

# Week 3: Learning from Random Samples

Brandon Stewart<sup>1</sup>

Princeton

September 24/26, 2018

---

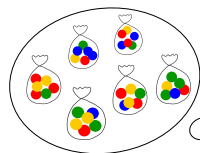
<sup>1</sup>These slides are heavily influenced by Matt Blackwell, Adam Glynn, Justin Grimmer and Jens Hainmueller. Some illustrations by Shay O'Brien.

# Where We've Been and Where We're Going...

- Last Week
  - ▶ random variables
  - ▶ joint distributions
- This Week
  - ▶ Monday: Point Estimation
    - ★ sampling and sampling distributions
    - ★ point estimates
    - ★ properties (bias, variance, consistency)
  - ▶ Wednesday: Interval Estimation
    - ★ confidence intervals
    - ★ comparing two groups
- Next Week
  - ▶ hypothesis testing
  - ▶ what is regression?
- Long Run
  - ▶ probability  $\rightarrow$  inference  $\rightarrow$  regression  $\rightarrow$  causal inference

Questions?

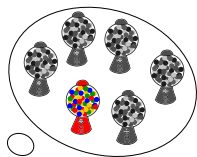
# Where We've Been and Where We're Going...



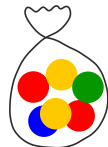
Probability



Data generating  
process



Inference



Observed  
data

## Where We've Been and Where We're Going. . .

**Inference:** given that we observe something in the data, what is our best guess of what it will be in other samples?

- For the last few classes we have talked about **probability**—that is, if we knew how the world worked, we are describing what kind of data we should expect.
- Now we want to move the other way. If we have a set of data, can we estimate the various parts of the probability distributions that we have talked about. Can we **estimate** the mean, the variance, the covariance, etc?
- Moving forward this is going to be very important. Why? Because we are going to want to estimate the population **conditional expectation** in regression.

# Primary Goals for This Week

We want to be able to interpret the numbers in this table (and a couple of numbers that can be derived from these numbers).

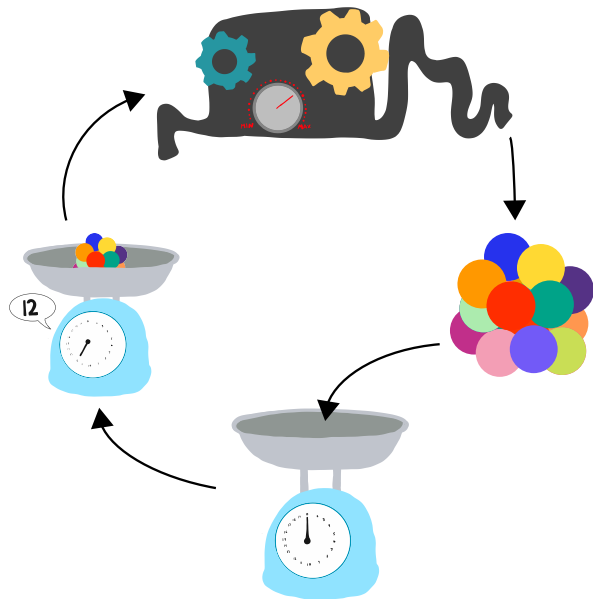
**Table 1. Mean Level of Anger Toward A Black Family Moving in Next Door, by Region (Whites Only)**

Region	Experimental Condition		Estimated Percent Angry
	Baseline	Black Family	
Non-South	2.28 <sup>a</sup> (.07)	2.24 (.05)	0
	425 <sup>b</sup>	461	
South	1.95 (.06)	2.37 (.08)	42
	139	136	

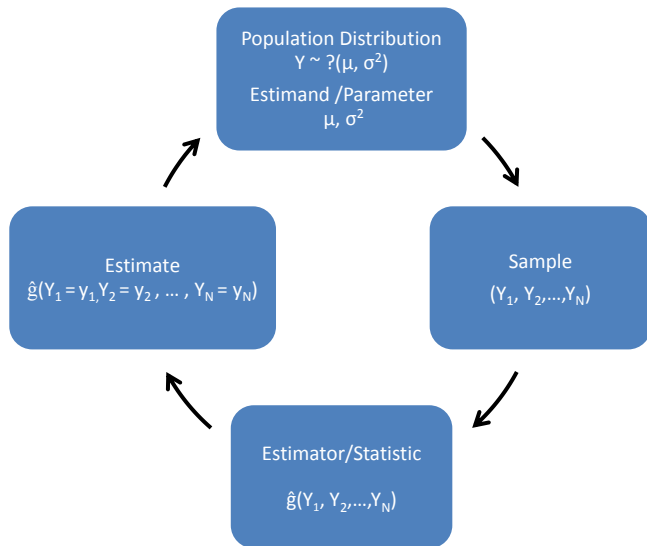
<sup>a</sup>Standard error of the estimate.

<sup>b</sup>Number of cases.

# An Overview



# An Overview



# 1 Populations, Sampling, Sampling Distributions

- Conceptual
- Mathematical

## 2 Overview of Point Estimation

## 3 Properties of Estimators

## 4 Review and Example

## 5 Fun With Hidden Populations

## 6 Interval Estimation

## 7 Large Sample Intervals for a Mean

- Simple Example
- Kuklinski Example

## 8 Small Sample Intervals for a Mean

## 9 Comparing Two Groups

## 10 Fun With Correlation

## 11 Appendix: $\chi^2$ and $t$ -distribution



# Populations

- Typically, we want to learn about the distribution of random variable (or set of random variables) for a **population** of interest.
- e.g. the distribution of votes for Hillary Clinton in the population of registered voters in the United States. This is an example of a **finite population**.
- Sometimes the population will be more abstract, such as the population of all possible television ads. This is an example of an **infinite population**.
- With either a finite or infinite population our main goal in inference is to learn about the **population distribution** or particular aspects of that distribution, like the mean or variance, which we call a **population parameter** (or just parameter).

# Population Distribution

- We sometimes call the population distribution the **data generating process** and represent it with a pmf or pdf,  $f(x; \theta)$ .
- Ideally we would place no restrictions on  $f$  and learn everything we can about it from the data. This **nonparametric** approach is difficult due to the fact that the space of possible distributions is vast!
- Instead, we will often make a **parametric** assumption and assume that the formula for  $f$  is known up to some unknown parameters.
- Thus,  $f$  has two parts: the known part which is the formula for the pmf/pdf (sometimes called the parametric model and comes from the distributional assumptions) and the unknown part, which are the parameters,  $\theta$ .

# Population Distribution

- For instance, suppose we have a binary r.v. such as intending to vote for Hillary Clinton ( $X = 1$ ). Then we might assume that the population distribution is Bernoulli with unknown probability of  $X = 1$ ,  $\theta$ .
- Thus we would have:

$$f(x; \theta) = \theta^x(1 - \theta)^{1-x}$$

Our goal is to learn about the probability of someone voting for Hillary Clinton from a sample of draws from this distribution.

- Probability tells us what types of samples we should expect for different values of  $\theta$ .
- For some problems, such as estimating the mean of a distribution, we actually won't need to specify a parametric model for the distribution allowing us to take an **agnostic** view of statistics.

# Using Random Samples to Estimate Population Parameters

Let's look at possible estimators  $\hat{\mu}$  for the (unobserved) mean income in the population  $\mu$ .

How can we use sample data to estimate  $\mu$ ?

If we think of a sample of size  $n$  as **randomly sampled** with replacement from the population, then  $Y_1, \dots, Y_n$  are **independently and identically distributed** (i.i.d.) random variables with  $E[Y_i] = \mu$  and  $V[Y_i] = \sigma^2$  for all  $i \in \{1, \dots, n\}$ .

In R, we can draw an i.i.d. random sample of size  $n$  from population by

```
sample(population, size=n, replace=TRUE)
```

where `population` is a vector that contains the  $Y_i$  values for all units in the population.

Our estimators,  $\hat{\mu}$ , are functions of  $Y_1, \dots, Y_n$  and will therefore be random variables with their own probability distributions.

# Why Samples?

- The population is infinite (e.g., ballot order randomization **process**).
- The population is large (e.g., **all** Southern adults with telephones).
- The population includes **counterfactuals**.

Even if our population is limited to the surveyed individuals in the Kuklinski et al. study, we might take the population of interest to include potential responses for each individual to **both** the treatment **and** baseline questions. For each individual we must sample one of these two responses (because we cannot credibly ask both questions).

This occurs whenever we are interested in making causal inferences.

# Estimands, Estimators, and Estimates

The goal of statistical inference is to learn about the unobserved population distribution, which can be characterized by **parameters**.

- **Estimands** are the parameters that we aim to estimate. Often written with greek letters (e.g.  $\mu, \theta$ , population mean) :  $\frac{1}{N} \sum_{i=1}^N y_i$
- **Estimators** are functions of sample data (i.e. **statistics**) which we use to learn about the estimands. Often denoted with a “hat” (e.g.  $\hat{\mu}, \hat{\theta}$ )
- **Estimates** are particular values of estimators that are realized in a given sample (e.g. sample mean):  $\frac{1}{n} \sum_{i=1}^n y_i$



# What Are We Estimating?

- Our goal is to learn about the **data generating process** that generated the sample
- We might assume that the data  $Y_1, \dots, Y_n$  are i.i.d. draws from the **population distribution** with p.m.f. or p.d.f. of a certain form  $f_Y()$  indexed by unknown parameters  $\theta$
- Even without a full probability model we can estimate particular properties of a distribution such as the mean  $E[Y_i] = \mu$  or the variance  $V[Y_i] = \sigma^2$
- An **estimator**  $\hat{\theta}$  of some parameter  $\theta$ , is a **function** of the sample  $\hat{\theta} = h(Y_1, \dots, Y_n)$  and thus is a **random variable**.

# Why Study Estimators?

- Two Goals:
  - ① **Inference**: How much uncertainty do we have in this estimate?
  - ② **Evaluate Estimators**: How do we choose which estimator to use?
- We will consider the hypothetical **sampling distribution** of estimates we would have obtained if we had drawn **repeated samples** of size  $n$  from the population.
- In real applications, we cannot draw repeated samples, so we attempt to approximate the sampling distribution (either by resampling or by mathematical formulas)



# Sampling Distribution of the Sample Mean

Consider a toy example where we have the population and can construct sampling distributions.

Repeated Sampling Procedure:

- 1 Take a simple random sample of size  $n = 4$ .
- 2 Calculate the sample mean.
- 3 Repeat steps 1 and 2 at least 10,000 times.
- 4 Plot the sampling distribution of the sample means (maybe as a histogram).

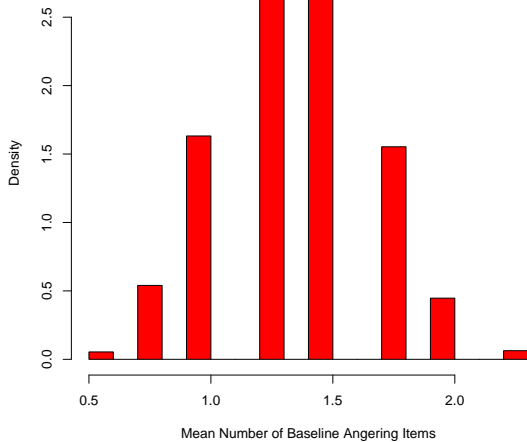
# Repeated Sampling Procedure

```
# population data
ypop <- c(rep(0,0),rep(1,17),rep(2,10),rep(3,4))

# simulate the sampling distribution of the sample mean

SamDistMeans <- replicate(10000, mean(sample(ypop,size=4,replace=TRUE)))
```

# Sampling Distribution of the Sample Mean



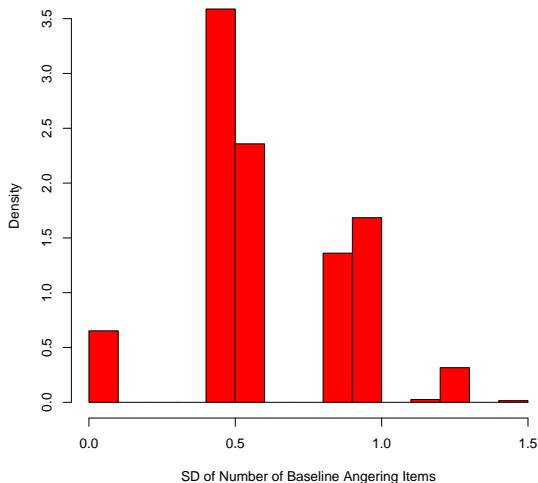
# Sampling Distribution of the Sample Standard Deviation

We can consider sampling distributions for other sample statistics (e.g., the sample standard deviation).

