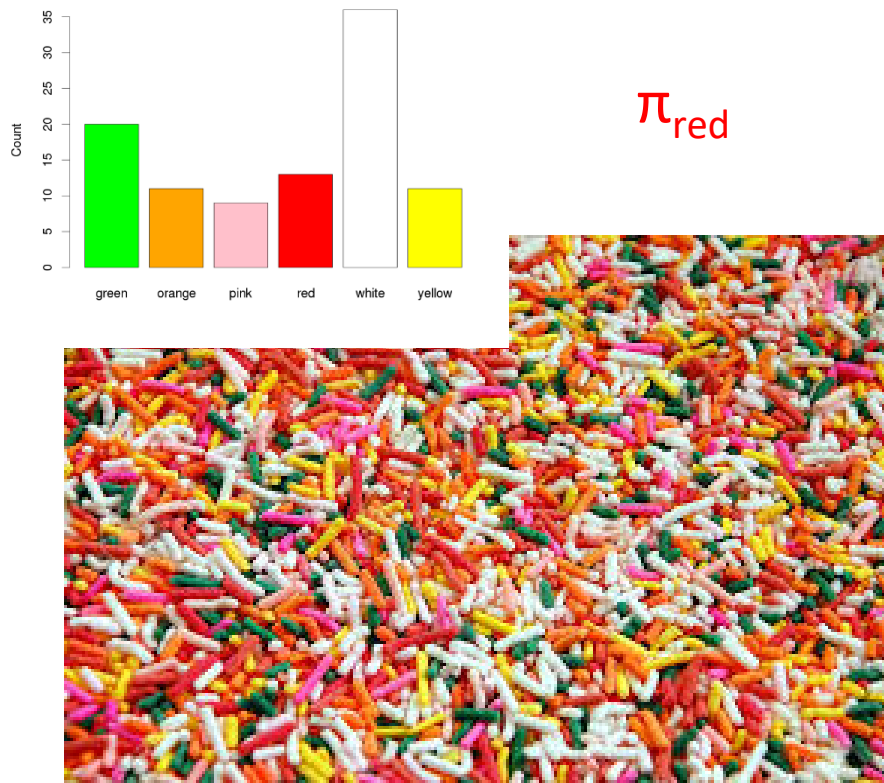
A: The mean length of 10 words (\bar{x})

i.e., each point in our **distribution** is a statistic!

