

# Motivating Example

---

# ROASTING COFFEE: INTRODUCTION

# ROASTING COFFEE: FIELD TRIP

# **Statistics, Parameters, and Estimators**

---

# Introduction to Properties of Estimators

---

# ESTIMATING CHARACTERISTICS OF POPULATIONS

- Data scientists are not given the joint density distribution.
- One of our principle tasks is to produce *estimates* about the *universe* given only the limited information that we get from a sample of data.

# ESTIMATING CHARACTERISTICS OF POPULATIONS (CONT.)

## Parameters, Statistics, and Estimators

- A **parameter**,  $\theta$ , is some summary of the joint distribution.
- A **statistic** is some function that maps data from  $\mathbb{R}^n \rightarrow \mathbb{R}$ .
- An **estimator**,  $\hat{\theta}$ , is some statistic of the data that we use to produce a guess for  $\theta$ .

# $\bar{X}$ IS AN ESTIMATOR

We refer to statistics that are “good” guesses for a population parameter as **estimators**.

- $\theta = E[X]$  is a parameter.
- $T_{(n)} =$  choose the 30<sup>th</sup> percentile is a statistic
- $\hat{\theta} = \bar{X}$  is an estimator for  $\theta$ .



# INTRODUCING ESTIMATORS

In this section, you will learn about the properties that make estimators good guesses for a population parameter.

## Properties of estimators

- **Consistency:** As sample size grows, does the estimator converge in probability to the true value?
- **Bias** and **unbiasedness:** Is the estimator systematically too low or too high?
- **Efficiency:** Does an estimator have relatively *large* or *small* sampling variance, standard error, and MSE?

# **Reading: Core Estimation Theory**

---

## READING: CORE ESTIMATION THEORY

**Note: This is a READING CALL, just placing it here for organization.** Read pages 102–105, stopping at 3.2.3, Variance Estimators.

# **Desirable Properties of Estimators**

---

# DESIRABLE PROPERTIES OF ESTIMATORS

Throughout, we refer to  $\theta$  as a population *feature* or *parameter*.

- $\theta$  has a fixed, true value, but this value is not directly observable.
- Without ground truth, how can we know if our estimator is doing its job?

# DESIRABLE PROPERTIES OF ESTIMATORS: UNBIASEDNESS

## Bias of an estimator

The bias of an estimator is the expected difference between the *true* population value and the estimator.

- The *bias* of  $\hat{\theta}$  is  $E[\hat{\theta}] - \theta$ .
- If  $E[\hat{\theta}] = \theta$ , then there is no bias, and the estimator is *unbiased*.

## DESIRABLE PROPERTIES OF ESTIMATORS: UNBIASEDNESS (CONT.)

If  $\hat{\theta}$  is unbiased, does that mean that, for any sample taken, the estimate produced by  $\hat{\theta} = \theta$ ?

- $\hat{\theta}$  as an *estimator* is a random variable.
- The value that it takes on given a sample is the *estimate*.

## **DESIRABLE PROPERTIES OF ESTIMATORS: UNBIASEDNESS (CONT.)**



# DESIRABLE PROPERTIES OF ESTIMATORS: EFFICIENCY

## Mean Squared Error of an Estimator

- The MSE of an estimator  $\hat{\theta}$  is

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

# DESIRABLE PROPERTIES OF ESTIMATORS: EFFICIENCY

# DESIRABLE PROPERTIES OF ESTIMATORS: CONSISTENCY

## Consistency

An estimator  $\hat{\theta}$  is consistent for  $\theta$  if  $\hat{\theta} \xrightarrow{P} \theta$ .

# DESIRABLE PROPERTIES OF ESTIMATORS: CONSISTENCY

# **Applying Properties of Estimators**

---

# EVALUATING ESTIMATORS

Which do you prefer?

## EVALUATING ESTIMATORS (CONT.)

## DESIRABLE PROPERTIES OF ESTIMATORS: REVIEW

- An *inconsistent* estimator is of little use.
- A *more efficient* estimator is preferable to a *less efficient* estimator, all else equal.
- An *unbiased* estimator is preferable to a *biased* estimator, all else equal.



