

# Statistics Review II

EC 320: Introduction to Econometrics

---

Emmett Saulnier

Winter 2022

# Housekeeping

- Computational problem set 1 was due yesterday
- Analytical problem set 1 is due on Friday at midnight
- I will post the next computational problem set soon
- "Where can I get help with R?"
  - Labs
  - RStudio has a nice set of resources here
  - Google
  - Office hours
  - Email me or Gio

In particular, [this video](#) about RMarkdown might help you start to figure out what is going on in Rmd files.

# Statistics Review

# Overview

A few key terms:

- **Population:** a (usually large) group of items or events we would like to know about.
- **Parameter:** a value that describes that population. The parameter of interest is the parameter that the researcher seeks to learn about.
- **Sample:** a survey of a subset of the population.

Usually we aim to draw observations **randomly** from the population, such that it becomes a **representative sample** of the population. More details on this later!

# Overview

## **Focus:** Populations vs Samples

- How can we make inferences about a **population** based on a small **sample** of the population?
- In particular, how do we learn about an unknown population *parameter* of interest?

**Challenge:** Usually cannot access information about the entire population.

**Solution:** Sample from the population and estimate the parameter.

- Draw  $n$  observations from the population, then use an estimator.

# Sampling

There are myriad ways to produce a sample,<sup>\*</sup> but we will restrict our attention to **simple random sampling**, where

1. Each observation is a random variable.
2. The  $n$  random variables are independent.
3. Life becomes much simpler for the econometrician.

<sup>\*</sup> Only a subset of these can help produce reliable statistics.

# Estimators

An **estimator** is a rule (or formula) for estimating an unknown population parameter given a sample of data.

- Each observation in the sample is a random variable.
- An estimator is a combination of random variables  $\implies$  it is a random variable.

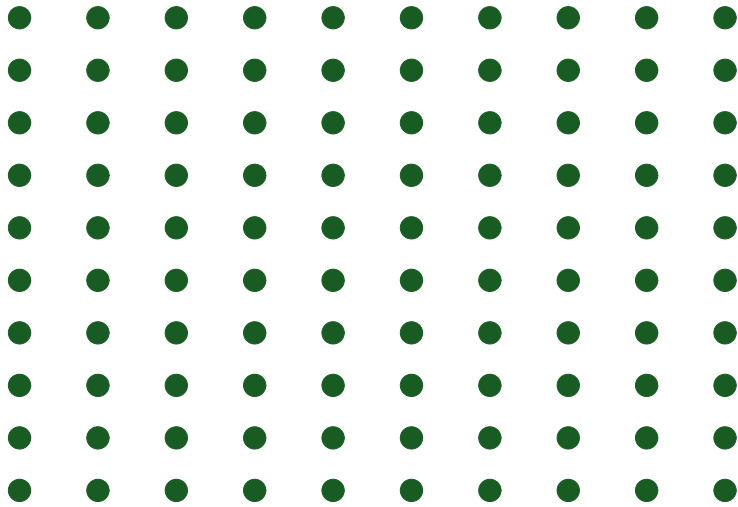
**Example:** Sample mean

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

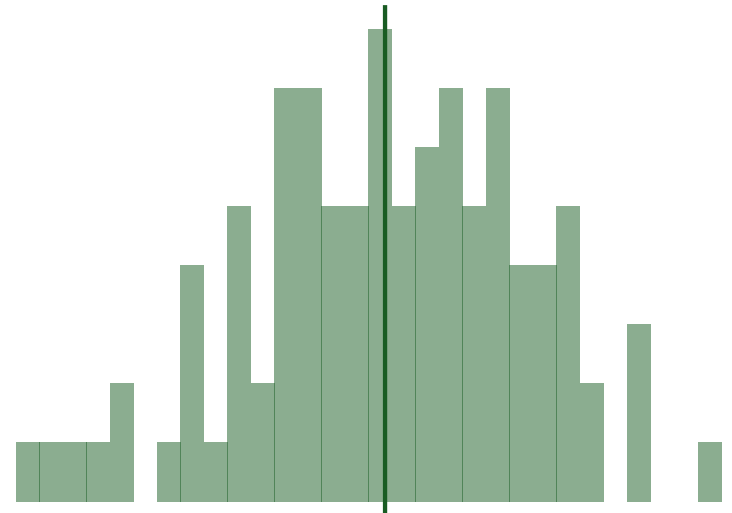
- $\bar{X}$  is an estimator for the population mean  $\mu$ .
- Given a sample,  $\bar{X}$  yields an **estimate**  $\bar{x}$  or  $\hat{\mu}$ , a specific number.