Repeated Sampling Procedure:

- 1 Take a simple random sample of size  $n = 4$ .
- 2 Calculate the sample **standard deviation**.
- 3 Repeat steps 1 and 2 at least 10,000 times.
- 4 Plot the sampling distribution of the sample standard deviations (maybe as a histogram).

# Sampling Distribution of the Sample Standard Deviation



# Standard Error

We refer to the standard deviation of a sampling distribution as a **standard error**.

Two Points of Potential Confusion:

- Each sampling distribution has its own standard deviation, and therefore its own standard error. (.35 for mean, .30 for sd)
- Some people refer to an estimated standard error as the standard error.

# 1 Populations, Sampling, Sampling Distributions

- Conceptual
- Mathematical

## 2 Overview of Point Estimation

## 3 Properties of Estimators

## 4 Review and Example

## 5 Fun With Hidden Populations

## 6 Interval Estimation

## 7 Large Sample Intervals for a Mean

- Simple Example
- Kuklinski Example

## 8 Small Sample Intervals for a Mean

## 9 Comparing Two Groups

## 10 Fun With Correlation

## 11 Appendix: $\chi^2$ and $t$ -distribution

## Notation for Sampling Distributions

Suppose we took a simple random sample with replacement from the population.

We say that  $X_1, X_2, \dots, X_n$  are identically and independently distributed from a population distribution with a mean ( $E[X_1] = \mu$ ) and a variance ( $V[X_1] = \sigma^2$ ).

Then we write  $X_1, X_2, \dots, X_n \sim_{i.i.d} ?(\mu, \sigma^2)$



# Describing the Sampling Distribution for the Mean

We would like a full description of the sampling distribution for the mean, but it will be useful to separate this description into three parts.

If we assume that  $X_1, \dots, X_n \sim_{i.i.d} ?(\mu, \sigma^2)$ , then we would like to identify the following things about  $\bar{X}_n$ .

- $E[\bar{X}_n]$
- $V[\bar{X}_n]$
- ?

## Expectation of $\bar{X}_n$

Again, let  $X_1, X_2, \dots, X_n$  be identically and independently distributed from a population distribution with a mean ( $E[X_1] = \mu$ ) and a variance ( $V[X_1] = \sigma^2$ ). Using the properties of expectation, calculate

$$\begin{aligned} E[\bar{X}_n] &= E\left[\frac{1}{n} \sum_{i=1}^n X_i\right] \\ &=? \end{aligned}$$

## Variance of $\bar{X}_n$

Again, let  $X_1, X_2, \dots, X_n$  be identically and independently distributed from a population distribution with a mean ( $E[X_1] = \mu$ ) and a variance ( $V[X_1] = \sigma^2$ ). Using the properties of variances, calculate

$$\begin{aligned} V[\bar{X}_n] &= V\left[\frac{1}{n} \sum_{i=1}^n X_i\right] \\ &=? \end{aligned}$$

## What about the “?”

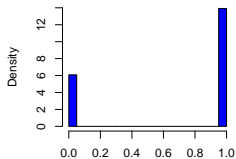
If  $X_1, \dots, X_n \sim i.i.d. N(\mu, \sigma^2)$ , then

$$\bar{X}_n \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

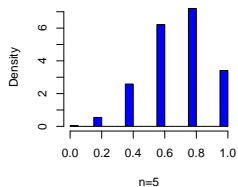
What if  $X_1, \dots, X_n$  are not normally distributed?

# Bernoulli (Coin Flip) Distribution

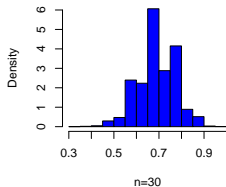
Population Distribution



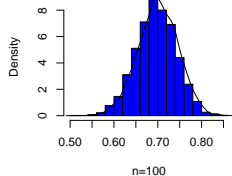
Sampling Distribution of the Mean



Sampling Distribution of the Mean

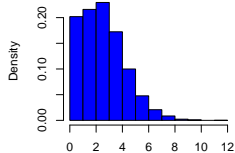


Sampling Distribution of the Mean

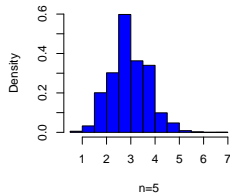


# Poisson (Count) Distribution

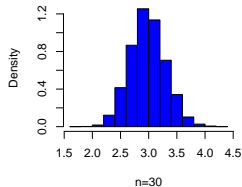
Population Distribution



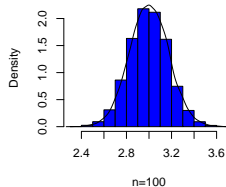
Sampling Distribution of the Mean



Sampling Distribution of the Mean

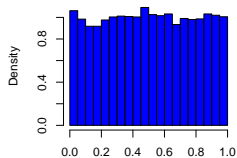


Sampling Distribution of the Mean

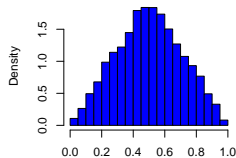


# Uniform Distribution

**Population Distribution**

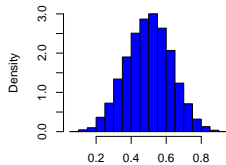


**Sampling Distribution of the Mean**



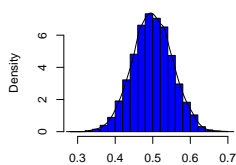
$n=2$

**Sampling Distribution of the Mean**



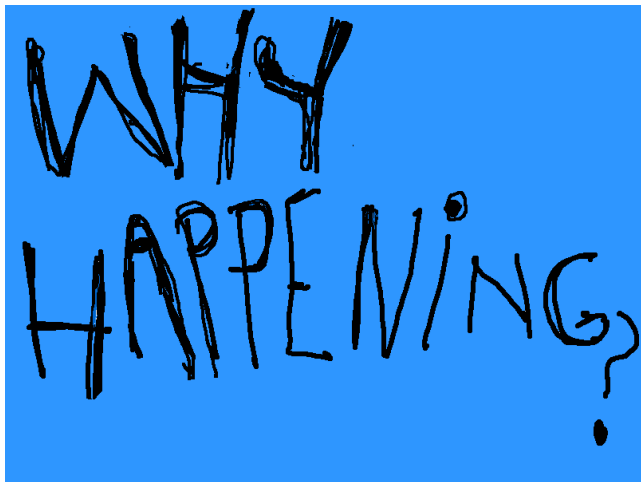
$n=5$

**Sampling Distribution of the Mean**



$n=30$

Why would this be true?



Images from *Hyperbole and a Half* by Allie Brosh.



# The Central Limit Theorem

In the previous slides, as  $n$  increases, the sampling distribution of  $\bar{X}_n$  appeared to become more bell-shaped. This is the basic implication of the **Central Limit Theorem**:

If  $X_1, \dots, X_n \sim_{i.i.d.} (\mu, \sigma^2)$  and  $n$  is large, then

$$\bar{X}_n \sim_{approx} N\left(\mu, \frac{\sigma^2}{n}\right)$$

## The Central Limit Theorem: What are we glossing over?

To understand the Central Limit Theorem mathematically we need a few basic definitions in place first.

### Definition (Convergence in Probability)

A sequence  $X_1, \dots, X_n$  of random variables **converges in probability** towards a real number  $a$  if, for all accuracy levels  $\varepsilon > 0$ ,

$$\lim_{n \rightarrow \infty} \Pr(|X_n - a| \geq \varepsilon) = 0$$

We write this as

$$X_n \xrightarrow{p} a \quad \text{or} \quad \text{plim}_{n \rightarrow \infty} X_n = a.$$

# The Central Limit Theorem: What are we glossing over?

## Definition (Law of Large Numbers)

Let  $X_1, X_2, \dots, X_n$  be a sequence of i.i.d. random variables, each with finite mean  $\mu$ . Then for all  $\varepsilon > 0$ ,

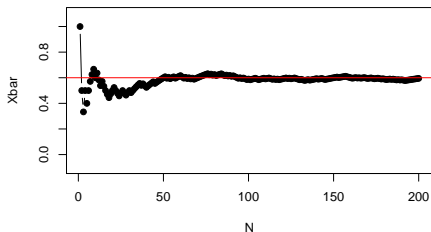
$$\bar{X}_n \xrightarrow{p} \mu \text{ as } n \rightarrow \infty$$

or equivalently,

$$\lim_{n \rightarrow \infty} \Pr(|\bar{X}_n - \mu| \geq \varepsilon) = 0$$

where  $\bar{X}_n$  is the sample mean.

Example: Mean of  $N$  independent tosses of a coin:



# The Central Limit Theorem: What are we glossing over?

## Definition (Convergence in Distribution)

Consider a sequence of random variables  $X_1, \dots, X_n$ , each with CDFs  $F_1, \dots, F_n$ . The sequence is said to **converge in distribution** to a limiting random variable  $X$  with CDF  $F$  if

$$\lim_{n \rightarrow \infty} F_n(x) = F(x),$$

for every point  $x$  at which  $F$  is continuous. We write this as

$$X_N \xrightarrow{d} X.$$

- As  $n$  grows, the distribution of  $X_n$  converges to the distribution of  $X$ .
- Convergence in probability is a special case of convergence in distribution in which the distribution converges to a **degenerate distribution** (i.e. a probability distribution which only takes a single value).

# The Central Limit Theorem: What are we glossing over?

## Definition (Lindeberg-Lévy Central Limit Theorem)

Let  $X_1, \dots, X_n$  a sequence of i.i.d. random variables each with mean  $\mu$  and variance  $\sigma^2 < \infty$ . Then, for **any** population distribution of  $X$ ,

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} \mathcal{N}(0, \sigma^2).$$

- As  $n$  grows, the  $\sqrt{n}$ -scaled sample mean converges to a normal random variable.
- CLT also implies that the **standardized** sample mean converges to a standard normal random variable:

$$Z_n \equiv \frac{\bar{X}_n - E[\bar{X}_n]}{\sqrt{V[\bar{X}_n]}} = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \xrightarrow{d} \mathcal{N}(0, 1).$$

- Note that CLT holds for a random sample from *any* population distribution (with finite mean and variance) — what a convenient result!

## Question:

As the number of observations in a dataset increases, which of the following is true?

- A) The distribution of  $X$  becomes more normally distributed.
- B) The distribution of  $\bar{X}$  becomes more normally distributed.
- C) Both statements are true.

# Point Estimation

Suppose we are primarily interested in specific characteristics of the population distribution.

For example, suppose we are primarily interested in  $E[X]$ .

We refer to characteristics of the population distribution (e.g.,  $E[X]$ ) as **parameters**. These are often denoted with a greek letter (e.g.  $\mu$ ).

We use a statistic (e.g.,  $\bar{X}$ ) to estimate a parameter, and we will denote this with a hat (e.g.  $\hat{\mu}$ ). A **statistic** is a function of the sample.

