# Review of confidence intervals and introduction to the bootstrap

# Overview

Review and continuation of sampling distributions and confidence intervals

The bootstrap

If there is time:

- Using the bootstrap to create confidence intervals in R

# Announcement

Homework 4 has been posted!

It is due on Gradescope on <span style="color:red">Sunday February 15th at 11pm</span>

- **Be sure to mark each question on Gradescope!**

The material this week is going to be a bit more conceptually challenging

**Please attend the practice sessions** and office hours to reinforce your understanding!

# Review of confidence intervals

# Question: What is a confidence interval?

A **confidence interval** is an interval <u>computed by a method</u> that will contain the *parameter* a specified percent of times

- i.e., if the estimation were repeated many times, the interval will have the parameter x% of the time

The **confidence level** is the percent of all intervals that contain the parameter
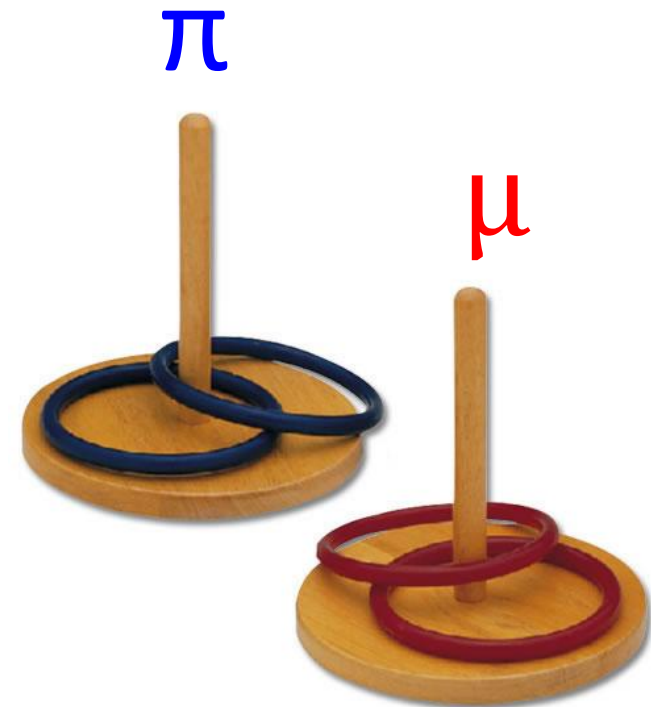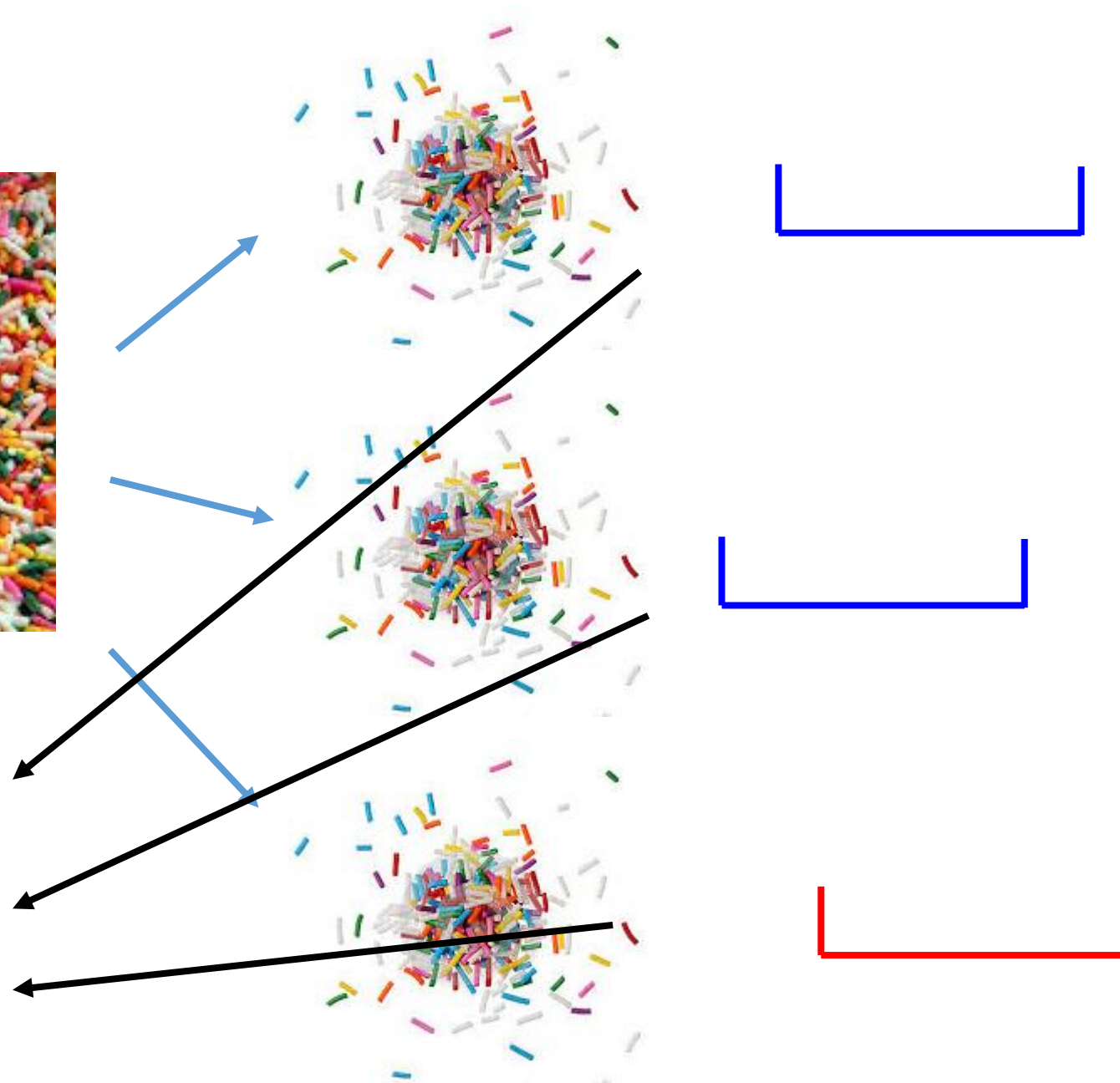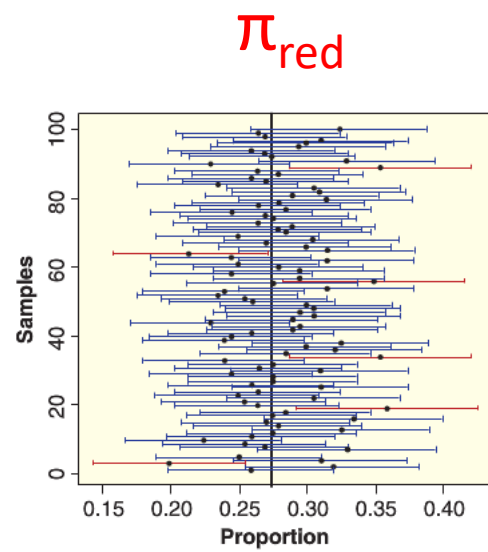
# Think ring toss…
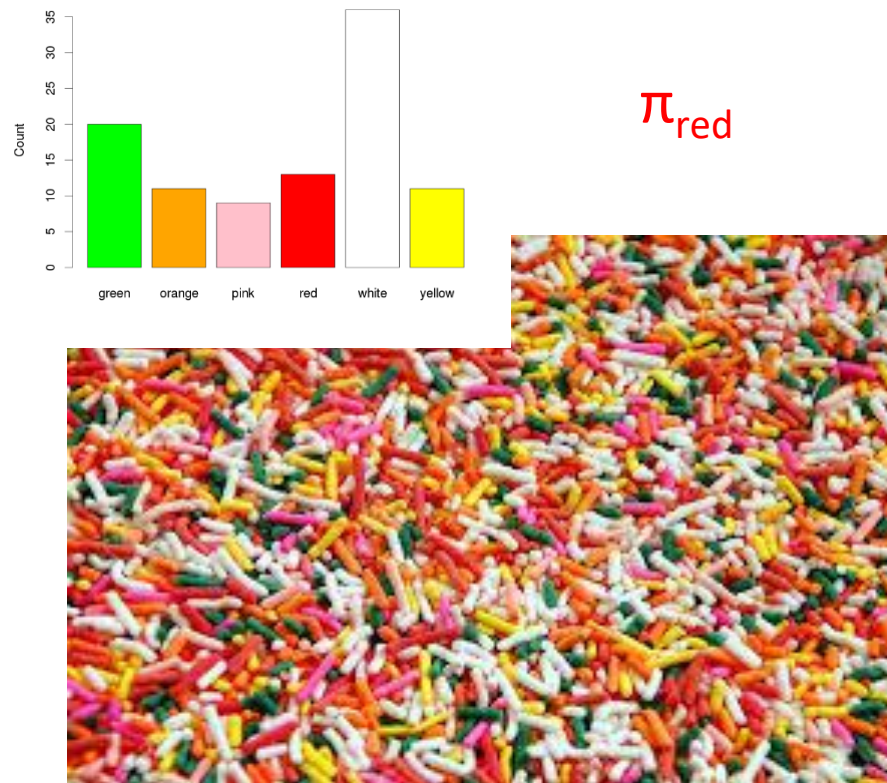
**A parameter exists in the world**