# Population vs. Sample

**Question:** Why do we care about *population vs. sample*?



**Population**



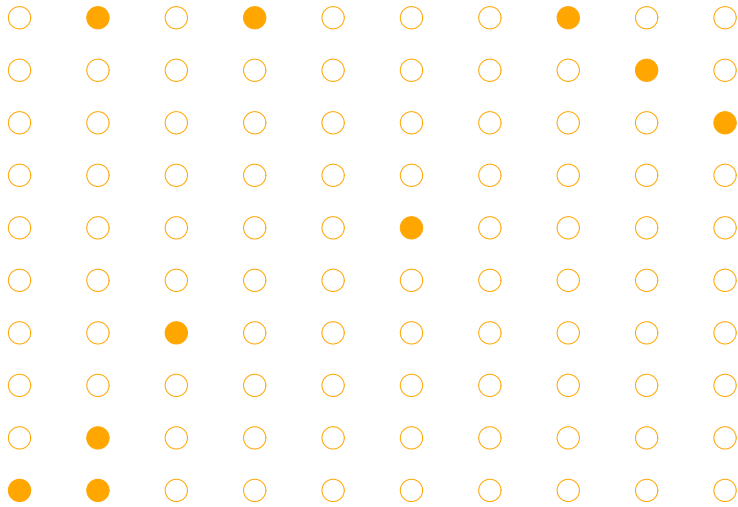
**Population relationship**

$$\mu = 3.75$$

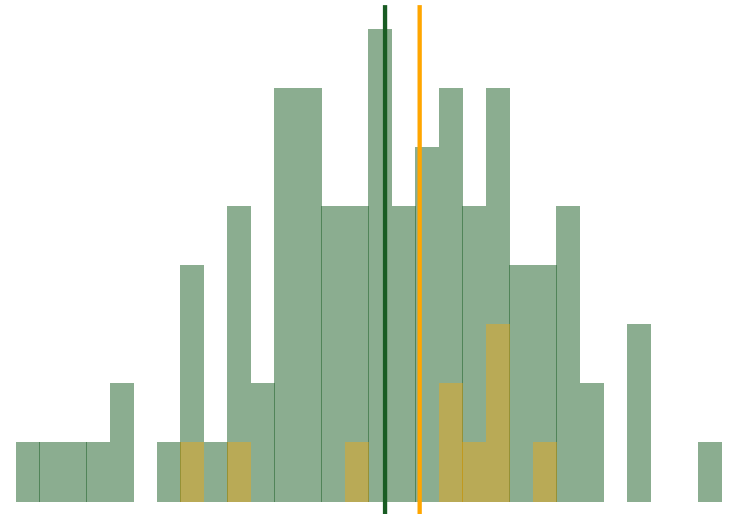


# Population vs. Sample

**Question:** Why do we care about *population vs. sample*?



**Sample 1:** 10 random individuals



**Population relationship**

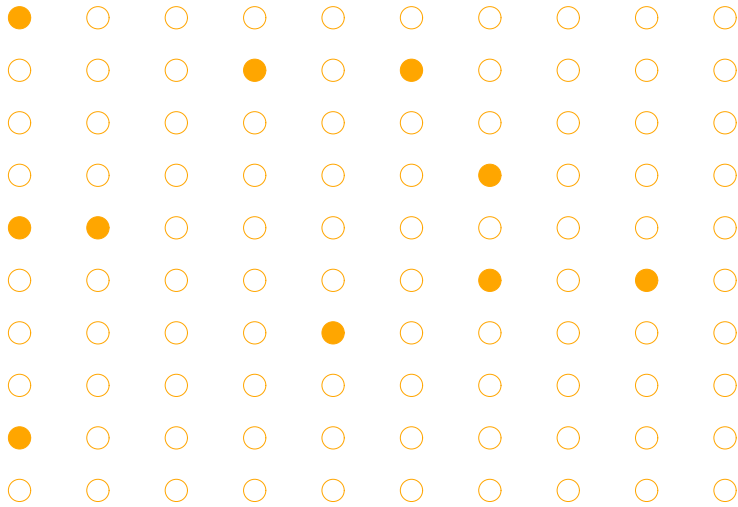
$$\mu = 3.75$$

**Sample relationship**

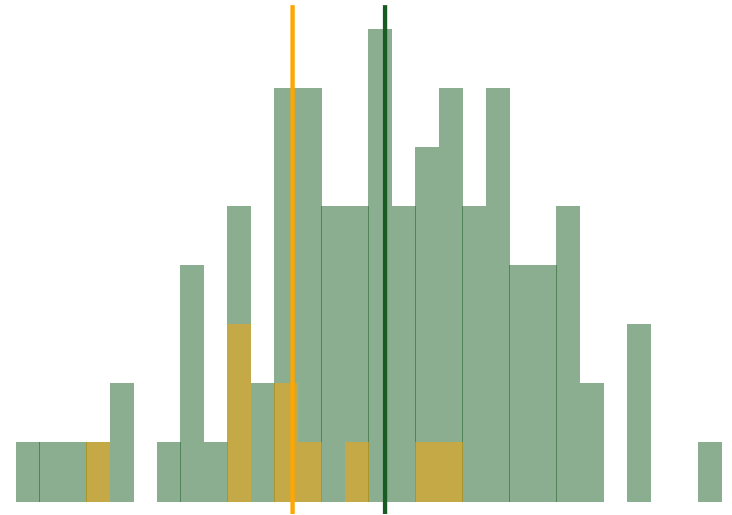
$$\hat{\mu} = 8.34$$

# Population vs. Sample

**Question:** Why do we care about *population vs. sample*?



**Sample 2:** 10 random individuals



**Population relationship**

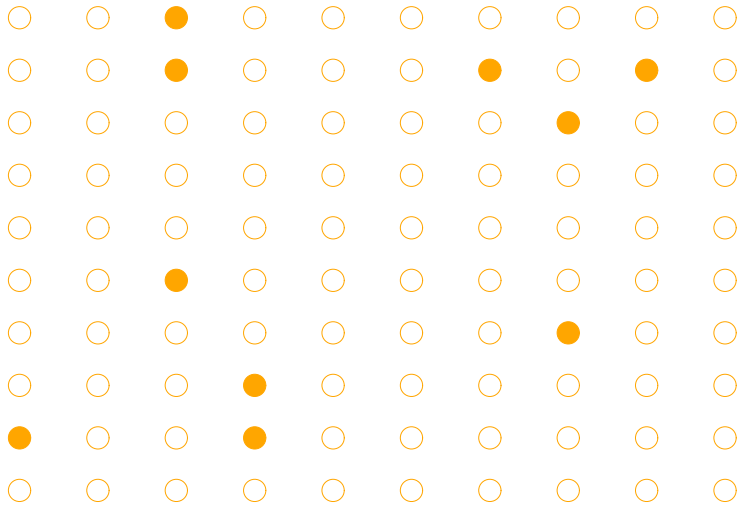
$$\mu = 3.75$$

**Sample relationship**

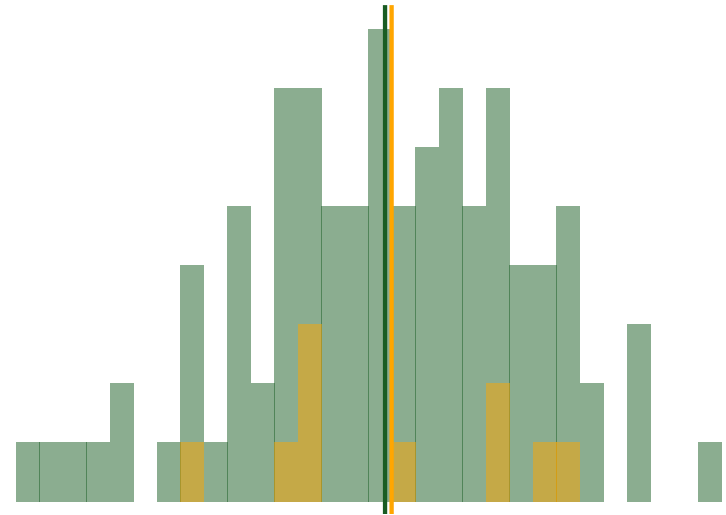
$$\hat{\mu} = -8.54$$