# Why Point Estimation?

- Estimating one number is typically easier than estimating many (or an infinite number of) numbers.
- The question of interest may be answerable with single characteristic of the distribution (e.g., if  $E[Y] - E[X]$  identifies the proportion angered by the sensitive item, then it may be sufficient to estimate  $E[Y]$  and  $E[X]$ )



## Estimators for $\mu$

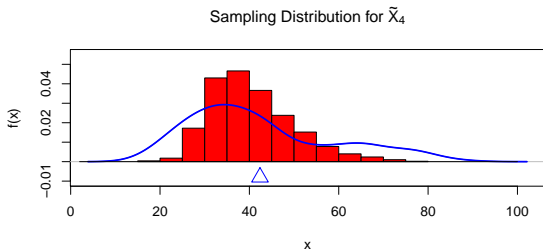
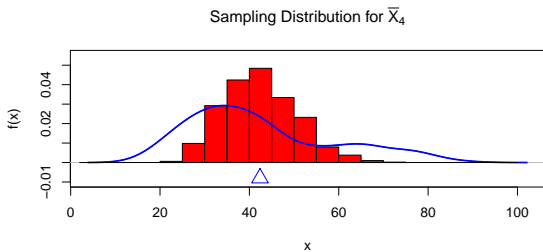
Some possible estimators  $\hat{\mu}$  for the balance point  $\mu$ :

A)  $\bar{X}_n = \frac{1}{n}(X_1 + \dots + X_n)$ , the sample average

B)  $\tilde{X}_n = \text{median}(X_1, \dots, X_n)$ , the sample median

Clearly, one of these estimators is better than the other, but how can we define “better”?

# Age population distribution in blue, sampling distributions in red



# Methods of Finding Estimators

We will primarily discuss the **Method of Least Squares** for finding estimators in this course. However, many of the estimators we discuss can also be derived by **Method of Moments** or **Method of Maximum Likelihood** (covered in Soc504).

When estimating simple features of a distribution we can use the **plug-in principle**, the idea that you write down the feature of the distribution you are interested in and estimate with the sample analog. Formally this is using the Empirical CDF to estimate features of the population.

# Desirable Properties of Estimators

Sometimes there are many possible estimators for a given parameter. Which one should we choose?

- We'd like an estimator that gets the right answer on average.
- We'd like an estimator that doesn't change much from sample to sample.
- We'd like an estimator that gets closer to the right answer (probabilistically) as the sample size increases.
- We'd like an estimator that has a known sampling distribution (approximately) when the sample size is large.

# Properties of Estimators

Estimators are random variables, for which randomness comes from **repeated sampling** from the population.

The distribution of an estimator due to repeated sampling is called the **sampling distribution**.

The properties of an estimator refer to the characteristics of its sampling distribution.

Finite-sample Properties (apply for any sample size):

- **Unbiasedness**: Is the sampling distribution of our estimator centered at the true parameter value?  $E[\hat{\mu}] = \mu$
- **Efficiency**: Is the variance of the sampling distribution of our estimator reasonably small?  $V[\hat{\mu}_1] < V[\hat{\mu}_2]$

Asymptotic Properties (kick in when  $n$  is large):

- **Consistency**: As our sample size grows to infinity, does the sampling distribution of our estimator converge to the true parameter value?
- **Asymptotic Normality**: As our sample size grows large, does the sampling distribution of our estimator approach a normal distribution?

# 1: Bias (Not Getting the Right Answer on Average)

## Definition

**Bias** is the expected difference between the **estimator** and the parameter. Over repeated samples, an unbiased estimator is right on average.

$$\begin{aligned}\text{Bias}(\hat{\mu}) &= E[\hat{\mu} - E[X]] \\ &= E[\hat{\mu}] - \mu\end{aligned}$$

Bias is **not** the difference between a particular estimate and the parameter. For example,

$$\text{Bias}(\bar{X}_n) \neq E[\bar{x}_n - E[X]]$$

An estimator is **unbiased** iff:

$$\text{Bias}(\hat{\mu}) = 0$$

## Example: Estimators for Population Mean

Candidate estimators:

- 1  $\hat{\mu}_1 = Y_1$  (the first observation)
- 2  $\hat{\mu}_2 = \frac{1}{2}(Y_1 + Y_n)$  (average of the first and last observation)
- 3  $\hat{\mu}_3 = 42$
- 4  $\hat{\mu}_4 = \bar{Y}_n$  (the sample average)

How do we choose between these estimators?

# Bias of Example Estimators

Which of these estimators are unbiased?

①  $E[Y_1 - \mu] = \mu - \mu = 0$

②  $E[\frac{1}{2}(Y_1 + Y_n) - \mu] = \frac{1}{2}(E[Y_1] + E[Y_n]) - \mu = \frac{1}{2}(\mu + \mu) - \mu = 0$

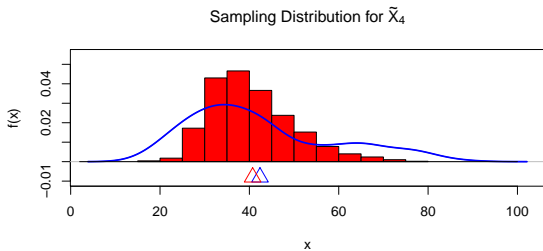
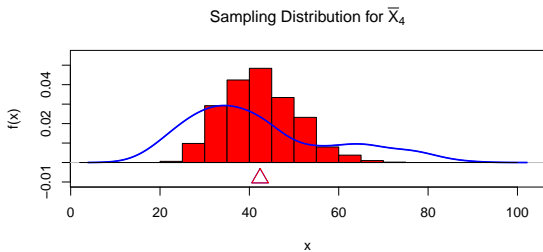
③  $E[42 - \mu] = 42 - \mu$

④  $E[\bar{Y}_n - \mu] = \frac{1}{n} \sum_1^n E[Y_i] - \mu = \mu - \mu = 0$

- Estimators 1, 2, and 4 are unbiased because they get the right answer on average.
- Estimator 3 is biased.



# Age population distribution in blue, sampling distributions in red



## 2: Efficiency (doesn't change much sample to sample)

- How should we choose among unbiased estimators?
- All else equal, we prefer estimators that have a sampling distribution with smaller variance.

### Definition (Efficiency)

If  $\hat{\theta}_1$  and  $\hat{\theta}_2$  are unbiased estimators of  $\theta$ , then  $\hat{\theta}_1$  is **more efficient** relative to  $\hat{\theta}_2$  iff

$$V[\hat{\theta}_1] < V[\hat{\theta}_2]$$

- Under repeated sampling, estimates based on  $\hat{\theta}_1$  are likely to be closer to  $\theta$
- Note that this does **not** imply that a particular estimate is always close to the true parameter value
- The standard deviation of the sampling distribution of an estimator,  $\sqrt{V[\hat{\theta}]}$ , is often called the **standard error** of the estimator

Aronow and Miller discuss efficiency in terms of MSE (more on this in a second).

## Variance of Example Estimators

What is the variance of our estimators?

①  $V[Y_1] = \sigma^2$

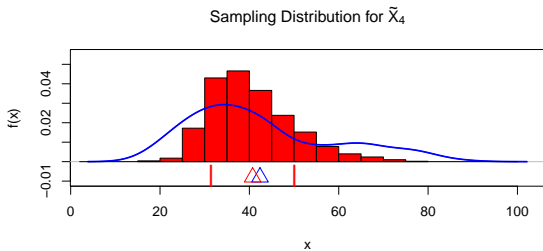
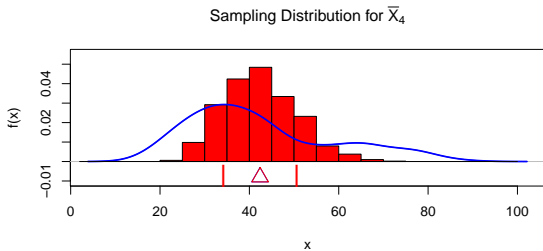
②  $V[\frac{1}{2}(Y_1 + Y_n)] = \frac{1}{4}V[Y_1 + Y_n] = \frac{1}{4}(\sigma^2 + \sigma^2) = \frac{1}{2}\sigma^2$

③  $V[42] = 0$

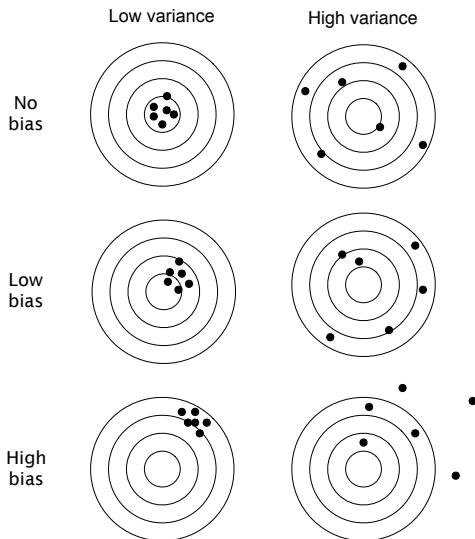
④  $V[\bar{Y}_n] = \frac{1}{n^2} \sum_1^n V[Y_i] = \frac{1}{n^2} n\sigma^2 = \frac{1}{n}\sigma^2$

Among the unbiased estimators, the sample average has the smallest variance. This means that Estimator 4 (the sample average) is likely to be closer to the true value  $\mu$ , than Estimators 1 and 2.

# Age population distribution in blue, sampling distributions in red



# Choosing Estimators



Salganik (2018), Figure 3.1

# Mean Squared Error

How can we choose between an unbiased estimator and a biased, but more efficient estimator?

## Definition (Mean Squared Error)

To compare estimators in terms of both efficiency and unbiasedness we can use the **Mean Squared Error** (MSE), the expected squared difference between  $\hat{\theta}$  and  $\theta$ :

$$MSE(\hat{\theta}) = E[(\hat{\theta} - \theta)^2] = \text{Bias}(\hat{\theta})^2 + V(\hat{\theta}) = [E[\hat{\theta}] - \theta]^2 + V(\hat{\theta})$$

## 3-4: Asymptotic Evaluations

(what happens as sample size increases)

- Unbiasedness and efficiency are **finite-sample properties** of estimators, which hold regardless of sample size
- Estimators also have **asymptotic properties**, i.e., the characteristics of sampling distributions when sample size becomes infinitely large
- To define asymptotic properties, consider a **sequence** of estimators at increasing sample sizes:

$$\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_n$$

- For example, the sequence of sample means ( $\bar{X}_n$ ) is defined as:

$$\bar{X}_1, \bar{X}_2, \dots, \bar{X}_n = X_1, \frac{X_1 + X_2}{2}, \dots, \frac{X_1 + \dots + X_n}{n}$$

- Asymptotic properties of an estimator are defined by the behavior of  $\hat{\theta}_1, \dots, \hat{\theta}_n$  when  $n$  goes to infinity.

# Stochastic Convergence

- When a sequence of random variables stabilizes to a certain probabilistic behavior as  $n \rightarrow \infty$ , the sequence is said to show **stochastic convergence**.
- Two types of stochastic convergence are of particular importance:
  - 1 **Convergence in probability**: values in the sequence eventually take a constant value  
(i.e. the **limiting distribution** is a point mass)
  - 2 **Convergence in distribution**: values in the sequence continue to vary, but the variation eventually comes to follow an unchanging distribution  
(i.e. the limiting distribution is a well characterized distribution)



# Convergence in Probability

## Definition (Convergence in Probability)

A sequence  $X_1, \dots, X_n$  of random variables **converges in probability** towards a real number  $a$  if, for all accuracy levels  $\varepsilon > 0$ ,

$$\lim_{n \rightarrow \infty} \Pr (|X_n - a| \geq \varepsilon) = 0$$

We write this as

$$X_n \xrightarrow{p} a \quad \text{or} \quad \text{plim}_{n \rightarrow \infty} X_n = a.$$

- As  $n$  increases, almost all of the PDF/PMF of  $X_n$  will be concentrated in the  $\varepsilon$ -interval around  $a$ ,  $[a - \varepsilon, a + \varepsilon]$
- A sufficient (but not necessary) condition for convergence in probability:

$$E[X_n] \rightarrow a \quad \text{and} \quad V[X_n] \rightarrow 0 \quad \text{as} \quad n \rightarrow \infty$$

- For example, the sample mean  $\bar{X}_n$  converges to the population mean  $\mu$  in probability because

$$E[\bar{X}_n] = \mu \quad \text{and} \quad V[\bar{X}_n] = \sigma^2/n \rightarrow 0 \quad \text{as} \quad n \rightarrow \infty$$

### 3: Consistency

(does it get closer to the right answer as sample size increases)

#### Definition

An estimator  $\theta_n$  is **consistent** if the sequence  $\theta_1, \dots, \theta_n$  converges in probability to the true parameter value  $\theta$  as sample size  $n$  grows to infinity:

$$\theta_n \xrightarrow{P} \theta \quad \text{or} \quad \text{plim}_{n \rightarrow \infty} \theta_n = \theta$$

- Often seen as a minimal requirement for estimators
- A consistent estimator may still perform badly in small samples
- Two ways to verify consistency:
  - ① Analytic: Often easier to check if  $E[\theta_n] \rightarrow \theta$  and  $V[\theta_n] \rightarrow 0$
  - ② Simulation: Increase  $n$  and see how the sampling distribution changes
- Does unbiasedness imply consistency?
- Does consistency imply unbiasedness?