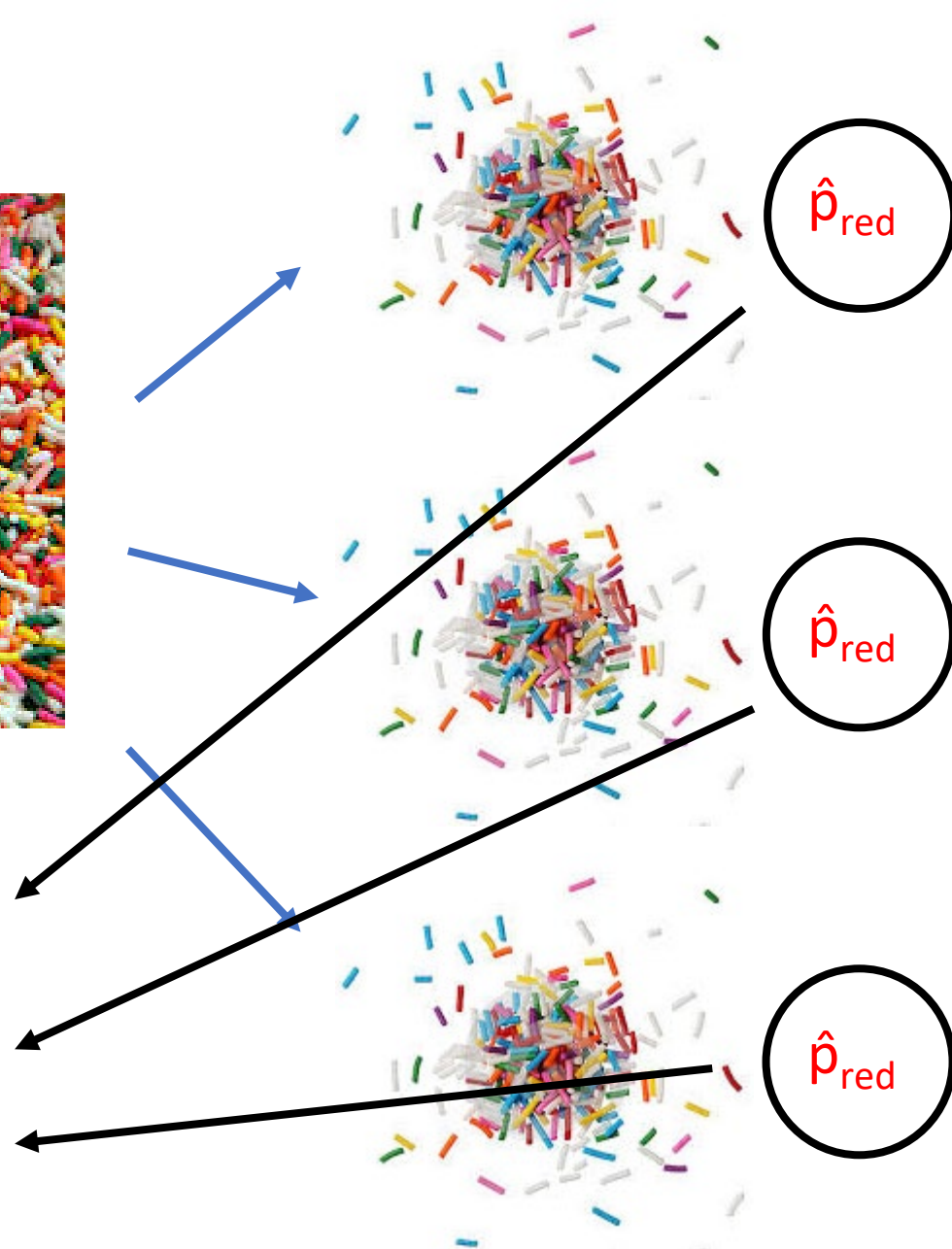
Sampling distribution

A **sampling distribution** is the distribution of sample statistics computed for different samples of the same size (n) from the same population.

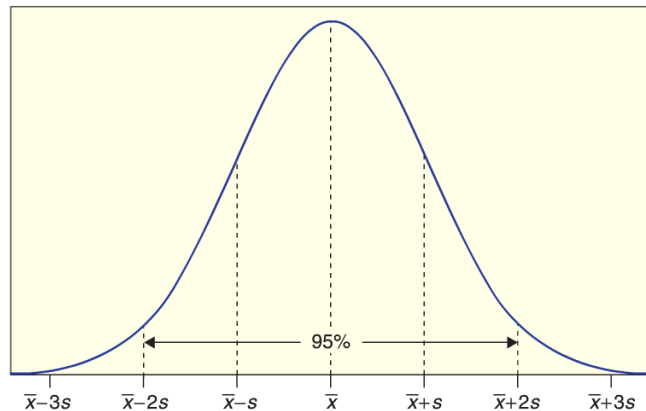
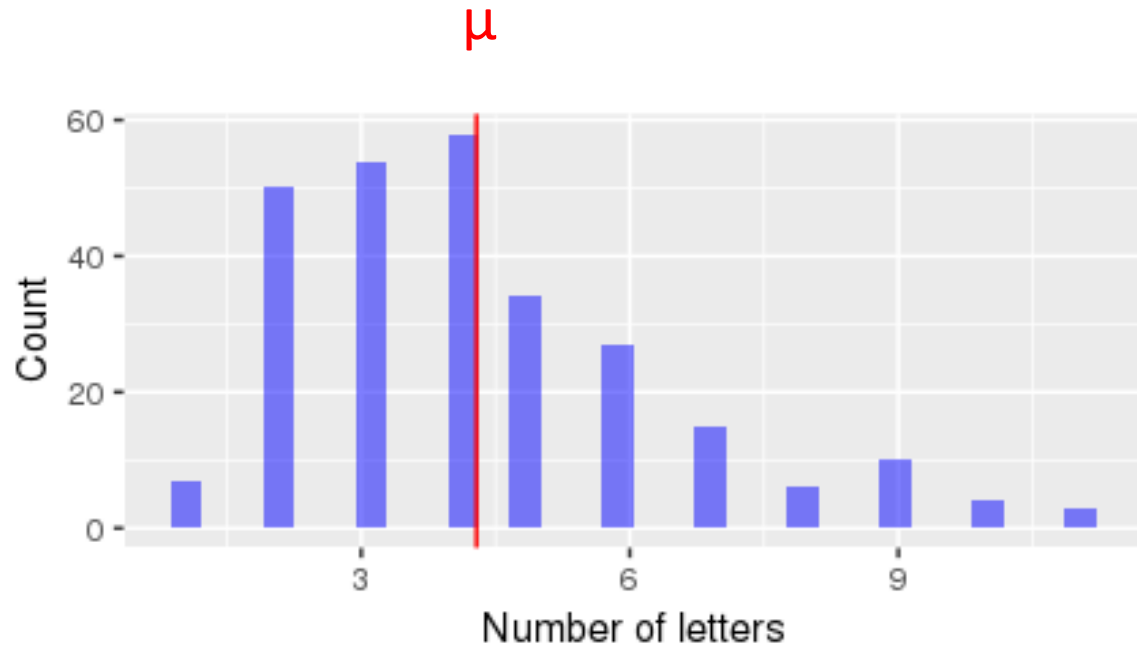
A sampling distribution shows us how the sample statistic varies from sample to sample.