**We toss intervals at it**
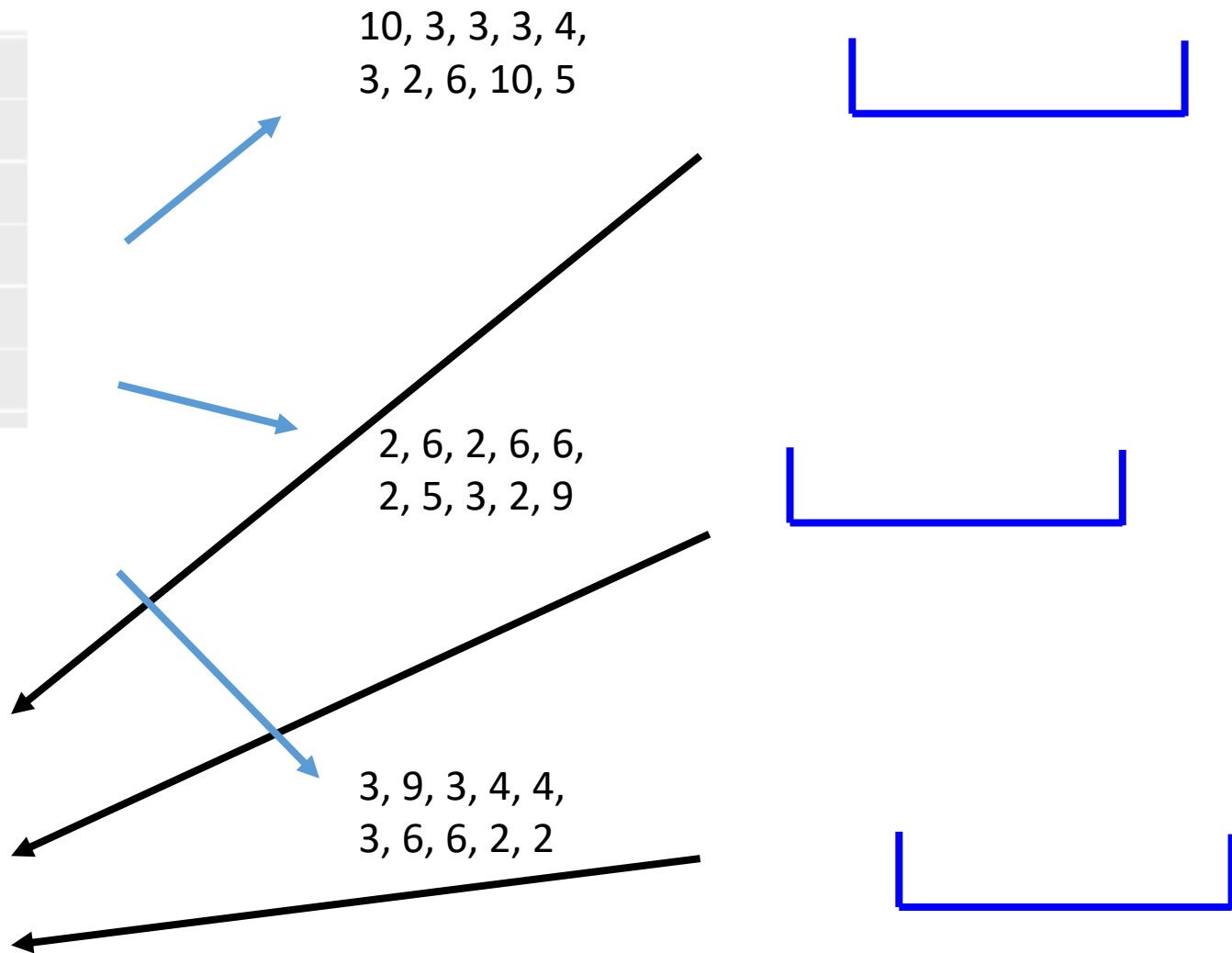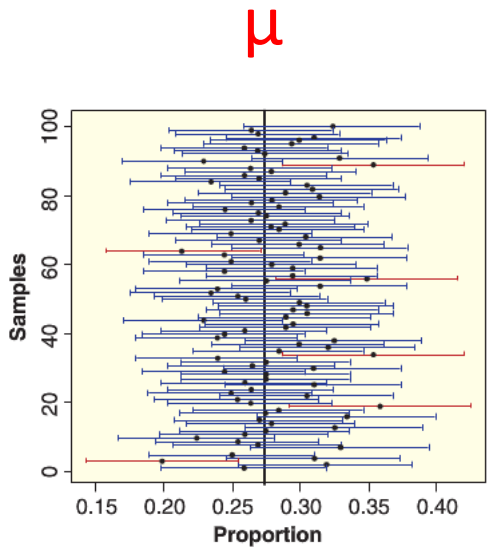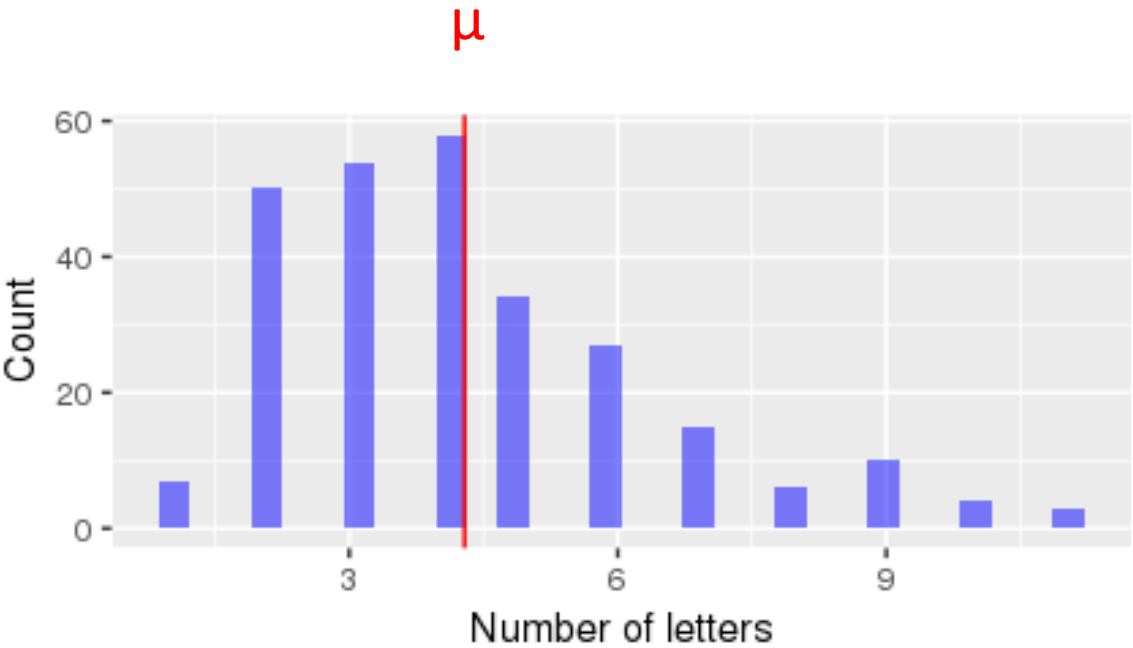- These are our confidence intervals