# Deriving Consistency of Estimators

Our candidate estimators:

- 1  $\hat{\mu}_1 = Y_1$
- 2  $\hat{\mu}_2 = 4$
- 3  $\hat{\mu}_3 = \bar{Y}_n \equiv \frac{1}{n}(Y_1 + \cdots + Y_n)$
- 4  $\hat{\mu}_4 = \tilde{Y}_n \equiv \frac{1}{n+5}(Y_1 + \cdots + Y_n)$

Which of these estimators are consistent for  $\mu$ ?

- 1  $E[\hat{\mu}_1] = \mu$  and  $V[\hat{\mu}_1] = \sigma^2$
- 2  $E[\hat{\mu}_2] = 4$  and  $V[\hat{\mu}_2] = 0$
- 3  $E[\hat{\mu}_3] = \mu$  and  $V[\hat{\mu}_3] = \frac{1}{n}\sigma^2$
- 4  $E[\hat{\mu}_4] = \frac{n}{n+5}\mu$  and  $V[\hat{\mu}_4] = \frac{n}{(n+5)^2}\sigma^2$

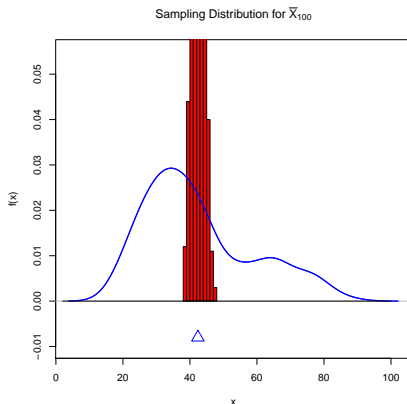
# Consistency

The sample mean is a consistent estimator for  $\mu$ .

$$\bar{X}_n \sim_{\text{approx}} N\left(\mu, \frac{\sigma^2}{n}\right)$$

As  $n$  increases,  $\frac{\sigma^2}{n}$  approaches 0.

$$n = 125100$$



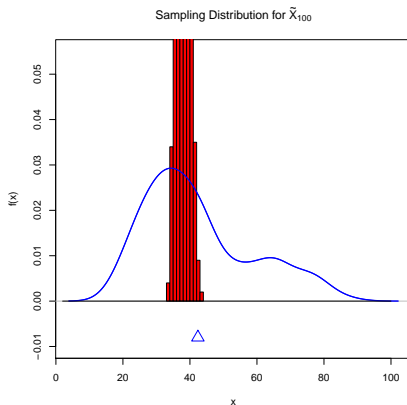
# Inconsistency

An estimator can be inconsistent in several ways:

- The sampling distribution collapses around the wrong value
- The sampling distribution never collapses around anything

# Inconsistency

Consider the median estimator:  $\tilde{X}_n = \text{median}(Y_1, \dots, Y_n)$  Is this estimator consistent for the expectation?  
 $n = 125100$



## 4: Asymptotic Distribution (known sampling distribution for large sample size)

We are also interested in the shape of the sampling distribution of an estimator as the sample size increases.

The sampling distributions of many estimators converge towards a normal distribution.

For example, we've seen that the sampling distribution of the sample mean converges to the normal distribution.

# Mean Squared Error

How can we choose between an unbiased estimator and a biased, but more efficient estimator?

## Definition (Mean Squared Error)

To compare estimators in terms of both efficiency and unbiasedness we can use the **Mean Squared Error** (MSE), the expected squared difference between  $\hat{\theta}$  and  $\theta$ :

$$MSE(\hat{\theta}) = E[(\hat{\theta} - \theta)^2] = \text{Bias}(\hat{\theta})^2 + V(\hat{\theta}) = [E[\hat{\theta}] - \theta]^2 + V(\hat{\theta})$$



# Review and Example

Gerber, Green, and Larimer (*American Political Science Review*, 2008)

Dear Registered Voter:

## WHAT IF YOUR NEIGHBORS KNEW WHETHER YOU VOTED?

Why do so many people fail to vote? We've been talking about the problem for years, but it only seems to get worse. This year, we're taking a new approach. We're sending this mailing to you and your neighbors to publicize who does and does not vote.

The chart shows the names of some of your neighbors, showing which have voted in the past. After the August 8 election, we intend to mail an updated chart. You and your neighbors will all know who voted and who did not.

## DO YOUR CIVIC DUTY — VOTE!

---

MAPLE DR	Aug 04	Nov 04	Aug 06
9995 JOSEPH JAMES SMITH	Voted	Voted	_____
9995 JENNIFER KAY SMITH		Voted	_____
9997 RICHARD B JACKSON		Voted	_____
9999 KATHY MARIE JACKSON		Voted	_____

## Basic Analysis

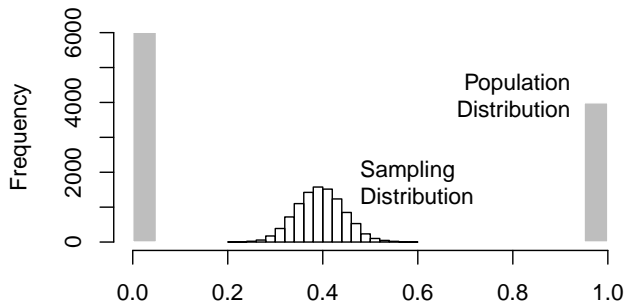
```
load("gerber_green_larimer.RData")
## turn turnout variable into a numeric
social$voted <- 1 * (social$voted == "Yes")
neigh.mean <- mean(social$voted[social$treatment == "Neighbors"])
neigh.mean
contr.mean <- mean(social$voted[social$treatment == "Civic Duty"])
contr.mean
neigh.mean - contr.mean
```

$$.378 - .315 = .063$$

Is this a “real” effect? Is it big?

# Population vs. Sampling Distribution

We want to think about the sampling distribution of the estimator.



But remember that we only get to see **one** draw from the sampling distribution. Thus ideally we want an estimator with good **properties**.

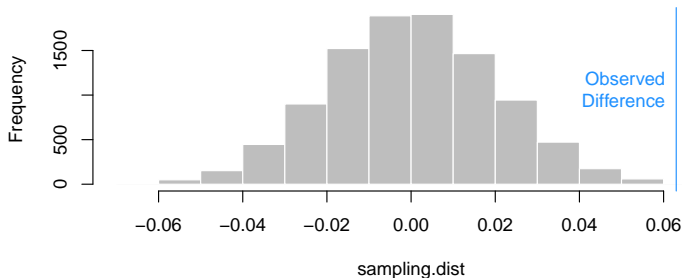
# Asymptotic Normality

Going back to the Gerber, Green, and Larimer result. . .

- The estimator is **difference in means**
- The estimate is 0.063
- Suppose we have an estimate of the estimator's **standard error**  
 $\widehat{SE}(\hat{\theta}) = 0.02$ .
- What if there was **no** difference in means in the population  
( $\mu_y - \mu_x = 0$ )?
- By asymptotic Normality  $(\hat{\theta} - 0)/SE(\hat{\theta}) \sim N(0, 1)$
- By the properties of Normals, we know that this implies that  
 $\hat{\theta} \sim \mathcal{N}(0, SE(\hat{\theta}))$

# Asymptotic Normality

We can plot this to get a feel for it.



Does the observed difference in means seem plausible if there really were no difference between the two groups in the population?

# Summary of Today

- Sampling distributions provide away for studying the properties of sample statistics.
- We must usually make assumptions and/or appeal to a large  $n$  in order to derive a sampling distribution.
- Choosing a point estimator may require tradeoffs between desirable properties.

Next Class: interval estimation

# Summary of Properties

Concept	Criteria	Intuition
Unbiasedness	$E[\hat{\mu}] = \mu$	Right on average
Efficiency	$V[\hat{\mu}_1] < V[\hat{\mu}_2]$	Low variance
Consistency	$\hat{\mu}_n \xrightarrow{P} \mu$	Converge to estimand as $n \rightarrow \infty$
Asymptotic Normality	$\hat{\mu}_n \overset{\text{approx.}}{\sim} N(\mu, \frac{\sigma^2}{n})$	Approximately normal in large $n$

# Fun with Hidden Populations



Dennis M. Feehan and Matthew J. Salganik (2016)  
“Generalizing the Network Scale-Up Method: A New  
Estimator for the Size of Hidden Populations”  
*Sociological Methodology*,  
<http://dx.doi.org/10.1177/0081175016665425>

Slides graciously provided by Matt Salganik.

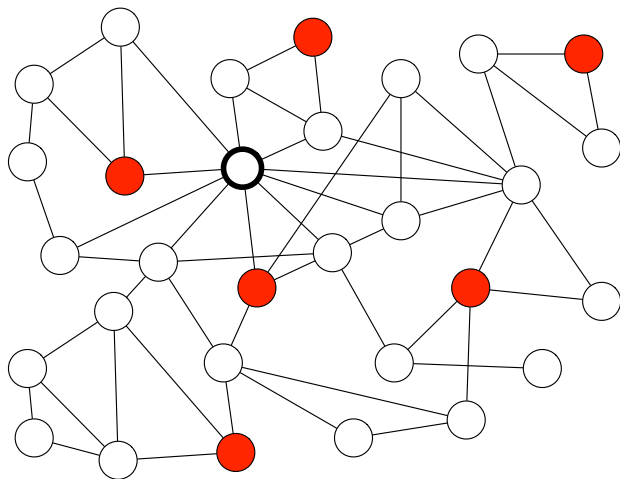


# Scale-up Estimator



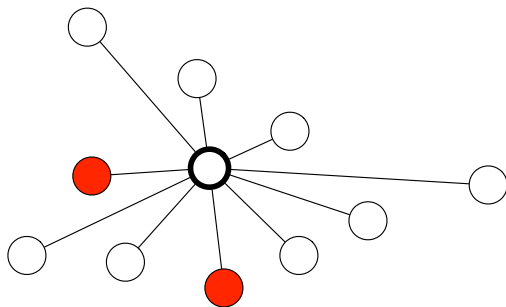
Basic insight from Bernard et al. (1989)

## Network scale-up method



$$\hat{N}_T = \frac{\sum_i y_{i,T}}{\sum_i \hat{d}_i} \times N$$

## Network scale-up method



$$\hat{N}_T = \frac{2}{10} \times 30 = 6$$

If  $\underbrace{y_{i,k} \sim \text{Bin}(d_i, N_k/N)}_{\text{basic scale-up model}}$ , then maximum likelihood estimator is

$$\hat{N}_T = \frac{\sum_i y_{i,T}}{\sum_i \hat{d}_i} \times N$$

- $\hat{N}_T$ : number of people in the target population
- $y_{i,T}$ : number of people in target population known by person  $i$
- $\hat{d}_i$ : estimated number of people known by person  $i$
- $N$ : number of people in the population