# Population vs. Sample

**Question:** Why do we care about *population vs. sample*?



**Sample 3:** 10 random individuals



**Population relationship**

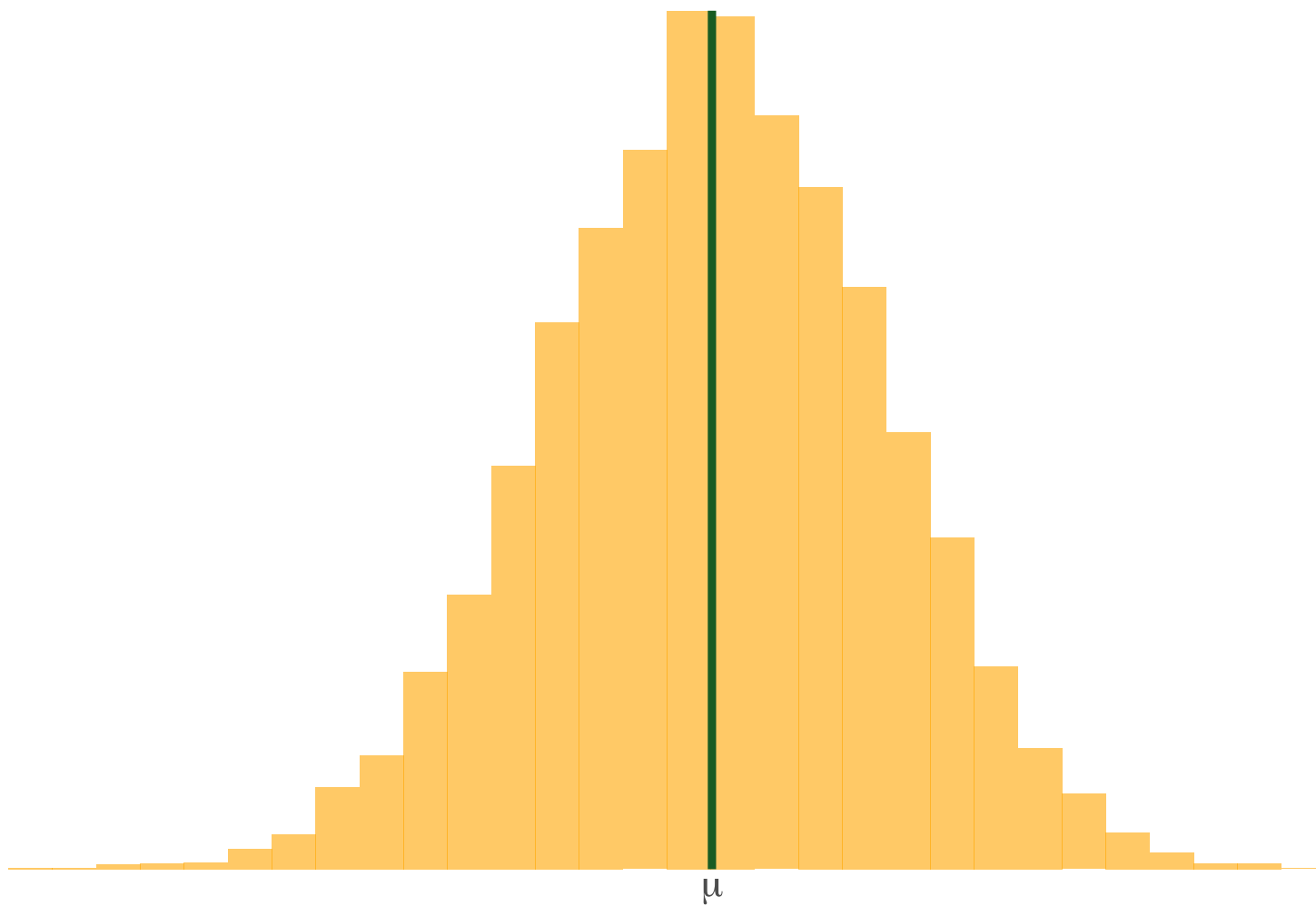
$$\mu = 3.75$$

**Sample relationship**

$$\hat{\mu} = 4.62$$

Let's repeat this **10,000 times** and then plot the estimates.

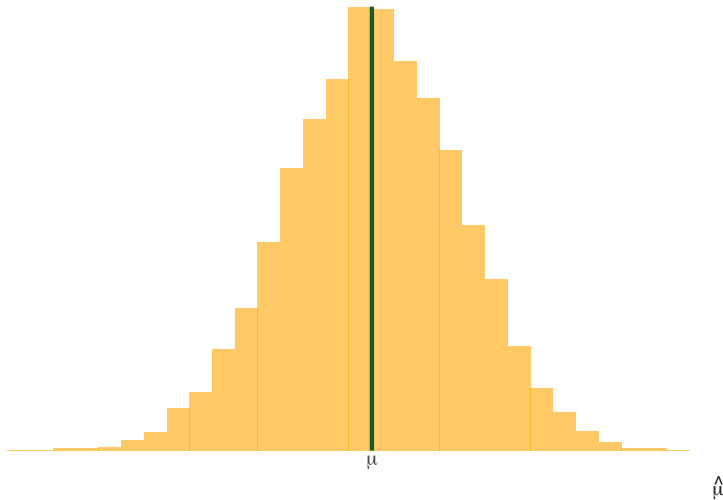
(This exercise is called a Monte Carlo simulation.)



$\hat{\mu}$

# Population vs. Sample

**Question:** Why do we care about *population vs. sample*?



- On average, the mean of the samples are close to the population mean.
- But...some individual samples can miss the mark.
- The difference between individual samples and the population creates **uncertainty**.

# Population vs. Sample

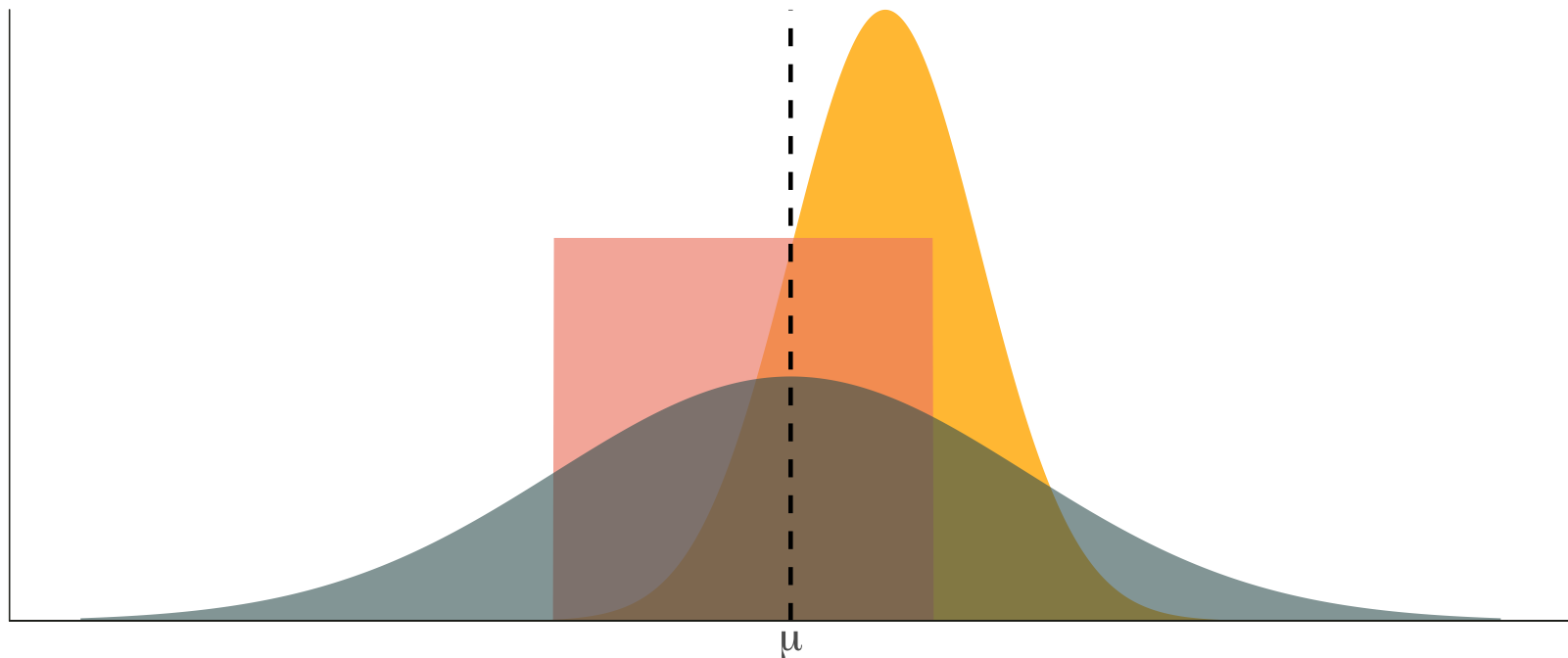
**Question:** Why do we care about *population vs. sample*?