**95% of those intervals capture the parameter**
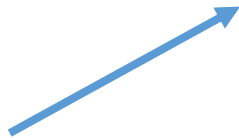- (95% is our confidence level)

$\pi$

$\mu$

$\pi_{red}$

$\pi_{red}$
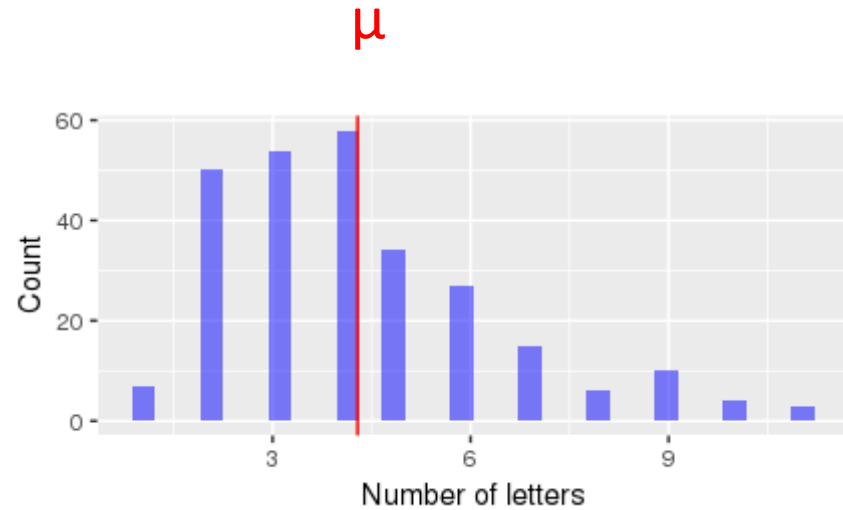
# Gettysburg address word length sampling distribution



10, 3, 3, 3, 4,
3, 2, 6, 10, 5

2, 6, 2, 6, 6,
2, 5, 3, 2, 9

3, 9, 3, 4, 4,
3, 6, 6, 2, 2

$\pi_{red}$
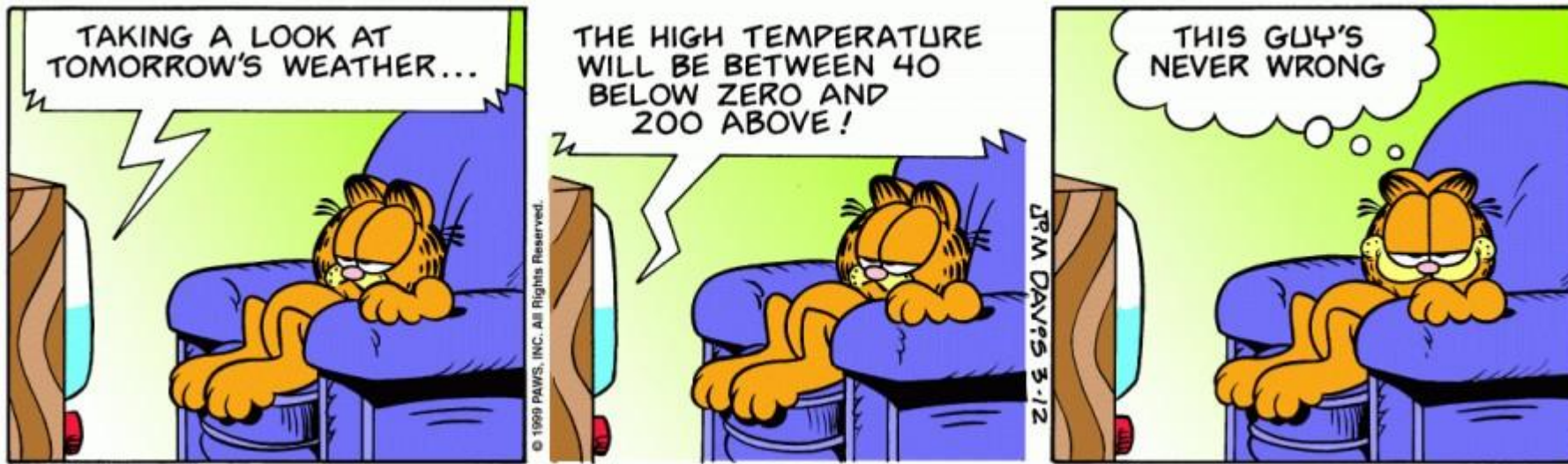
$\mu$

10, 3, 3, 3, 4,
3, 2, 6, 10, 5

$\pi$

$\mu$

For a 95% confidence level, 95% of the intervals we create will have the parameter in them