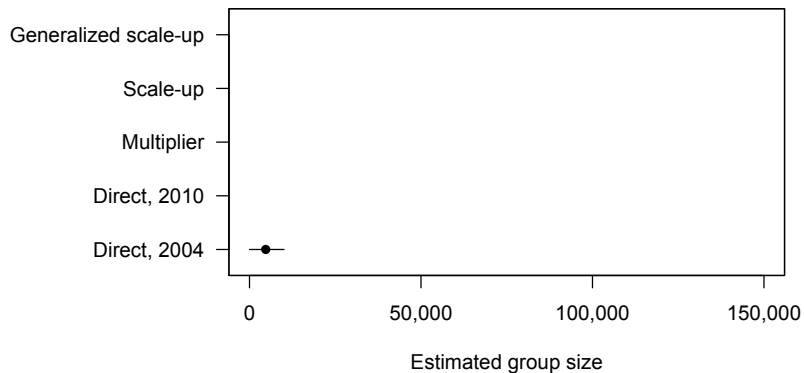
See Killworth et al., (1998)

Target population	Location	Citation
Mortality in earthquake	Mexico City, Mexico	Bernard et al. (1989)
Rape victims	Mexico City, Mexico	Bernard et al. (1991)
HIV prevalence, rape, & homelessness	U.S.	Killworth et al. (1998)
Heroin use	14 U.S. cities	Kadushin et al. (2006)
Choking incidents in children	Italy	Snidero et al. (2007, 2009)
Groups most at-risk for HIV/AIDS	Ukraine	Paniotto et al. (2009)
Heavy drug users	Curitiba, Brazil	Salganik et al. (2011)
Men who have sex with men	Japan	Ezoe et al. (2012)
Groups most at risk for HIV/AIDS	Almaty, Kazakhstan	Scutelnicuic (2012a)
Groups most at risk for HIV/AIDS	Moldova	Scutelnicuic (2012b)
Groups most at risk for HIV/AIDS	Thailand	Aramrattan (2012)
Groups most at risk for HIV/AIDS	Chongqing, China	Guo (2012)
Groups most at risk for HIV/AIDS	Rwanda	Rwanda Biomedical Center (2012)

## Does it Work? Under What Conditions?

- Feehan and Salganik study the properties of the estimator
- They show that for the estimator to be **unbiased** and **consistent** requires a particular assumption that average personal network size is the same in the hidden population as the remainder.
- This was unknown up to this point!
- Analyzing the estimator let them see that the problem can be addressed by collecting a new kind of data on the visibility of hidden population (which can easily be collected with respondent driven sampling)

## Heavy Drug Users, Curitiba, Brazil



## Meta points

- Studying estimators can not only expose problems but suggest solutions
- Another example of creative and interesting ideas coming from the applied people
- Formalizing methods is important because it is what allows them to be studied- it was a long time before anyone discovered the bias/consistency concerns!



## Appendix: More Details on Network scale-up method

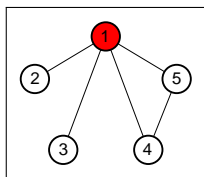
Two important advantages of network scale-up method:

- only requires a random sample of the general population and is therefore easier to standardize across place and time
- built-in validation

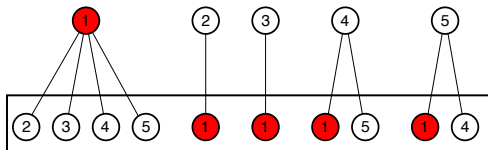
but there are problems too ...

If you have all the data, does it get the right answer?

Issue 1: Set of egos can be different from sequence of alters.



$$p = 0.2$$



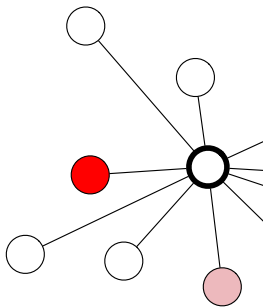
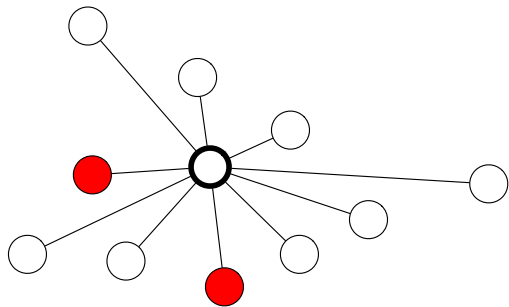
$$p_{alter} = 0.4$$

$$p_{alter} = p \times \frac{\text{avg. degree (target pop.)}}{\text{avg. degree (general pop.)}} = p\delta$$

Estimates will be biased by a factor of  $\delta$  ("degree ratio")

Is sampling the only source of error?

Issue 2: Ego is not aware of everything about all of their alters.



Estimates will be biased by a factor of  $\tau$  ("information transmission rate")

Generalized scale-up estimator:

$$\hat{p} = \frac{\sum_i y_i}{\sum_i \hat{d}_i} \cdot \left(\frac{1}{\hat{\delta}}\right) \cdot \left(\frac{1}{\hat{\tau}}\right)$$

## Game of contacts: Context

We developed the [game of contacts](#) to estimate transmission rate and degree ratio. To estimate the number of [heavy drug users](#) in Curitiba, Brazil (city of 1.8 million people), we did a two-part study:

- 1 “game of contacts” to estimate transmission rate and degree ratio (sample of 294 heavy drug users)
- 2 scale-up survey (sample of 500 people in general population)

Results combined to produce estimates that are compared to estimates from other methods

## Game of contacts

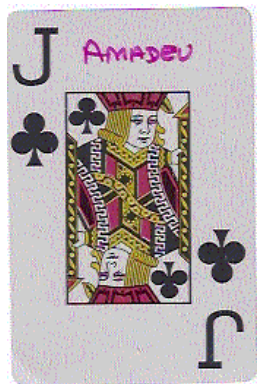
Use a variation of approach from McCarty et al. (1997). Interviewer shuffles a deck of 24 playing cards . . .





## Game of contacts

A card is pulled from the deck and the respondent is asked:



How many people do you know named [Amadeu]?

## Game of contacts

The respondent will pick up this many blocks and place them:



Record answers; clear board; repeated for 24 names.

# Game of contacts: Results

294 participants, 4,173 alters

“selective exposure” and “selective disclosure” (Kitts, 2003)

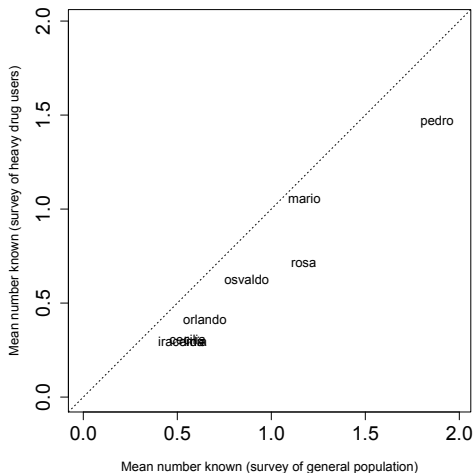
	Use	~Use
Aware		
~Aware		

Transmission rate  $\hat{\tau} = 0.76$ ,  $[0.72, 0.80]$

Other data checks in paper

## Game of contacts: Degree ratio

Ask the same questions in the game of contacts and the scale-up survey (e.g. “How many people do you know named Pedro?”)



## References

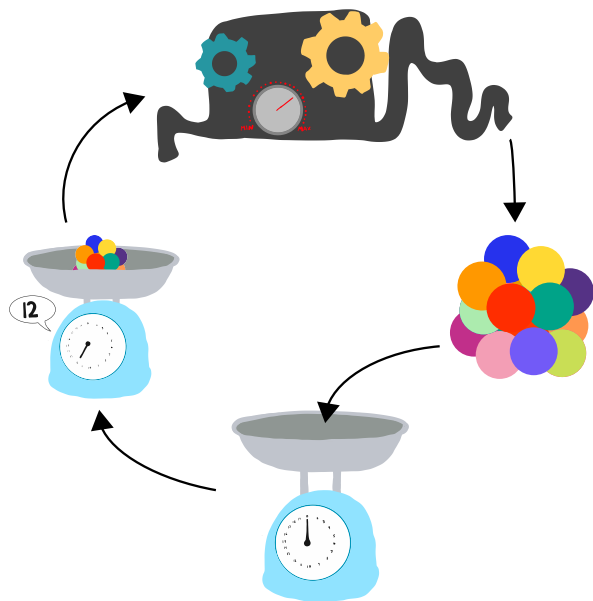
- Kuklinski et al. 1997 “Racial prejudice and attitudes toward affirmative action” *American Journal of Political Science*  
<http://www.jstor.org/stable/2111770>
- Gerber, Green, and Larimer. 2008. “Social pressure and voter turnout: Evidence from a large-scale field experiment.” *American Political Science Review* 102: 33-48.  
<https://doi.org/10.1017/S000305540808009X>.
- Feehan and Salganik 2017 “Generalizing the Network Scale-Up Method: A New Estimator for the Size of Hidden Populations” *Sociological Methodology*,  
<http://dx.doi.org/10.1177/0081175016665425>.

# Where We've Been and Where We're Going...

- Last Week
  - ▶ random variables
  - ▶ joint distributions
- This Week
  - ▶ Monday: Point Estimation
    - ★ sampling and sampling distributions
    - ★ point estimates
    - ★ properties (bias, variance, consistency)
  - ▶ Wednesday: Interval Estimation
    - ★ confidence intervals
    - ★ comparing two groups
- Next Week
  - ▶ hypothesis testing
  - ▶ what is regression?
- Long Run
  - ▶ probability  $\rightarrow$  inference  $\rightarrow$  regression  $\rightarrow$  causal inference

Questions?

# Last Time



- 1 Populations, Sampling, Sampling Distributions
  - Conceptual
  - Mathematical
- 2 Overview of Point Estimation
- 3 Properties of Estimators
- 4 Review and Example
- 5 Fun With Hidden Populations
- 6 Interval Estimation**
- 7 Large Sample Intervals for a Mean
  - Simple Example
  - Kuklinski Example
- 8 Small Sample Intervals for a Mean
- 9 Comparing Two Groups
- 10 Fun With Correlation
- 11 Appendix:  $\chi^2$  and  $t$ -distribution



# What is Interval Estimation?

- A **point estimator**  $\hat{\theta}$  estimates a scalar population parameter  $\theta$  with a single number.
- However, because we are dealing with a random sample, we might also want to report **uncertainty** in our estimate.
- An **interval estimator** for  $\theta$  takes the following form:

$$[\hat{\theta}_{lower}, \hat{\theta}_{upper}]$$

where  $\hat{\theta}_{lower}$  and  $\hat{\theta}_{upper}$  are random quantities that vary from sample to sample.

- The interval represents the range of possible values within which we estimate the true value of  $\theta$  to fall.
- An **interval estimate** is a realized value from an interval estimator. The estimated interval typically forms what we call a **confidence interval**, which we will define shortly.

## Normal Population with Known $\sigma^2$

Suppose we have an i.i.d. random sample of size  $n$ ,  $X_1, \dots, X_n$ , from  $X \sim N(\mu, 1)$ .