**Answer:** Uncertainty matters.

- $\hat{\mu}$  is a random variable that depends on the sample.
- In practice, we don't know whether our sample is similar to the population or not.
- Individual samples may have means that differ greatly from the population.
- We will have to keep track of this uncertainty.

# Properties of Estimators

Imagine that we want to estimate an unknown parameter  $\mu$ , and we know the distributions of three competing estimators. **Which one should we use?**





# Properties of Estimators

**Question:** What properties make an estimator reliable?

**Answer 1: Unbiasedness.**

On average (after *many* samples), does the estimator tend toward the correct value?

**More formally:** Does the mean of estimator's distribution equal the parameter it estimates?

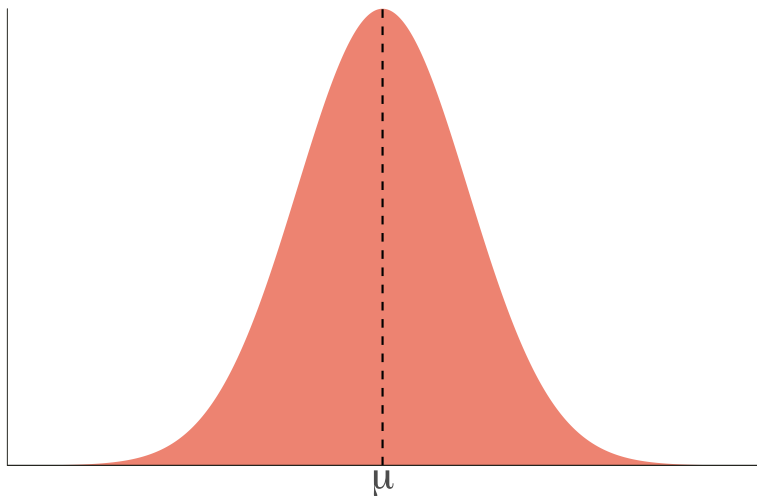
$$\text{Bias}_{\mu}(\hat{\mu}) = \mathbb{E}[\hat{\mu}] - \mu$$

# Properties of Estimators

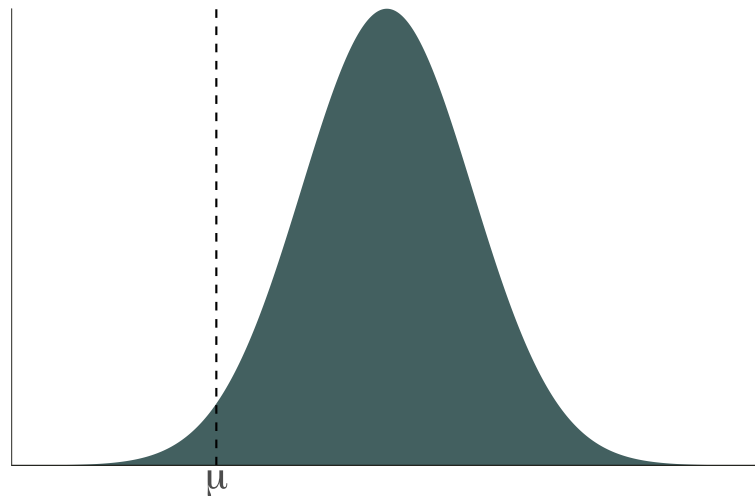
**Question:** What properties make an estimator reliable?

**Answer 1: Unbiasedness.**

**Unbiased estimator:**  $\mathbb{E}[\hat{\mu}] = \mu$



**Biased estimator:**  $\mathbb{E}[\hat{\mu}] \neq \mu$



# Unbiasedness Example

Is the sample mean  $\frac{1}{n} \sum_{i=1}^n x_i = \hat{\mu}$  an unbiased estimator of the population mean  $E(x_i) = \mu$ ?

$$\begin{aligned}\mathbb{E}[\hat{\mu}] &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n x_i\right] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[x_i] \\ &= \frac{1}{n} \sum_{i=1}^n \mu \\ &= \mu\end{aligned}$$

# Properties of Estimators

**Question:** What properties make an estimator reliable?

**Answer 2: Efficiency (low variance).**

The central tendencies (means) of competing distributions are not the only things that matter. We also care about the **variance** of an estimator.

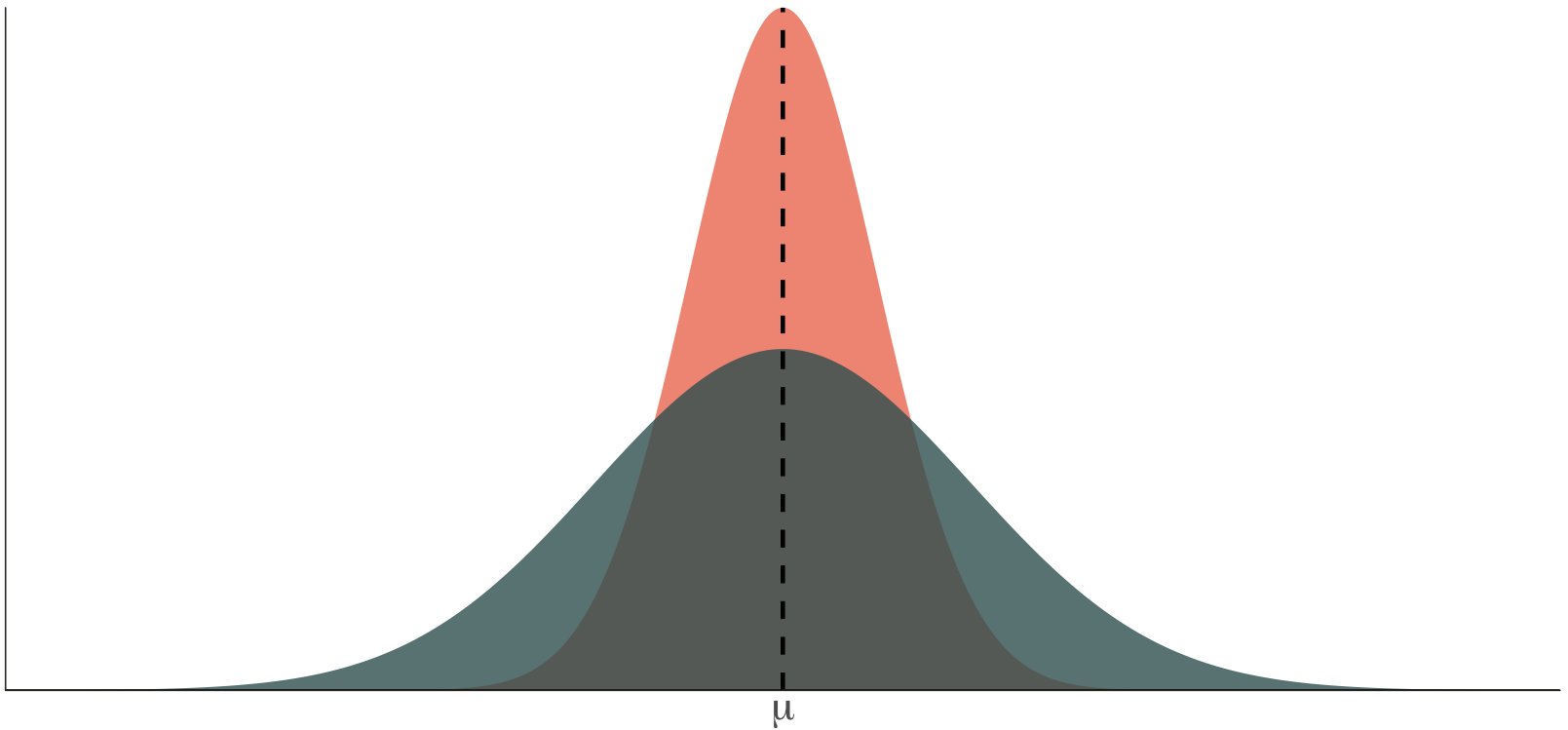
$$\text{Var}(\hat{\mu}) = \mathbb{E}\left[(\hat{\mu} - \mathbb{E}[\hat{\mu}])^2\right]$$