# Confidence Intervals

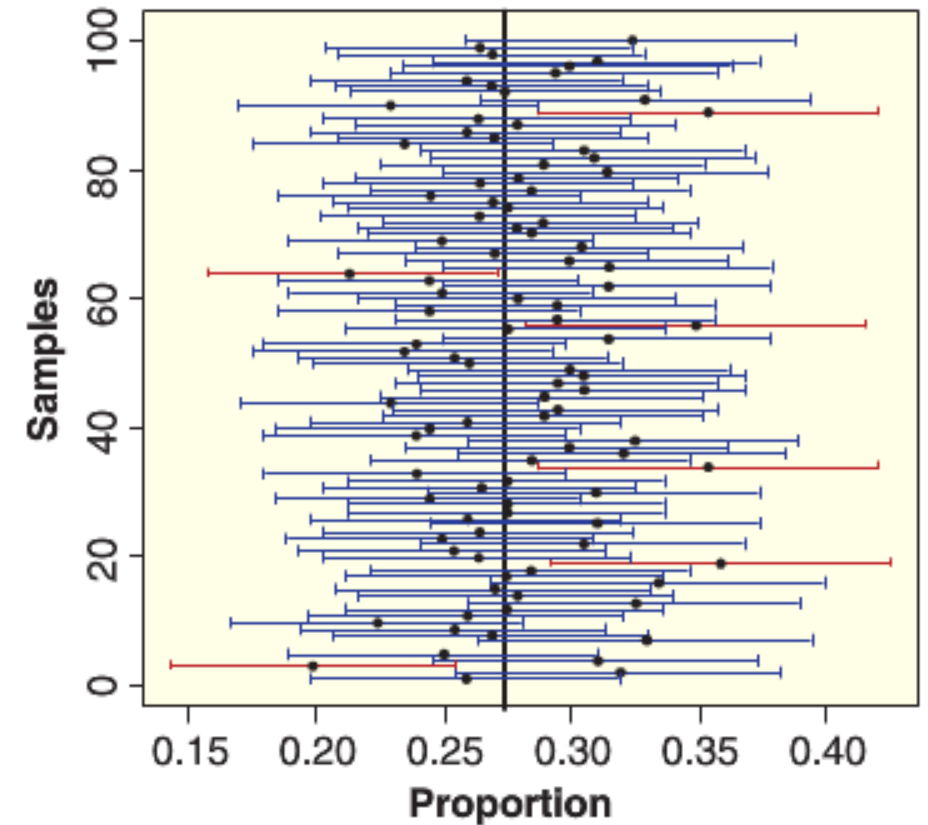There is a tradeoff between the **confidence level** (percent of times we capture the parameter) and the **confidence interval size!**

# Confidence Intervals

Q: For any given confidence interval, do we know whether it has really captured the parameter?

But we do know that if we create 100 intervals, ~95 of these intervals will have the parameter in it
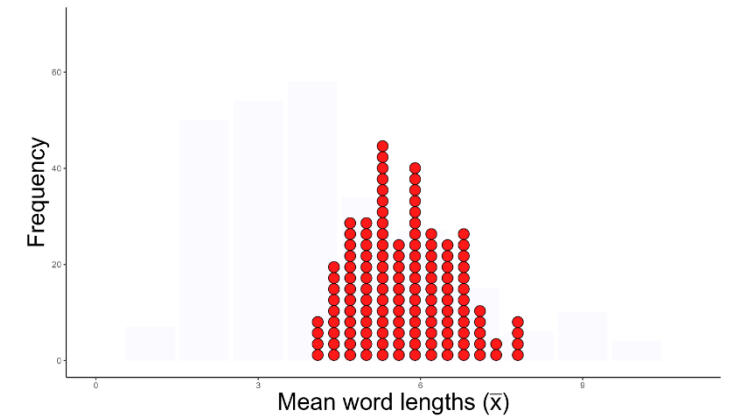
(for a 95% confidence interval)

# Review of sampling distributions

# Sampling distributions

Q$_2$: What is a sampling distribution?


Q$_3$: What does a sampling distribution show us?

# Gettysburg address word length sampling distribution



μ

Count
40 -
20 -
0 -

3      6      9

Number of letters

E[x̄] = μ

Sampling distribution!

x̄−3s   x̄−2s   x̄−s   x̄   x̄+s   x̄+2s   x̄+3s

95%

10, 3, 3, 3, 4,
3, 2, 6, 10, 5

x̄ = 5

2, 6, 2, 6, 6,
2, 5, 3, 2, 9

x̄ = 4.3

3, 9, 3, 4, 4,
3, 6, 6, 2, 2

x̄ = 4.2

Unbiased: E[x̄] = μ

# The standard error

Q$_4$: What is the **standard error**?

Q$_5$: What symbol do we use to denote the standard error?

# Sampling distribution in R

$Q_6$: Suppose we have a function called get_sample() that can generate samples from a population