# Random Sampling

---

# **Reading: Independent and Identically Distributed**

---

## READING: INDEPENDENT AND IDENTICALLY DISTRIBUTED

- Read Section 3.0 and 3.1 in *Foundations of Agnostic Statistics* (pages 91–96 in our copy of the book.)
  - Work to place the formal math definition into plain understanding. We will discuss this when you come back.
  - Notice how  $\mu$  and  $\sigma^2$  are concepts that are immediately familiar, but they now represent population parameters.

# **Independent and Identically Distributed**

---

# INDEPENDENT AND IDENTICALLY DISTRIBUTED

## Definition: independent and identically distributed (IID)

- Undertake some process an arbitrary number of times—call it *sampling*—to create a value from a phenomenon.
- If each instance of sampling draws from the same probability distribution, then we say the collection of values are **identically distributed**.
- If none of the instances of sampling provide information about other instances of sampling, then we say the collection of values are **independent**.

# INDEPENDENT AND IDENTICALLY DISTRIBUTED (CONT.)

# MAKING ASSUMPTIONS

- What is a statistical assumption?
- When are they important?
- When can I violate them?
- What do I do if they are violated?

# MAKING PREDICTIONS



# MAKING PREDICTIONS

- When data are independent and identically distributed (IID), it is possible to learn desirable things from a sample:
  - Accurate characterizations of the probability distribution function and values
  - Reliable characterizations of certainty and uncertainty
- We can learn about unseen *population parameters* from a sample and then use our domain knowledge to evaluate whether these population parameters apply to some circumstance.

# WHAT IS A POPULATION?

## U.S. Senator fundraising

There are 100 U.S. Senators.

- One could reason about a finite population with 100 elements.
- Or, one could assume a continuous distribution for fundraising.

## **Learnosity: Is This IID?**

---

# Is This IID?

**Note: This is a learnosity activity, just placing it here for organization.**

For each of these, would you say that this sample is distributed IID? What do you understand about the process that brings you to your conclusion?

- Coffee
  - **Goal:** Understand how roasted a batch of coffee is.
  - **Sampling process:** Dip into a drum to pull a set of 30 beans.
- Conduct a draft
  - **Goal:** Randomly select people for military service.
  - **Sampling process:** Draw balls from an urn with a day of the year. Everybody born on that day goes to war.
- Produce training data for machine vision

## IS THIS IID? (CONT.)

- Teach a computer to *really* read
  - **Goal:** Teach a computer to understand theme and plot (the work of School of Information professor David Bamman).
  - **Sampling process:** Feed a neural network each word (in each sentence [in each paragraph (in each chapter)]) of the English language literary canon.
- Any others?
  - **Goal:**
  - **Sampling process:**

# Reading: Sample Statistics

---

## READING: SAMPLE STATISTICS

- Read pages 96–98.
- Stop before theorem 3.2.5.

# Sample Statistics

---



# SAMPLE STATISTICS

Sample statistics are:

- Functions that are applied to samples of random variables
- *Themselves* random variables

# Conduct Sampling

---

## CONDUCT SAMPLING

**Note: This is a learnosity activity, just placing it here for organization.**

The goal is that this will solidify the understanding that students have from the reading and will move us forward into a conversation about expected value and the sample variance of the sample average.

- Students will be provided with data and starter code that produces a population.
- From this population, they will have sliders they can pull that will increase or decrease the number of samples that they take, the number of draws per sample, and the underlying population variation.

# **The Sample Mean as an Estimator**

---

# THE SAMPLE MEAN

## Definition: Sample Mean

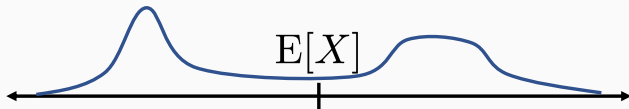
For IID random variables,  $X_1, ..X_n$ , the *sample mean* is

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

- The sample mean is an *estimator* for  $E[X]$ .