Lower variance estimators produce estimates closer to the mean in each sample.

# Properties of Estimators

**Question:** What properties make an estimator reliable?

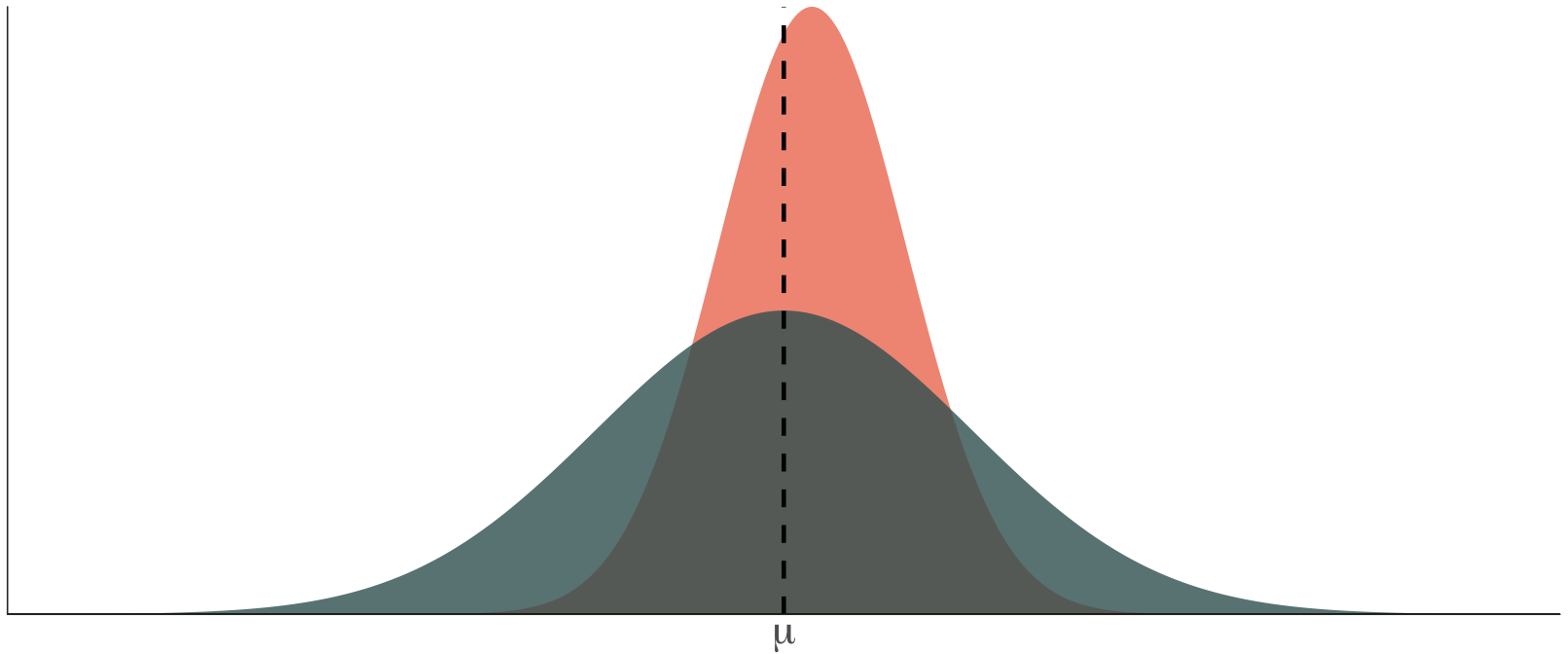
**Answer 2: Low Variance (a.k.a. Efficiency).**



# The Bias-Variance Tradeoff

Should we be willing to take a bit of bias to reduce the variance?

In econometrics, we generally prefer unbiased estimators. Some other disciplines think more about this tradeoff.



# Unbiased Estimators

In addition to the sample mean, there are several other unbiased estimators we will use often.

- **Sample variance** to estimate variance  $\sigma^2$ .
- **Sample covariance** to estimate covariance  $\sigma_{XY}$ .
- **Sample correlation** to estimate the population correlation coefficient  $\rho_{XY}$ .

# Unbiased Estimators

Sample variance  $S_X^2$  is an unbiased estimator of the population variance  $\sigma^2$

$$S_X^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

Sample covariance  $S_{XY}$  is an unbiased estimator of the population covariance  $\sigma_{XY}$

$$S_{XY} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}).$$

Sample correlation  $r_{XY}$  is an unbiased estimator of the population correlation coefficient  $\rho_{XY}$

$$r_{XY} = \frac{S_{XY}}{\sqrt{S_X^2} \sqrt{S_Y^2}}.$$



# Population Distributions

Suppose we have some estimator  $\hat{\theta}$  for a parameter  $\theta$ :

- We don't know  $\theta$ , but say we assume that  $\hat{\theta}$  follows some probability distribution  $p(\hat{\theta})$
- Next, suppose we hypothesize some value for  $\theta$ , say  $\theta = 2.5$
- Now we use our estimator  $\hat{\theta}$  to calculate an estimate for  $\theta$ . Say that we get  $\hat{\theta} = 45$
- We "know" the distribution of  $\hat{\theta}$ , so we know the probability of getting  $\hat{\theta} = 45$  if really  $\theta = 2.5$ . So we can say "if  $\theta$  really was 2.5, then the probability of getting  $\hat{\theta} = 45$  is super super low. Thus the probability that  $\theta = 2.5$  is super super low".
- We are able to make a statement about the true value of  $\theta$  just by knowing the distribution of our preferred estimator  $\hat{\theta}$

Sweet, but what distribution should we be assuming?