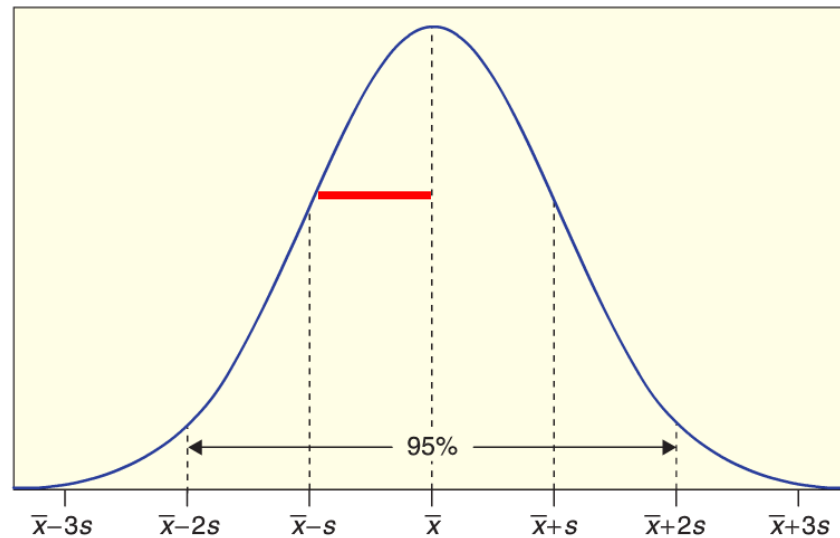How could we estimate the SE of the mean using R?

```
sampling_dist <- do_it (10000) *  {

        curr_sample <- get_sample()
        mean(curr_sample)       ⟵  What symbol should
                                    we use for this quantity?
}

SE_mean <- sd(sampling_dist)
```
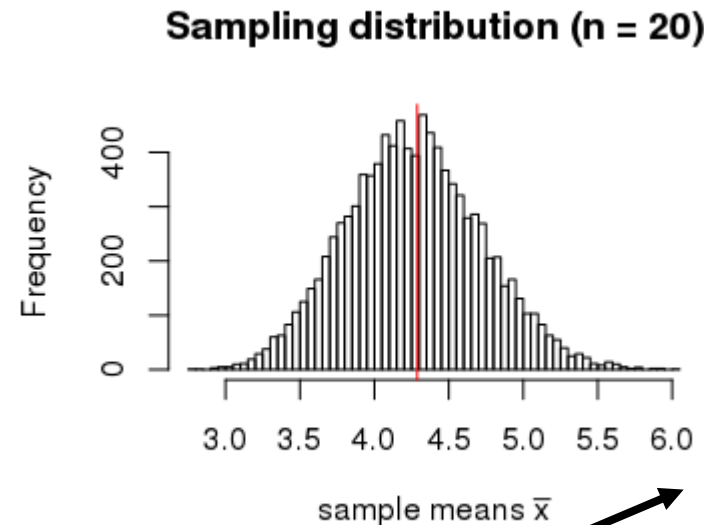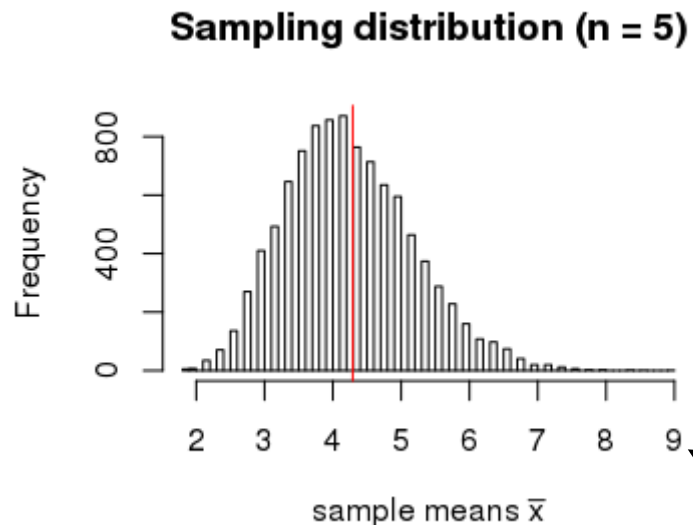
# Q$_7$: What are two ways that sampling distribution changes with a larger sample size n?

As the sample size n increases:

1. The sampling distribution becomes more like a normal distribution

2. The sampling distribution statistics become more concentrated around population parameter
   - i.e., the SE becomes smaller
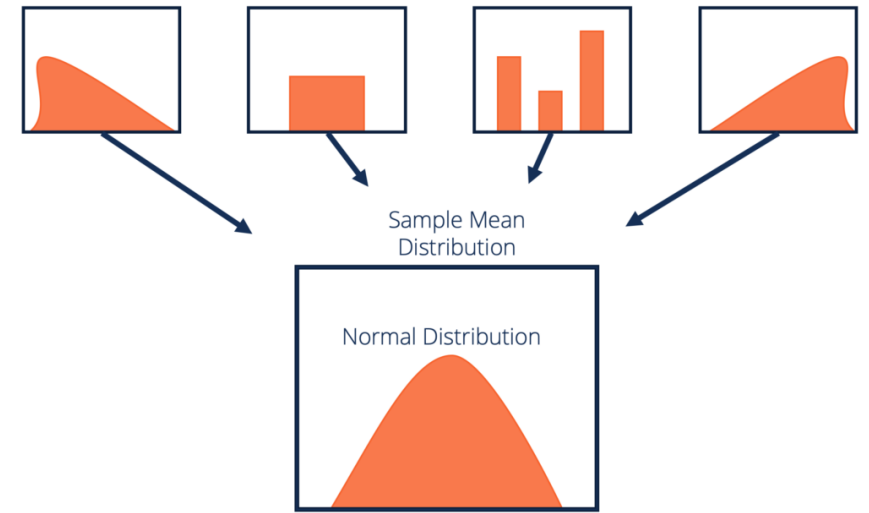


x-axis range 9 vs. 6

# Central limit theorem

For random samples with a sufficiently large sample size (n), the distribution of sample statistics for a <span style="color:red">mean ($\bar{x}$)</span> or a <span style="color:red">proportion ($\hat{p}$)</span> is:

- Normally distributed

- Centered at the value of the population parameter

# Normal distributions

Q$_9$: For a normal distribution, what percentage of points lie within 2 standard deviations for the population mean?