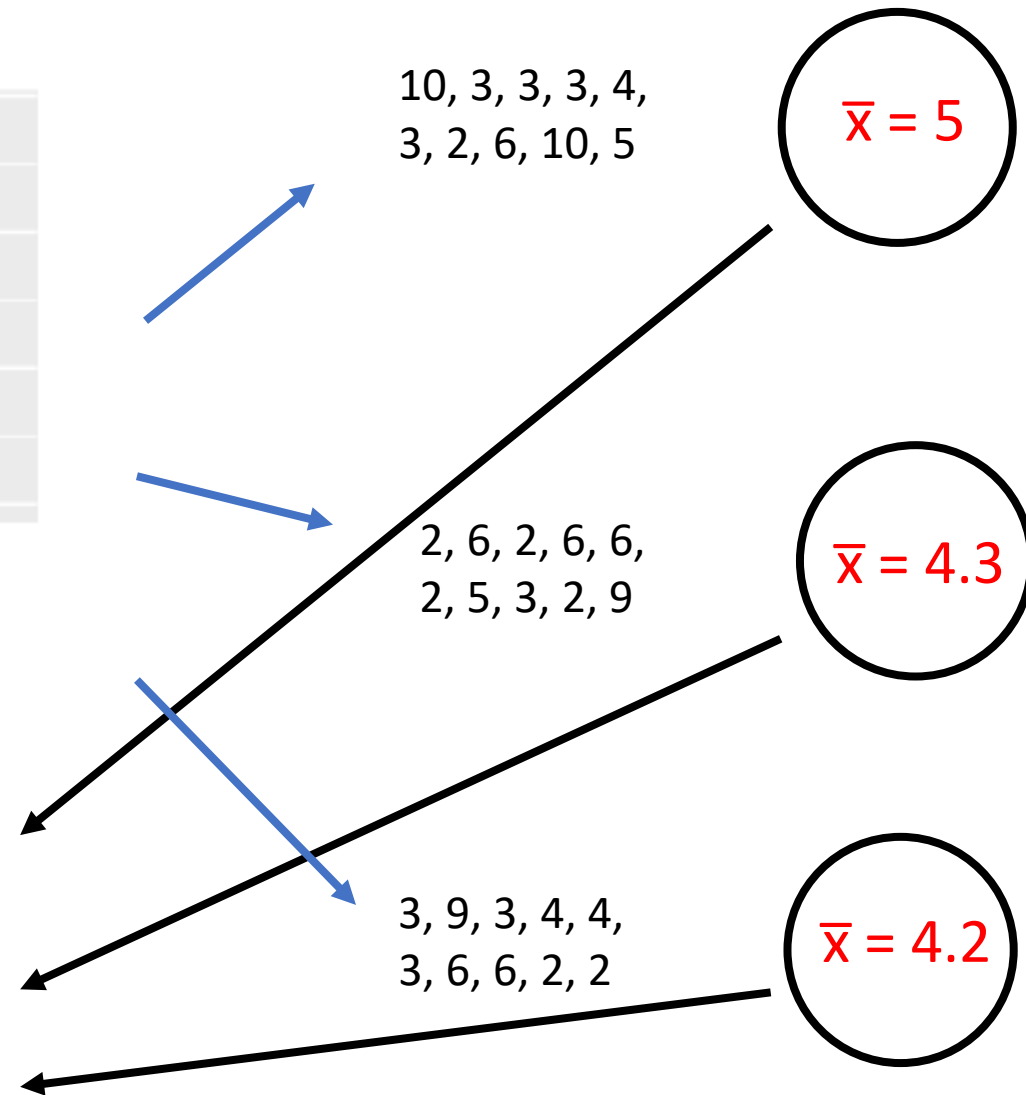
Sampling distribution!



Gettysburg address word length sampling distribution



Sampling distribution!



[Gettysburg sampling distribution app](#)