# The Central Limit Theorem

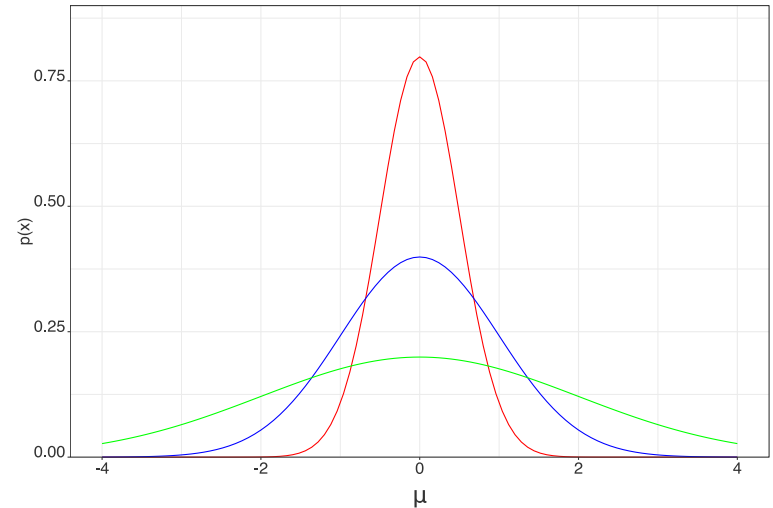
## Theorem

Let  $x_1, x_2, \dots, x_n$  be a random sample from a population with mean  $\mathbb{E}[X] = \mu$  and variance  $\text{Var}(X) = \sigma^2 < \infty$ , let  $\bar{X}$  be the sample mean. Then, as  $n \rightarrow \infty$ , the function  $\frac{\sqrt{n}(\bar{X} - \mu)}{S_x}$  converges to a **Normal Distribution** with mean 0 and variance 1.

- CLT states that when  $n \rightarrow \infty$ , the sample mean will be normally distributed.
- The Law of Large Number (LLN) states that as  $n \rightarrow \infty$ , the sample converges on the population mean.
- The only unknown parameter is  $\mu$ , and we can get around that by making probabilistic statements about it.

# Normal Distribution

- Continuous distribution where  $x_i$  can take the value of any real number
- Domain spans the entire real line
- Centered on the distribution mean  $\mu$
- The width of the distribution (fatness of its tails) is moderated  $\sigma^2$

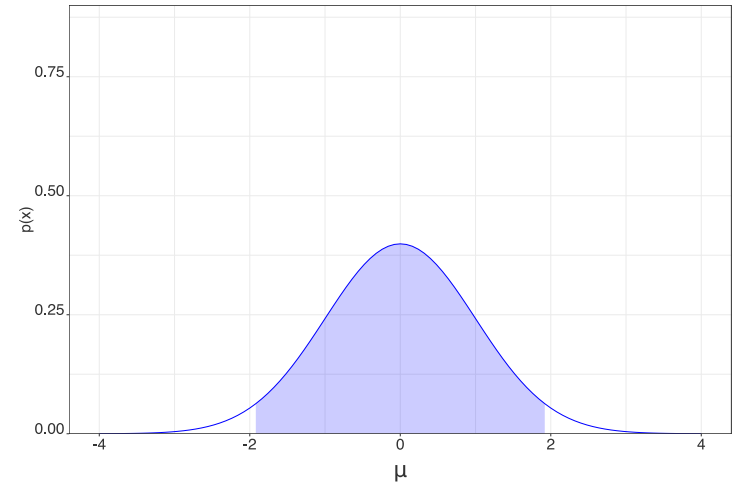


The greater the variance, the wider the range of values that commonly appear, hence greater probability density mass.

# Normal Distribution

**Rule 1:** The probability that the random variable takes a value  $x_i$  is 0 for any  $x_i \in \mathbb{R}$

**Rule 2:** The probability that the random variable takes a value within  $[x_i, x_j]$  range, where  $x_i \neq x_j$ , is the area under  $p(x)$  between those two values



The area above represents  $p(x) = 0.95$ . The values  $\{-1.96, 1.96\}$  represent the 95% confidence interval for  $\mu$ .

# Hypothesis Testing

How do we assess an estimate of the population mean?

- Is it meaningfully different than existing evidence on the population mean?
- Is it *statistically distinguishable* from previously hypothesized values of the population mean?
- Is the estimate extreme enough to update our prior beliefs about the population mean?

We can conduct statistical tests to address these questions.

# Hypothesis Testing

**Null hypothesis ( $H_0$ ):**  $\mu = \mu_0$

**Alternative hypothesis ( $H_1$ ):**  $\mu \neq \mu_0$

There are four possible outcomes of our test:

1. We **fail to reject** the null hypothesis and the null is true.
2. We **reject** the null hypothesis and the null is false.
3. We **reject** the null hypothesis, but the null is actually true (**Type I error**).
4. We **fail to reject** the null hypothesis, but the null is actually false (**Type II error**).

# Hypothesis Testing

We **fail to reject** the null hypothesis and the null is true.

- The defendant was acquitted and he didn't do the crime.

We **reject** the null hypothesis and the null is false.

- The defendant was convicted and he did the crime.

We **reject** the null hypothesis, but the null is actually true.

- The defendant was convicted, but he didn't do the crime!
- **Type I error** (a.k.a. *false positive*)

We **fail to reject** the null hypothesis, but the null is actually false.

- The defendant was acquitted, but he did the crime!
- **Type II error** (a.k.a. *false negative*)

# Hypothesis Testing

$\hat{\mu}$  is random: it could be anything, even if  $\mu = \mu_0$  is true.

- But if  $\mu = 0$  is true, then  $\hat{\mu}$  is unlikely to take values far from zero.
- As the variance of  $\hat{\mu}$  shrinks, we are even less likely to observe "extreme" values of  $\hat{\mu}$  (assuming  $\mu = \mu_0$ ).