# Sampling distributions

Q$_{10}$: For a **sampling distribution** that is a normal distribution, what percentage of *statistics* lie within 2 standard deviations (SE) of the parameter value?



Q$_{11}$: If we had a *statistic value* and the value of the *SE*, could we compute a 95% confidence interval?

# Confidence intervals

Q$_{11}$: Suppose we are going to randomly chosen statistic value

And we are going create an interval centered at the statistic value

How large would the interval need to be to overlap with the parameter 95% of the time?



Confidence interval

# Confidence intervals



Parameter

2 SE

2 SE

0.20　0.25　0.30　0.35

$Q_{12}$: What is a formula can we use to calculate 95% confidence intervals?

95% confidence interval:  stat ± 2 · SE ⟶ $Q_{13}$: What is this quantity called?

# Confidence intervals

Parameter

$Q_{14}$: How frequently do 95% confidence intervals **fail** to capture the parameter of interest?

- 5% of the time

2 SE   2 SE

0.20   0.25   0.30   0.35

stat

2 · SE   2 · SE

Confidence interval

95% confidence interval:  stat  ±  2  · SE

# Confidence intervals for other confidence levels

Q$_{15}$: How could we get a 99.7% confidence interval confidence level?



99.7%

-3  -2  -1      +1  +2  +3

stat

3 · SE          3 · SE

Confidence interval

# Confidence intervals for other confidence levels

Q$_{16}$: How could we get a 68% confidence interval confidence level?



Confidence interval

# Confidence intervals for other confidence levels

Q[16]: How could we get a confidence interval for the q[th] confidence level?



CI = stat ± q* · SE

library(SDS1000)
cnorm(0.95)
 [1]  1.96

# Sampling distributions

Q$_{18}$: Could we calculate the SE by repeatedly sampling from a population to create sampling distribution, and then take the sd of this sampling distribution?





SCIENCE!!!

# Sampling distributions

Q$_{19}$: If we can't calculate the sampling distribution, what else can we do?

Bootstraps

1. Estimate SE with $\hat{SE}$ *from a single sample of data*
2. Then use $\bar{x} \pm 2 \cdot \hat{SE}$ to get the 95% CI

# Confident intervals

**If you don't feel confident about confidence intervals...**

Go to the practice sessions

Review material

Ask questions on Ed Discussion

Come to office hours!

# The bootstrap

# The bootstrap

The bootstrap is a method to estimate the standard error

- $\hat{SE}$ is an estimate for SE
- We will use the symbol SE* as the **bootstrap estimate** for SE  (rather than $\hat{SE}$ )

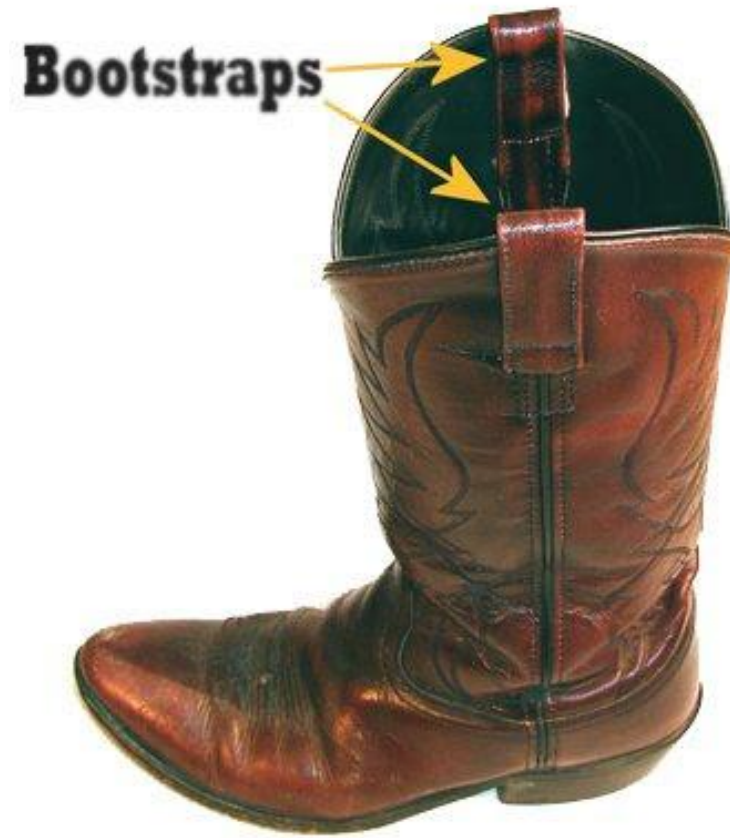1. Estimate SE with SE*

2. Then use  $\bar{x}$ ± 2 · SE*  to get the 95% CI

# Plug-in principle

Suppose we get one sample of size $n$ from a population

We <u>pretend that this sample is the population</u>   (plug-in principle)

1. We then sample $n$ points <u>with replacement</u> from ***our sample***, and compute our statistic of interest

2. We repeat this process 1000's of times and get a *bootstrap* sample distribution

3. The standard deviation of this bootstrap distribution (SE* bootstrap) is a good approximate for standard error SE from the real sampling distribution

# Gettysburg address word length bootstrap distribution

**The sample (n = 10)**
10, 3, 3, 3, 4, 3, 2, 6, 4, 5

μ

Count — Number of letters

3, 3, 3, 5, 3,
4, 5, 2, 2, 10

$\bar{x}* = 4$

3, 3, 2, 3, 6,
4, 6, 5, 3, 6

$\bar{x}* = 4.1$

5, 3, 2, 3, 3,
3, 10, 3, 4, 3

$\bar{x}* = 3.9$

SE*

95%

$\bar{x}-3s$  $\bar{x}-2s$  $\bar{x}-s$  $\bar{x}$  $\bar{x}+s$  $\bar{x}+2s$  $\bar{x}+3s$

Bootstrap distribution!