# THE SAMPLE MEAN IS AN ESTIMATOR

1. Population Distribution



2. Sample



3. Sampling Distribution of the Statistic



# EVALUATING THE SAMPLE MEAN AS AN ESTIMATOR

First Questions:

1. Is the sample mean biased?
2. How efficient is the sample mean?

# UNBIASEDNESS OF THE SAMPLE MEAN

## Theorem: The Sample Mean is Unbiased

For IID random variables,  $X_1, \dots, X_n$ , the sample mean  $\bar{X}$  is an unbiased estimator for the population mean  $E[X]$ .



# EFFICIENCY OF THE SAMPLE MEAN

## **Theorem: Sampling Variance of the Sample Mean**

For IID random variables,  $X_1, \dots, X_n$ , with population variance,  $V[X]$ , the variance of the sample mean is

$$V[\bar{X}] = \frac{V[X]}{n}$$

# LEARNOSITY: SAMPLING VARIANCE OF THE SAMPLE MEAN

**Note: This is a learnosity activity, just placing it here for organization.**

- If you increase the sample size, does the sampling variance of the sample mean increase, decrease, or stay the same?
- At what rate does this change? (Options include faster than the data, at the same rate as the data, slower than the data.)
- If you increase the sample size, does the population variance,  $V[X]$ , change? (Hint: It is a population parameter.)

# Consistency and the Continuous Mapping Theorem

---

# CONSISTENCY

*"If you can't get it right as  $n$  goes to infinity, you shouldn't be in this business."*

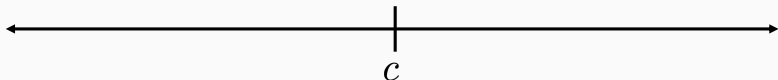
- Clive W.J. Granger

How can we formalize "right as  $n$  goes to infinity?"

- Estimators may converge at different rates.
- Deterministic guarantees are not possible.

# INTUITION FOR CONVERGENCE IN PROBABILITY

- Let  $T_{(1)}$  be the statistic with 1 datapoint.
- Let  $T_{(2)}$  be the statistic with 2 datapoints.
- Let  $T_{(3)}$  be the statistic with 3 datapoints.
- $\vdots$



# CONVERGENCE IN PROBABILITY

## Definition: Convergence in Probability

Let  $(T_{(1)}, T_{(2)}, T_{(3)}, \dots)$  be a sequence of random variables and let  $c \in \mathbb{R}$ .  $T_{(n)}$  *converges in probability* to  $c$  if for all  $\epsilon > 0$ ,

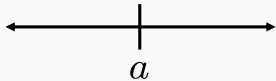
$$\lim_{n \rightarrow \infty} P\left[T_{(n)} \in (c - \epsilon, c + \epsilon)\right] = 1$$

We write this as  $T_{(n)} \xrightarrow{p} c$ .

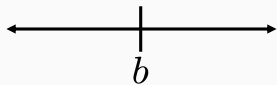
- An estimator  $\hat{\theta}$  is *consistent* for  $\theta$ , if  $\hat{\theta} \xrightarrow{p} \theta$ .

# CONTINUOUS MAPPING INTUITION

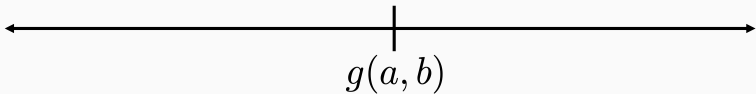
$S$



$T$



$g(S, T)$



## The Continuous Mapping Theorem

Let  $(S_{(1)}, S_{(2)}, S_{(3)}, \dots)$  and  $(T_{(1)}, T_{(2)}, T_{(3)}, \dots)$  be two sequences of random variables. Let  $g : \mathbb{R}^2 \rightarrow \mathbb{R}$  be a continuous function. If  $S_{(n)} \xrightarrow{p} a \in \mathbb{R}$  and  $T_{(n)} \xrightarrow{p} b \in \mathbb{R}$ , then  $g(S_{(n)}, T_{(n)}) \xrightarrow{p} g(a, b)$



# **Reading: Weak Law of Large Numbers**

---

## READING: WEAK LAW OF LARGE NUMBERS

Read page 100, beginning at theorem 3.2.8, through the end of page 102.