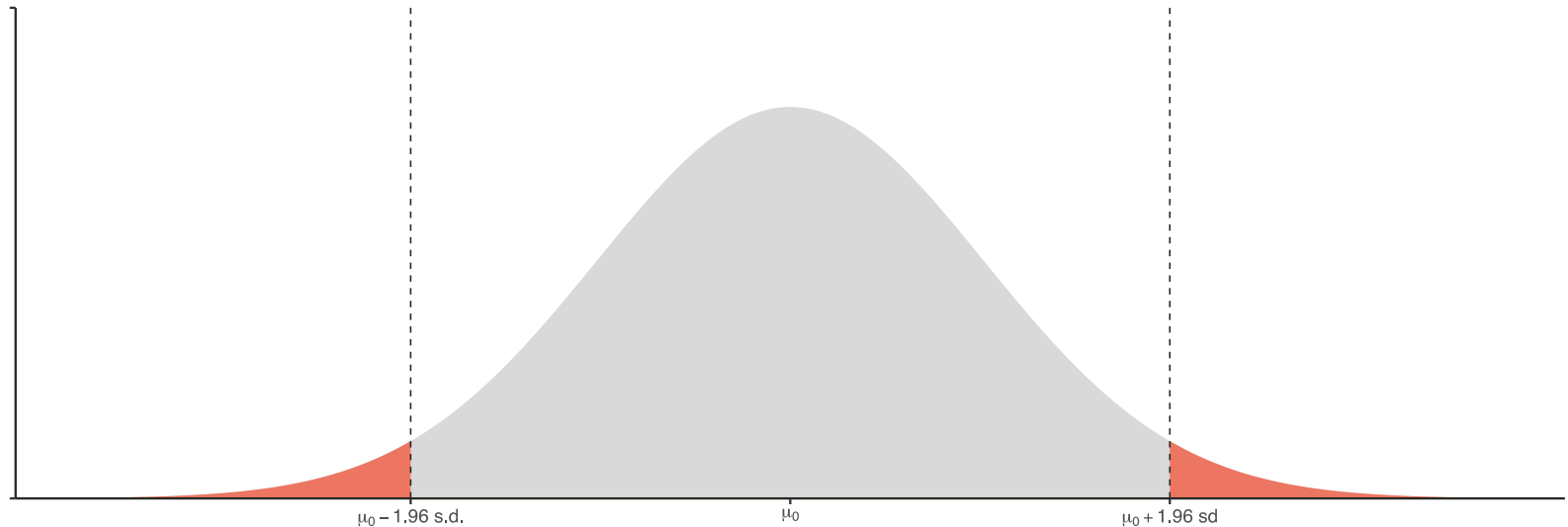
Our test should take extreme values of  $\hat{\mu}$  as evidence against the null hypothesis, but it should also weight them by what we know about the variance of  $\hat{\mu}$ .

- For now, we'll assume that the variable of interest  $X$  is normally distributed with mean  $\mu$  and standard deviation  $\sigma$ .



# Hypothesis Testing

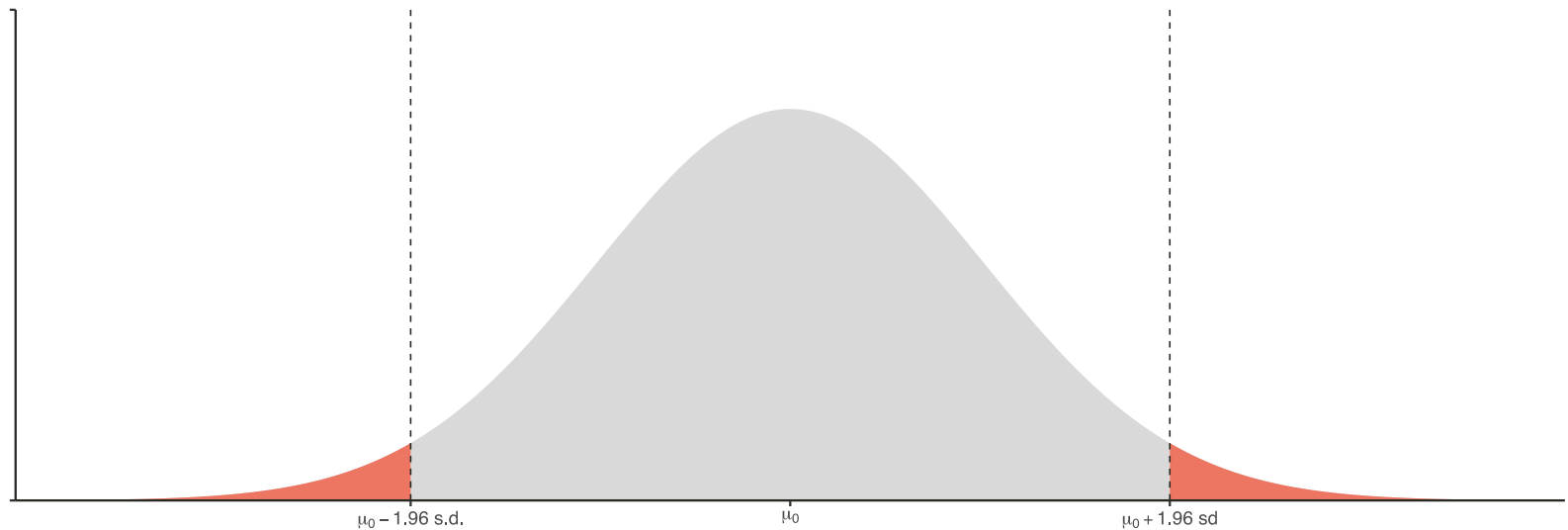
Reject  $H_0$  if  $\hat{\mu}$  lies in the **rejection region**.



- The area of the rejection region is defined by the **significance level** of the test.
- In a 5% test, the area is 0.05.
- Significance level = tolerance for Type I error.

# Hypothesis Testing

Reject  $H_0$  if  $|z| = \left| \frac{\hat{\mu} - \mu_0}{\text{sd}(\hat{\mu})} \right| > 1.96$ .



What happens to  $z$  as  $|\hat{\mu} - \mu_0|$  increases?

What happens to  $z$  as  $\text{sd}(\hat{\mu})$  increases?

# Hypothesis Testing

The formula for the  $z$  statistic assumes that we know  $\text{sd}(\hat{\mu})$ .

- In practice, we don't know  $\text{sd}(\hat{\mu})$ , so we have to estimate it.

If the variance of  $X$  is  $\sigma^2$ , then

$$\sigma_{\hat{\mu}}^2 = \frac{\sigma^2}{n}.$$

- We can estimate  $\sigma^2$  with the sample variance  $S_X^2$ .

The sample variance of the sample mean is

$$S_{\hat{\mu}}^2 = \frac{1}{n(n-1)} \sum_{i=1}^n (X_i - \bar{X})^2.$$

# Hypothesis Testing

The **standard error** of  $\hat{\mu}$  is the square root of  $S_{\hat{\mu}}^2$ :

$$\text{SE}(\hat{\mu}) = \sqrt{\frac{1}{n(n-1)} \sum_{i=1}^n (X_i - \bar{X})^2}.$$

- Standard error = sample standard deviation of an estimator.