From previous lecture, we know that the sampling distribution of the sample average is:

$$\bar{X}_n \sim N(\mu, \sigma^2/n) = N(\mu, 1/n)$$

Therefore, the standardized sample average is distributed as follows:

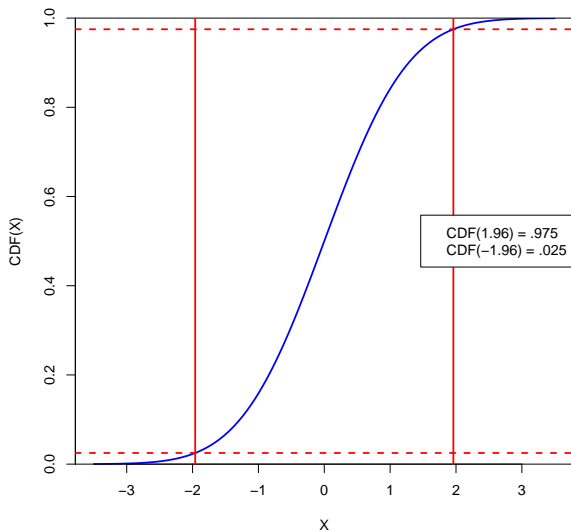
$$\frac{\bar{X}_n - \mu}{1/\sqrt{n}} \sim N(0, 1)$$

This implies

$$\Pr\left(-1.96 < \frac{\bar{X}_n - \mu}{1/\sqrt{n}} < 1.96\right) = .95$$

Why?

# CDF of the Standard Normal Distribution



## Constructing a Confidence Interval with Known $\sigma^2$

So we know that:

$$\Pr\left(-1.96 < \frac{\bar{X}_n - \mu}{1/\sqrt{n}} < 1.96\right) = .95$$

Rearranging yields:

$$\Pr\left(\bar{X}_n - 1.96/\sqrt{n} < \mu < \bar{X}_n + 1.96/\sqrt{n}\right) = .95$$

This implies that the following interval estimator

$$\left[\bar{X}_n - 1.96/\sqrt{n}, \bar{X}_n + 1.96/\sqrt{n}\right]$$

contains the true population mean  $\mu$  with probability 0.95.

We call this estimator a 95% **confidence interval** for  $\mu$ .

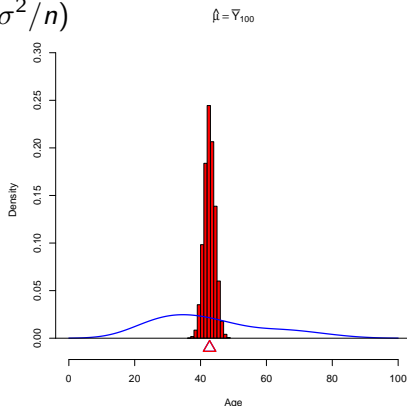
# Kuklinski Example

$$\bar{Y} \sim_{\text{approx}} N(\mu, \sigma^2/n)$$

Suppose the 1,161 respondents in the Kuklinski data set were the population, with  $\mu = 42.7$  and  $\sigma^2 = 257.9$ .

If we sampled 100 respondents, the sampling distribution of  $\bar{Y}_{100}$  is:

$$\bar{Y}_{100} \sim_{\text{approx}} N(42.7, 2.579)$$



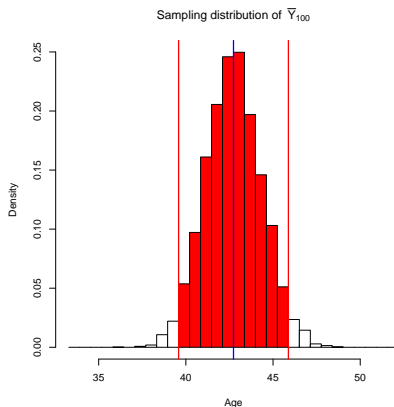
# The standard error of $\bar{Y}$

The **standard error** of the sample mean is the standard deviation of the sampling distribution for  $\bar{Y}$ :

$$SE(\bar{Y}) = \sqrt{V(\bar{Y})} = \frac{\sigma}{\sqrt{n}}$$

What is the probability that  $\bar{Y}$  falls within 1.96 SEs of  $\mu$ ?

But the 1,161 is actually the sample (not the population).



## Normal Population with Unknown $\sigma^2$

In practice, it is rarely the case that we somehow know the true value of  $\sigma^2$ .

Suppose now that we have an i.i.d. random sample of size  $n$   $X_1, \dots, X_n$  from  $X \sim N(\mu, \sigma^2)$ , where  $\sigma^2$  is **unknown**. Then, as before,

$$\bar{X}_n \sim N(\mu, \sigma^2/n) \quad \text{and so} \quad \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \sim N(0, 1).$$

Previously, we then constructed the interval:

$$[\bar{X}_n - z_{\alpha/2}\sigma/\sqrt{n}, \bar{X}_n + z_{\alpha/2}\sigma/\sqrt{n}]$$

But we can **not** directly use this now because  $\sigma^2$  is unknown.

Instead, we need an **estimator** of  $\sigma^2$ ,  $\hat{\sigma}^2$ .

# Estimators for the Population Variance

Two possible estimators of population variance:

$$S_{0n}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

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

Which do we prefer? Let's check properties of these estimators.

- ① Unbiasedness: We can show (after some algebra) that

$$E[S_{0n}^2] = \frac{n-1}{n} \sigma^2 \quad \text{and} \quad E[S_{1n}^2] = \sigma^2$$

- ② Consistency: We can show that

$$S_{0n}^2 \xrightarrow{P} \sigma^2 \quad \text{and} \quad S_{1n}^2 \xrightarrow{P} \sigma^2$$

$S_{1n}^2$  (unbiased and consistent) is commonly called the **sample variance**.



## Estimating $\sigma$ and the SE

Returning to Kulinski et. al. . .

We will use the sample variance:

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

and thus the sample standard deviation can be written as

$$S = \sqrt{S^2}$$

We will plug in  $S$  for  $\sigma$  and our estimated standard error will be

$$\widehat{SE}[\hat{\mu}] = \frac{S}{\sqrt{n}}$$

## 95% Confidence Intervals

If  $X_1, \dots, X_n$  are i.i.d. and  $n$  is large, then

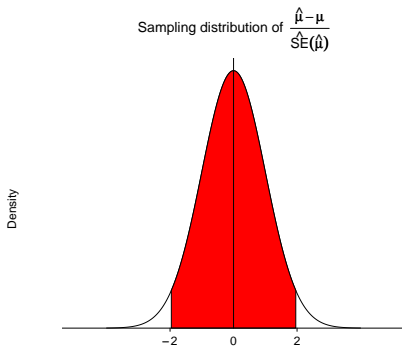
$$\hat{\mu} \sim N(\mu, (\widehat{SE}[\hat{\mu}])^2)$$

$$\hat{\mu} - \mu \sim N(0, (\widehat{SE}[\hat{\mu}])^2)$$

$$\frac{\hat{\mu} - \mu}{\widehat{SE}[\hat{\mu}]} \sim N(0, 1)$$

We know that

$$P\left(-1.96 \leq \frac{\hat{\mu} - \mu}{\widehat{SE}[\hat{\mu}]} \leq 1.96\right) = 95\%$$



## 95% Confidence Intervals

We can work backwards from this:

$$P\left(-1.96 \leq \frac{\hat{\mu} - \mu}{\widehat{SE}[\hat{\mu}]} \leq 1.96\right) = 95\%$$

$$P\left(-1.96\widehat{SE}[\hat{\mu}] \leq \hat{\mu} - \mu \leq 1.96\widehat{SE}[\hat{\mu}]\right) = 95\%$$

$$P\left(\hat{\mu} - 1.96\widehat{SE}[\hat{\mu}] \leq \mu \leq \hat{\mu} + 1.96\widehat{SE}[\hat{\mu}]\right) = 95\%$$

The random quantities in this statement are  $\hat{\mu}$  and  $\widehat{SE}[\hat{\mu}]$ .  
Once the data are observed, nothing is random!

# What does this mean?

We can simulate this process using the Kuklinski data:

- 1) Draw a sample of size 100:

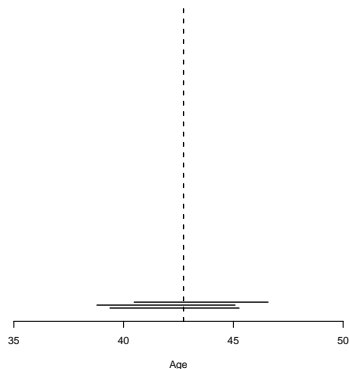


- 2) Calculate  $\hat{\mu}$  and  $\widehat{SE}[\hat{\mu}]$ :

$$\hat{\mu} = 43.53 \quad \widehat{SE}[\hat{\mu}] = 1.555$$

- 3) Construct the 95% CI:

$$(40.5, 46.6)$$

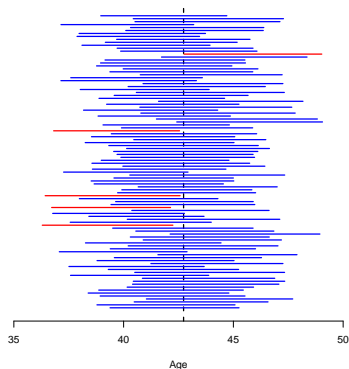


## What does this mean?

By repeating this process, we generate the sampling distribution of the 95% CIs.

Most of the CIs **cover** the true  $\mu$ ; some do not.

In the long run, we expect 95% of the CIs generated to contain the true value.



# Interpreting a Confidence Interval

This can be tricky, so let's break it down.

- Imagine we implement the interval estimator  $\bar{X}_n \pm 1.96/\sqrt{n}$  for a particular sample and obtain the estimate of  $[2.5, 4]$ .
- Does this mean that there is a .95 probability that the true parameter value  $\mu$  lies between these two particular numbers? **No!**
- Confidence intervals are easy to construct, but difficult to interpret:
  - ▶ Each confidence interval estimate from a particular sample either contains  $\mu$  or not
  - ▶ However, if we were to repeatedly calculate the interval estimator over many random samples from the same population, 95% of the time the constructed confidence intervals would cover  $\mu$
  - ▶ Therefore, we refer to .95 as the **coverage probability**

# What makes a good confidence interval?

- 1 The **coverage probability**: how likely it is that the interval covers the truth.
- 2 The **length** of the confidence interval:
  - ▶ Infinite intervals  $(-\infty, \infty)$  have coverage probability 1
  - ▶ For a probability, a confidence interval of  $[0, 1]$  also have coverage probability 1
  - ▶ Zero-length intervals, like  $[\bar{Y}, \bar{Y}]$ , have coverage probability 0
- Best: for a fixed confidence level/coverage probability, find the smallest interval

## Is 95% all there is?

Our 95% CI had the following form:  $\hat{\mu} \pm 1.96 \widehat{SE}[\hat{\mu}]$

Remember where 1.96 came from?

$$P\left(-1.96 \leq \frac{\hat{\mu} - \mu}{\widehat{SE}[\hat{\mu}]} \leq 1.96\right) = 95\%$$

What if we want a different percentage?

$$P\left(-z \leq \frac{\hat{\mu} - \mu}{\widehat{SE}[\hat{\mu}]} \leq z\right) = (1 - \alpha)\%$$

How can we find  $z$ ?

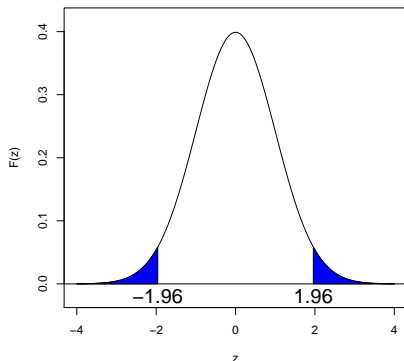


# Normal PDF

We know that  $z$  comes from the probability in the tails of the standard normal distribution.

When  $(1 - \alpha) = 0.95$ , we want to pick  $z$  so that 2.5% of the probability is in each tail.

This gives us a value of 1.96 for  $z$ .

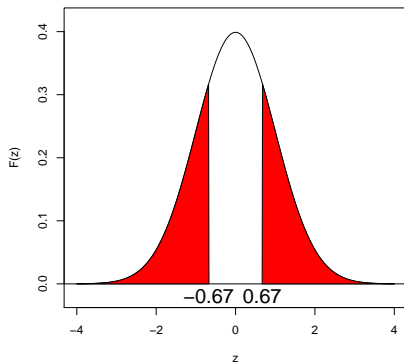


# Normal PDF

What if we want a 50% confidence interval?