Notice there is no 9's in the bootstrap samples

# Bootstrap process



"Fake sampling distribution"

# 95% Confidence Intervals

When a bootstrap distribution for a sample statistic is approximately normal, we can estimate a 95% confidence interval using:

$$Statistic \ \pm \ 2 \cdot SE^*$$

Where SE* is the standard error estimated using the bootstrap

# Findings CIs for many different parameters

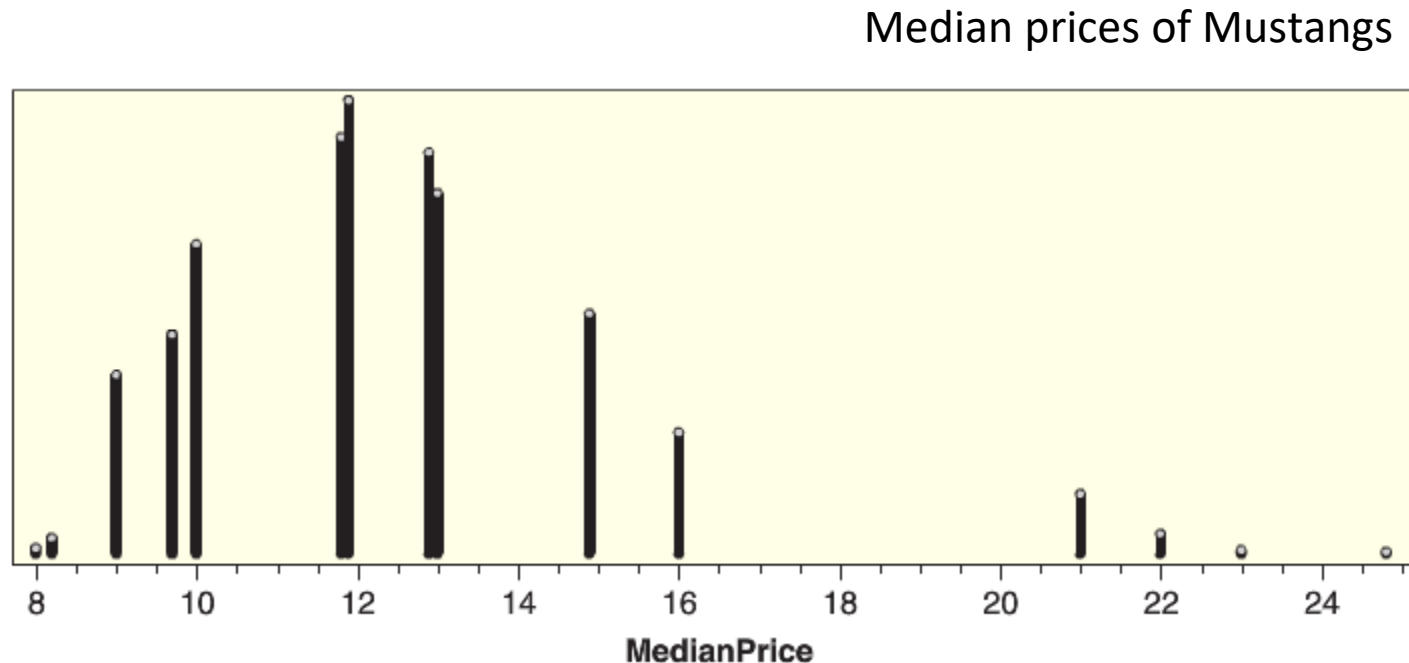The bootstrap method works for constructing confidence intervals for many different types of parameters!

# Caution: the bootstrap does not always work

Always look at the bootstrap distribution, if it is poorly behaved (e.g., heavily skewed, has isolated clumps of values, etc.), you should not trust the intervals it produces.

Median prices of Mustangs



MedianPrice

# Calculating bootstrap confidence intervals in R

# What are the steps needed to create a bootstrap SE?

1. Start with a sample

2. Repeat steps 10,000 times

   a. Resample the points in the sample to get a bootstrap sample
   b. Compute the statistic of interest on the bootstrap sample

3. Take the standard deviation of the bootstrap distribution to get SE*

# Sampling with replacement from a vector

my_sample <- c(3, 1, 4, 1, 5, 9)

To get a sample of size n = 6 with replacement:

boot_sample  <-  sample(my_sample,  6,  replace = TRUE)

# Sampling distribution in R

```
my_sample <- c(21, 29, 25, 19, 24, 22, 25, 26, 25, 29)


bootstrap_dist <-  do_it(10000) * {

        curr_boot <- sample(my_sample , 10, replace = TRUE)
        mean(curr_boot)

}


SE_boot <- sd(bootstrap_dist)
```

# Bootstrap confidence interval in R

obs_mean <- mean(my_sample)

CI_lower <-  obs_mean  - 2 * SE_boot

CI_upper <-  obs_mean  + 2 * SE_boot