When we use  $\text{SE}(\hat{\mu})$  in place of  $\text{sd}(\hat{\mu})$ , the  $z$  statistic becomes a  $t$  statistic:

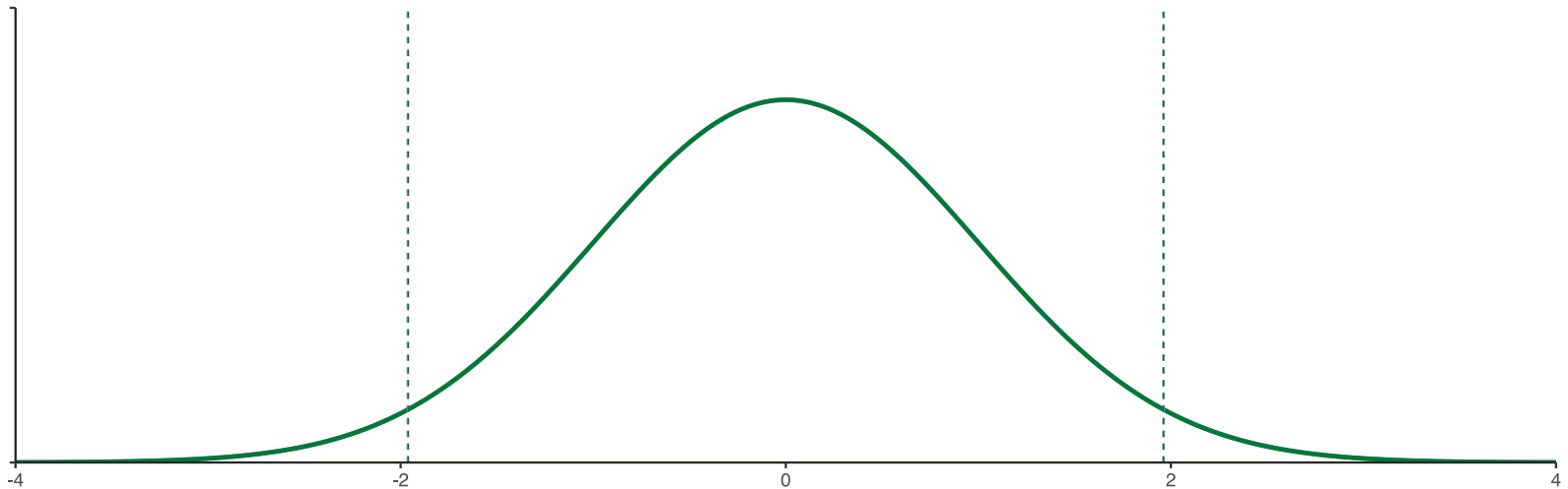
$$t = \frac{\hat{\mu} - \mu_0}{\text{SE}(\hat{\mu})}.$$

- Unlike the standard deviation of  $\hat{\mu}$ ,  $\text{SE}(\hat{\mu})$  varies from sample to sample.
- **Consequence:**  $t$  statistics do not necessarily have a normal distribution.

# Hypothesis Testing

## Normal distribution vs. t distribution

- A normal distribution has the same shape for any sample size.
- The shape of the t distribution depends the **degrees of freedom**.



- Degrees of freedom = 500.

# Hypothesis Testing

## t Tests (two-sided)

To conduct a t test, compare the  $t$  statistic to the appropriate **critical value** of the t distribution.

- To find the critical value in a t table, we need the degrees of freedom and the significance level  $\alpha$ .

Reject  $H_0$  at the  $\alpha \cdot 100$ -percent level if

$$|t| = \left| \frac{\hat{\mu} - \mu_0}{\text{SE}(\hat{\mu})} \right| > t_{\text{crit}}.$$

# Hypothesis Testing

## On Your Own

As the term progresses, we will encounter additional flavors of hypothesis testing and other related concepts.

You may find it helpful to review the following topics from Math 243:

- Confidence intervals
- One-sided  $t$  tests
- $p$  values

# Working with Data



# Data

## Experimental data

Data generated in controlled, laboratory settings.

Ideal for **causal identification**, but difficult to obtain in the social sciences.

- Intractable logistical problems
- Too expensive
- Morally repugnant

Experiments outside the lab: **randomized control trials** and **A/B testing**.

# Data

## Observational data

Data generated in non-experimental settings.

- Surveys
- Censuses
- Administrative records
- Environmental data
- Financial and sales transactions
- Social media

Mainstay of economic research, but **poses challenges** to causal identification.

# Tidy Data

Search:

|   | State      | Population | Murders |
|---|------------|------------|---------|
| 1 | Alabama    | 4779736    | 135     |
| 2 | Alaska     | 710231     | 19      |
| 3 | Arizona    | 6392017    | 232     |
| 4 | Arkansas   | 2915918    | 93      |
| 5 | California | 37253956   | 1257    |
| 6 | Colorado   | 5029196    | 65      |

Showing 1 to 6 of 51 entries

Previous

Next

**Rows** represent **observations**.

**Columns** represent **variables**.

Each **value** is associated with an **observation** and a **variable**.

# Cross Sectional Data

**Sample of individuals from a population at a point in time.**

Ideally, collected using **random sampling**.

- Random sampling + sufficient sample size = representative sample.
- Random sampling simplifies data analysis, but non-random samples are common (and difficult to work with).

Used extensively in applied microeconomics.\*

**Main focus of this course.**

\* Applied microeconomics = Labor, health, education, public finance, development, industrial organization, and urban economics.

# Time Series Data

## Observations of variables over time.

- Quarterly US GDP
- Annual US infant mortality rates
- Daily Amazon stock prices

Complication: Observations are not independent draws.

- GDP this quarter highly related to GDP last quarter.

Used extensively in empirical macroeconomics.

Requires more-advanced methods (EC 421 and EC 422).

# Pooled Cross Sectional Data

**Cross sections from different points in time.**

Useful for studying policy changes and relationship that change over time.

Requires more-advanced methods (EC 421 and many 400-level applied micro classes).

# Panel or Longitudinal Data

## Time series for each cross-sectional unit.

- Example: daily attendance data for a sample of students.

Difficult to collect, but useful for causal identification.

- Can control for *unobserved* characteristics.

Requires more-advanced methods (EC 421 and many 400-level applied micro classes).

# Messy Data

**Analysis-ready datasets are rare.** Most data are "messy."

The focus of this class is data analysis, but **data wrangling** is a non-trivial part of a data scientist/analyst's job.

R has a suite of packages that facilitate data wrangling.

- `readr`, `tidyr`, `dplyr`, `ggplot2` + others.
- Known collectively as the `tidyverse`.



# tidyverse

## The **tidyverse**: A package of packages

**readr**: Functions to import data.

**tidyr**: Functions to reshape messy data.

**dplyr**: Functions to work with data.

**ggplot2**: Functions to visualize data.

# Workflow

- Step 1: Load packages with `pacman`
- Step 2: Import data with `readr`
- Step 3: Reshape data with `tidyr`
- Step 4: Manipulate data with `dplyr`
- Step 5: Visualize and analyze data with `ggplot2`

# Why Bother?

**Q:** Why not just use **MS Excel** for data wrangling?

**A: Reproducibility**

- Easier to retrace your steps with R.

**A: Portability**

- Easy to re-purpose R code for new projects.

**A: Scalability**

- Excel chokes on big datasets.

**A: R Saves time** (eventually)

- Lower marginal costs in exchange for higher fixed costs.