When  $(1 - \alpha) = 0.50$ , we want to pick  $z$  so that 25% of the probability is in each tail.

This gives us a value of 0.67 for  $z$ .



## $(1 - \alpha)\%$ Confidence Intervals

In general, let  $z_{\alpha/2}$  be the value associated with  $(1 - \alpha)\%$  coverage:

$$P\left(-z_{\alpha/2} \leq \frac{\hat{\mu} - \mu}{\widehat{SE}[\hat{\mu}]} \leq z_{\alpha/2}\right) = (1 - \alpha)\%$$
$$P\left(\hat{\mu} - z_{\alpha/2}\widehat{SE}[\hat{\mu}] \leq \mu \leq \hat{\mu} + z_{\alpha/2}\widehat{SE}[\hat{\mu}]\right) = (1 - \alpha)\%$$

We usually construct the  $(1 - \alpha)\%$  confidence interval with the following formula.

$$\hat{\mu} \pm z_{\alpha/2}\widehat{SE}[\hat{\mu}]$$

# 1 Populations, Sampling, Sampling Distributions

- Conceptual
- Mathematical

## 2 Overview of Point Estimation

## 3 Properties of Estimators

## 4 Review and Example

## 5 Fun With Hidden Populations

## 6 Interval Estimation

## 7 Large Sample Intervals for a Mean

- Simple Example
- Kuklinski Example

## 8 Small Sample Intervals for a Mean

## 9 Comparing Two Groups

## 10 Fun With Correlation

## 11 Appendix: $\chi^2$ and $t$ -distribution

## The problem with small samples

Up to this point, we have relied on large sample sizes to construct confidence intervals.

If the sample is large enough, then the sampling distribution of the sample mean follows a normal distribution.

If the sample is large enough, then the sample standard deviation ( $S$ ) is a good approximation for the population standard deviation ( $\sigma$ ).

When the sample size is small, we need to know something about the distribution in order to construct confidence intervals with the correct coverage (because we can't appeal to the CLT or assume that  $S$  is a good approximation of  $\sigma$ ).

# Canonical Small Sample Example

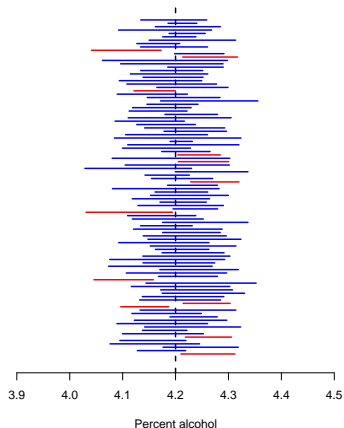
What happens if we use the large-sample formula?

The percent alcohol in Guinness beer is distributed  $N(4.2, 0.09)$ .

Take 100 six-packs of Guinness and construct CIs of the form

$$\hat{\mu} \pm 1.96\widehat{SE}[\hat{\mu}]$$

In this sample, only 88 of the 100 CIs cover the true value.



## The $t$ distribution

If  $X$  is normally distributed, then  $\bar{X}$  is normally distributed even in small samples. Assume

$$X \sim N(\mu, \sigma^2)$$

If we know  $\sigma$ , then

$$\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$$

We rarely know  $\sigma$  and have to use an estimate instead:

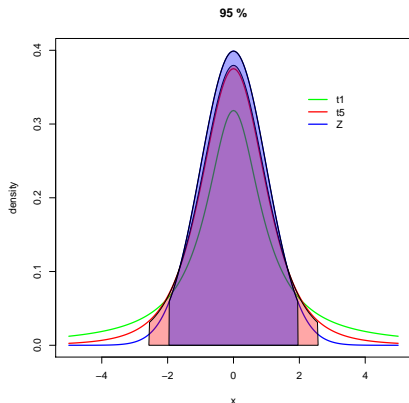
$$\frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} \sim ??t_{n-1}$$

# The $t$ distribution

Since we have to estimate  $\sigma$ , the distribution of  $\frac{\bar{X}-\mu}{\frac{s}{\sqrt{n}}}$  is still bell-shaped but is more spread out.

As the sample size increases, our estimates of  $\sigma$  improve and extreme values of  $\frac{\bar{X}-\mu}{\frac{s}{\sqrt{n}}}$  become less likely.

Eventually the  $t$  distribution converges to the standard normal.





## $(1 - \alpha)\%$ Confidence Intervals

In general, let  $t_{\alpha/2}$  be the value associated with  $(1 - \alpha)\%$  coverage:

$$P\left(-t_{\alpha/2} \leq \frac{\hat{\mu} - \mu}{\widehat{SE}[\hat{\mu}]} \leq t_{\alpha/2}\right) = (1 - \alpha)\%$$
$$P\left(\hat{\mu} - t_{\alpha/2}\widehat{SE}[\hat{\mu}] \leq \mu \leq \hat{\mu} + t_{\alpha/2}\widehat{SE}[\hat{\mu}]\right) = (1 - \alpha)\%$$

We usually construct the  $(1 - \alpha)\%$  confidence interval with the following formula.

$$\hat{\mu} \pm t_{\alpha/2}\widehat{SE}[\hat{\mu}]$$

## Small Sample Example

When we generated 95% CIs with the large sample formula

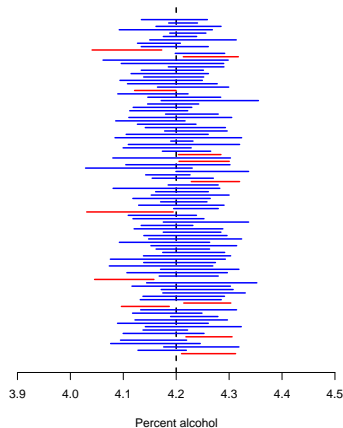
$$\hat{\mu} \pm 1.96\widehat{SE}[\hat{\mu}]$$

only 88 out of 100 intervals covered the true value.

When we use the correct small-sample formula

$$\hat{\mu} \pm t_{\alpha/2}\widehat{SE}[\hat{\mu}]2.57\widehat{SE}[\hat{\mu}]$$

95 of the 100 CIs in this sample cover the truth.



## Another Rationale for the $t$ -Distribution

Does  $\bar{X}_n \sim N(\mu, S_n^2/n)$ , which would imply  $\frac{\bar{X}_n - \mu}{S_n/\sqrt{n}} \sim N(0, 1)$ ?

No, because  $S_n$  is a random variable instead of a parameter (like  $\sigma$ ).

Thus, we need to derive the sampling distribution of the new random variable. It turns out that  $T_n$  follows **Student's  $t$ -distribution** with  $n - 1$  **degrees of freedom**.

### Theorem (Distribution of $t$ -Value from a Normal Population)

*Suppose we have an i.i.d. random sample of size  $n$  from  $N(\mu, \sigma^2)$ . Then, the sample mean  $\bar{X}_n$  standardized with the estimated standard error  $S_n/\sqrt{n}$  satisfies,*

$$T_n \equiv \frac{\bar{X}_n - \mu}{S_n/\sqrt{n}} \sim \tau_{n-1}$$

► Appendix

- 1 Populations, Sampling, Sampling Distributions
  - Conceptual
  - Mathematical
- 2 Overview of Point Estimation
- 3 Properties of Estimators
- 4 Review and Example
- 5 Fun With Hidden Populations
- 6 Interval Estimation
- 7 Large Sample Intervals for a Mean
  - Simple Example
  - Kuklinski Example
- 8 Small Sample Intervals for a Mean
- 9 Comparing Two Groups**
- 10 Fun With Correlation
- 11 Appendix:  $\chi^2$  and  $t$ -distribution

# Kuklinski Example Returns

The Kuklinski et al. (1997) article compares responses to the baseline list with responses to the treatment list.

- How should we estimate the difference between the two groups?
- How should we obtain a confidence interval for our estimate?

# Comparing Two Groups

We will often assume the following when comparing two groups,

- $X_{11}, X_{12}, \dots, X_{1n_1} \sim i.i.d.?(μ_1, σ_1^2)$
- $X_{21}, X_{22}, \dots, X_{2n_2} \sim i.i.d.?(μ_2, σ_2^2)$
- The two samples are independent of each other.

We will usually be interested in comparing  $μ_1$  to  $μ_2$ , although we will sometimes need to compare  $σ_1^2$  to  $σ_2^2$  in order to make the first comparison.

## Sampling Distribution for $\bar{X}_1 - \bar{X}_2$

What is the expected value of  $\bar{X}_1 - \bar{X}_2$ ?

$$\begin{aligned} E[\bar{X}_1 - \bar{X}_2] &= E[\bar{X}_1] - E[\bar{X}_2] \\ &= \frac{1}{n_1} \sum E[X_{1i}] - \frac{1}{n_2} \sum E[X_{2j}] \\ &= \frac{1}{n_1} \sum \mu_1 - \frac{1}{n_2} \sum \mu_2 \\ &= \mu_1 - \mu_2 \end{aligned}$$

## Sampling Distribution for $\bar{X}_1 - \bar{X}_2$

What is the variance of  $\bar{X}_1 - \bar{X}_2$ ?

$$\begin{aligned}\text{Var}[\bar{X}_1 - \bar{X}_2] &= \text{Var}[\bar{X}_1] + \text{Var}[\bar{X}_2] \\ &= \frac{1}{n_1^2} \sum \text{Var}[X_{1i}] + \frac{1}{n_2^2} \sum \text{Var}[X_{2j}] \\ &= \frac{1}{n_1^2} \sum \sigma_1^2 + \frac{1}{n_2^2} \sum \sigma_2^2 \\ &= \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}\end{aligned}$$



## Sampling Distribution for $\bar{X}_1 - \bar{X}_2$

What is the distributional form for  $\bar{X}_1 - \bar{X}_2$ ?

- $\bar{X}_1$  is distributed  $\sim N(\mu_1, \frac{\sigma_1^2}{n_1})$ .
- $\bar{X}_2$  is distributed  $\sim N(\mu_2, \frac{\sigma_2^2}{n_2})$ .
- $\bar{X}_1 - \bar{X}_2$  is distributed  $\sim N(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2})$ .

## CIs for $\mu_1 - \mu_2$

Using the same type of argument that we used for the univariate case, we write a  $(1 - \alpha)\%$  CI as the following:

$$\bar{X}_1 - \bar{X}_2 \pm z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

## Interval estimation of the population proportion

- Let's say that we have a sample of iid Bernoulli random variables,  $Y_1, \dots, Y_n$ , where each takes  $Y_i = 1$  with probability  $\pi$ . Note that this is also the **population proportion** of ones. We have shown in previous weeks that the expectation of one of these variable is just the probability of seeing a 1:  $E[Y_i] = \pi$ .
- The **variance of a Bernoulli random variable** is a simple function of its mean:  $\text{Var}(Y_i) = \pi(1 - \pi)$ .
- **Problem** Show that the sample proportion,  $\hat{\pi} = \frac{1}{n} \sum_{i=1}^n Y_i$ , of the above iid Bernoulli sample, is unbiased for the true population proportion,  $\pi$ , and that the sampling variance is equal to  $\frac{\pi(1-\pi)}{n}$ .
- Note that if we have an estimate of the population proportion,  $\hat{\pi}$ , then we also have an estimate of the sampling variance:  $\frac{\hat{\pi}(1-\hat{\pi})}{n}$ .
- Given the facts from the previous problem, we just apply the same logic from the population mean to show the following confidence interval:

$$P \left( \hat{\pi} - z_{\alpha/2} \times \sqrt{\frac{\hat{\pi}(1 - \hat{\pi})}{n}} \leq \pi \leq \hat{\pi} + z_{\alpha/2} \times \sqrt{\frac{\hat{\pi}(1 - \hat{\pi})}{n}} \right) = (1 - \alpha)$$

## Gerber, Green, and Larimer experiment

Let's go back to the Gerber, Green, and Larimer experiment from last class. Here are the results of their experiment:

**TABLE 2. Effects of Four Mail Treatments on Voter Turnout in the August 2006 Primary Election**

	Experimental Group				
	Control	Civic Duty	Hawthorne	Self	Neighbors
Percentage Voting	29.7%	31.5%	32.2%	34.5%	37.8%
N of Individuals	191,243	38,218	38,204	38,218	38,201

- Let's use what we have learned up until now and the information in the table to calculate a 95% confidence interval for the difference in proportions voting between the Neighbors group and the Civic Duty group.
- You may assume that the samples within each group are iid and the two samples are independent.

## Calculating the CI for social pressure effect

- We know distribution of sample proportion turned among Civic Duty group  $\hat{\pi}_C \sim N(\pi_C, (\pi_C(1 - \pi_C))/n_C)$
- Sample proportions are just sample means, so we can do difference in means:

$$\hat{\pi}_N - \hat{\pi}_C \sim N\left(\pi_N - \pi_C, \sqrt{SE_N^2 + SE_C^2}\right)$$

- Replace the variances with our estimates:

$$\hat{\pi}_N - \hat{\pi}_C \sim N\left(\pi_N - \pi_C, \sqrt{\widehat{SE}_N^2 + \widehat{SE}_C^2}\right)$$

- Apply usual formula to get 95% confidence interval:

$$(\hat{\pi}_N - \hat{\pi}_C) \pm 1.96 \times \sqrt{\widehat{SE}_N^2 + \widehat{SE}_C^2}$$

- Remember that we can calculate the sample variance for a sample proportion like so:  $(\hat{\pi}_C(1 - \hat{\pi}_C))/n_C$

## Gerber, Green, and Larimer experiment

**TABLE 2. Effects of Four Mail Treatments on Voter Turnout in the August 2006 Primary Election**

	Experimental Group				
	Control	Civic Duty	Hawthorne	Self	Neighbors
Percentage Voting	29.7%	31.5%	32.2%	34.5%	37.8%
N of Individuals	191,243	38,218	38,204	38,218	38,201

Now, we calculate the 95% confidence interval:

$$(\hat{\pi}_N - \hat{\pi}_C) \pm 1.96 \times \sqrt{\frac{\hat{\pi}_N(1 - \hat{\pi}_N)}{n_N} + \frac{\hat{\pi}_C(1 - \hat{\pi}_C)}{n_C}}$$

```
n.n <- 38201
samp.var.n <- (0.378 * (1 - 0.378))/n.n
n.c <- 38218
samp.var.c <- (0.315 * (1 - 0.315))/n.c
se.diff <- sqrt(samp.var.n + samp.var.c)
## lower bound
(0.378 - 0.315) - 1.96 * se.diff
## [1] 0.05626701
## upper bound
(0.378 - 0.315) + 1.96 * se.diff
## [1] 0.06973299
```

Thus, the confidence interval for the effect is [0.056267, 0.069733].

# Summary of Interval Estimation

- Interval estimates provide a means of assessing uncertainty.
- Interval estimators have sampling distributions.
- Interval estimates should be interpreted in terms of repeated sampling.

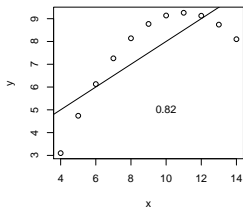
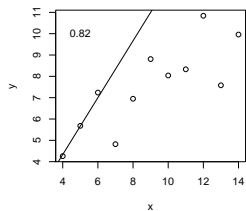
## Next Week

- Hypothesis testing
- What is regression?
- Reading
  - ▶ Aronow and Miller 3.4.2 (testing)
  - ▶ Aronow and Miller 4.1.1 (bivariate regression)
  - ▶ “Momentous Sprint at the 2156 Olympics” by Andrew J Tatem et al. *Nature* 2004
  - ▶ Optional: Imai Ch 2

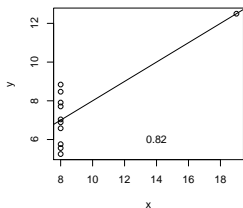
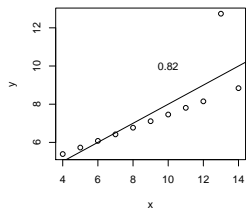


- 1 Populations, Sampling, Sampling Distributions
  - Conceptual
  - Mathematical
- 2 Overview of Point Estimation
- 3 Properties of Estimators
- 4 Review and Example
- 5 Fun With Hidden Populations
- 6 Interval Estimation
- 7 Large Sample Intervals for a Mean
  - Simple Example
  - Kuklinski Example
- 8 Small Sample Intervals for a Mean
- 9 Comparing Two Groups
- 10 Fun With Correlation**
- 11 Appendix:  $\chi^2$  and  $t$ -distribution

# Fun with Anscombe's Quartet

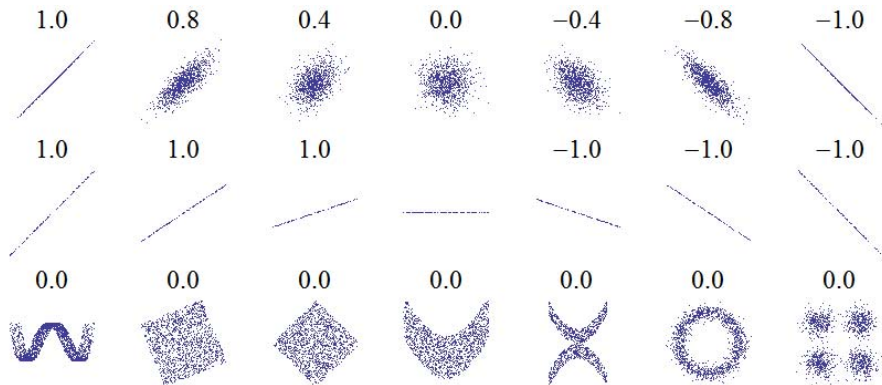


All yield same regression model!



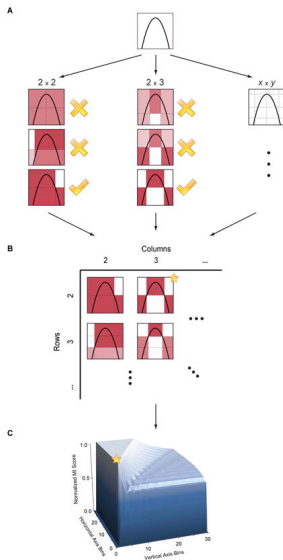
# Fun Beyond Correlation

Reshef, D. N.; Reshef, Y. A.; Finucane, H. K.; Grossman, S. R.; McVean, G.; Turnbaugh, P. J.; Lander, E. S.; Mitzenmacher, M.; Sabeti, P. C. (2011). "Detecting Novel Associations in Large Data Sets". *Science* 334 (6062): 1518-1524.



# Fun with Correlation Part 2

Enter the **Maximal Information Coefficient**



## Fun with Correlation Part 2

### Concerns with MIC

- low power
- originality?
- heuristic binning mechanism
- issues with equitability criterion

This is still an open issue!

- 1 Populations, Sampling, Sampling Distributions
  - Conceptual
  - Mathematical
- 2 Overview of Point Estimation
- 3 Properties of Estimators
- 4 Review and Example
- 5 Fun With Hidden Populations
- 6 Interval Estimation
- 7 Large Sample Intervals for a Mean
  - Simple Example
  - Kuklinski Example
- 8 Small Sample Intervals for a Mean
- 9 Comparing Two Groups
- 10 Fun With Correlation
- 11 Appendix:  $\chi^2$  and  $t$ -distribution

# A Sketch of why the Student $t$ -distribution

▶ Back

- We have need the distribution of  $\frac{\bar{X}_n - \mu}{S_n / \sqrt{n}}$
- $S_n = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$  is now a random variable
- Because we already derived the sampling distribution for  $\frac{\bar{X}_n - \mu}{\sigma^2 / \sqrt{n}}$ , we want to derive the sampling distribution for  $\frac{S_n}{\sigma^2}$  because the  $\sigma^2$  term will cancel.
- Some math will show our distribution is going to be of the form  $\sum Z^2$  where  $Z \sim N(0, 1)$ .
- Let's figure out what distribution that will be

## $\chi^2$ Distribution

Suppose  $Z \sim \text{Normal}(0, 1)$ .

Consider  $X = Z^2$

$$\begin{aligned}F_X(x) &= P(X \leq x) \\&= P(Z^2 \leq x) \\&= P(-\sqrt{x} \leq Z \leq \sqrt{x}) \\&= \frac{1}{\sqrt{2\pi}} \int_{-\sqrt{x}}^{\sqrt{x}} e^{-\frac{z^2}{2}} dt \\&= F_Z(\sqrt{x}) - F_Z(-\sqrt{x})\end{aligned}$$

The pdf then is

$$\frac{\partial F_X(x)}{\partial x} = f_Z(\sqrt{x}) \frac{1}{2\sqrt{x}} + f_Z(-\sqrt{x}) \frac{1}{2\sqrt{x}}$$



## Definition

Suppose  $X$  is a continuous random variable with  $X \geq 0$ , with pdf

$$f(x) = \frac{1}{2^{n/2}\Gamma(n/2)} x^{n/2-1} e^{-x/2}$$

Then we will say  $X$  is a  $\chi^2$  distribution with  $n$  degrees of freedom. Equivalently,

$$X \sim \chi^2(n)$$

## $\chi^2$ Properties

Suppose  $X \sim \chi^2(n)$

$$E[X] = E \left[ \sum_{i=1}^N Z_i^2 \right]$$

$$= \sum_{i=1}^N E[Z_i^2]$$

$$\text{var}(Z_i) = E[Z_i^2] - E[Z_i]^2$$

$$1 = E[Z_i^2] - 0$$

$$E[X] = N$$

## $\chi^2$ Properties

Suppose  $X \sim \chi^2(n)$

$$\begin{aligned}\text{var}(X) &= \sum_{i=1}^N \text{var}(Z_i^2) \\ &= \sum_{i=1}^N (E[Z_i^4] - E[Z_i^2]^2) \\ &= \sum_{i=1}^N (3 - 1) = 2N\end{aligned}$$

# Student's $t$ -Distribution

## Definition

Suppose  $Z \sim \text{Normal}(0, 1)$  and  $U \sim \chi^2(n)$ . Define the random variable  $Y$  as,

$$Y = \frac{Z}{\sqrt{\frac{U}{n}}}$$

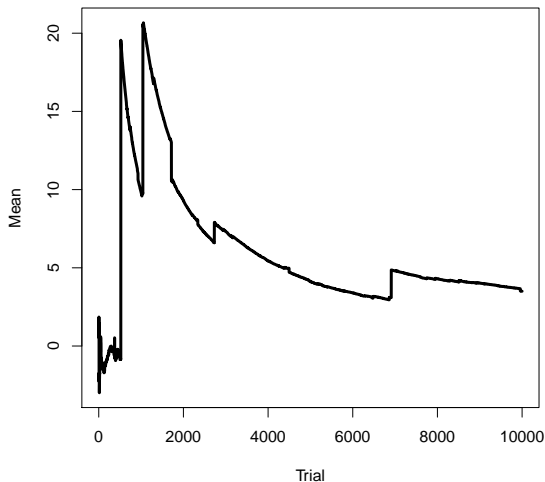
If  $Z$  and  $U$  are independent then  $Y \sim t(n)$ , with pdf

$$f(x) = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{\pi n} \Gamma(\frac{n}{2})} \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}}$$

We will use the  $t$ -distribution extensively for **test-statistics**

# Student's $t$ -Distribution, Properties

Suppose  $n = 1$ , **Cauchy** distribution



If  $X \sim \text{Cauchy}(1)$ , then:

# Student's $t$ -Distribution, Properties

Suppose  $n > 2$ , then

$$\text{var}(X) = \frac{n}{n-2}$$

As  $n \rightarrow \infty$   $\text{var}(X) \rightarrow 1$ .