# **Weak Law of Large Numbers**

---

# THE WEAK LAW OF LARGE NUMBERS

## Theorem: The Weak Law of Large Numbers

Let  $(X_1, X_2, X_3, \dots)$  be a sequence of i.i.d. random variables with finite variance. Let  $\bar{X}_{(n)} = \frac{1}{n} \sum_{i=1}^n X_i$ . Then

$$\bar{X}_{(n)} \xrightarrow{p} E[X]$$

- **Equivalently:** The sample mean is *consistent* for the population mean.

# CONSEQUENCES OF WLLN

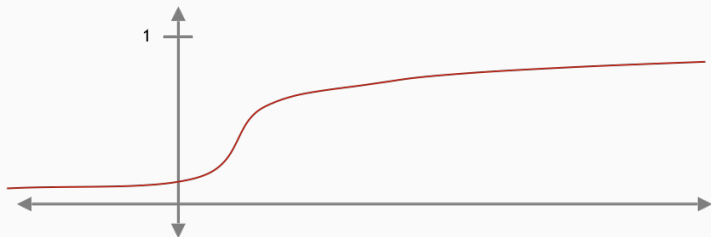
## **The sample mean is consistent.**

- The error converges in probability to zero.
- The sample mean is accurate, as long as we can make  $n$  large enough.

## **As a building block to study other estimators.**

- Combine WLLN with Continuous Mapping Theorem.

# APPLYING THE WLLN



**Objective:** Estimate the cdf  $F$  at a point  $x$ .

**Idea:** Use empirical cdf

# **Proof of the Weak Law of Large Numbers**

---

# PROOF OF THE WEAK LAW OF LARGE NUMBERS

Given random variable  $X$  and  $\epsilon > 0$ .



# PROOF OF THE WEAK LAW OF LARGE NUMBERS

Given random variable  $X$  and  $\epsilon > 0$ .

Let  $D = |\bar{X} - E[X]|$

$$V[\bar{X}] = E[D^2] = V[X]/n$$

$$\begin{aligned} V[\bar{X}] &= E[D^2] = E[D^2|D < \epsilon]P(D < \epsilon) + E[D^2|D \geq \epsilon]P(D \geq \epsilon) \\ &\geq 0 + \epsilon^2 P(D \geq \epsilon) \end{aligned}$$

$$P(D \geq \epsilon) \leq \frac{V[X]}{\epsilon^2 n} \xrightarrow{p} 0$$

# Simulating the WLLN

---

**Note: This is a learnosity activity, we're just including it here for organization.** Students will work through the notebook called WLLN.Rmd.

# **Reading: Estimating Population Variance**

---

## READING: ESTIMATING POPULATION VARIANCE

- Read pages 105–108, stopping before the beginning of the next section.
- You are going to use *the plug-in principle*, which is a general approach for designing estimators in a sample.

# Estimating Population Variance

---

# WHY ESTIMATE VARIANCE?



- How much coffee do students drink on average?
- How much does coffee intake differ from student to student?

# APPLYING THE PLUG-IN PRINCIPLE

Population Variance:  $V[X] = E[X^2] - E[X]^2$

## The Plug-In Variance Estimator

Given a sample  $\mathbf{X} = (X_1, X_2, \dots, X_n)$ , the *plug-in variance estimator* is,

$$\tilde{V}(\mathbf{X}) = \overline{X^2} - \bar{X}^2$$



## CONSISTENCY OF THE PLUG-IN ESTIMATOR

### **Theorem: Consistency of the Plug-In Variance Estimator**

Let  $(X_1, X_2, X_3, \dots)$  be a sequence of i.i.d. random variables with finite variance  $V[X]$ . Then  $\tilde{V}(\mathbf{X}) = \overline{X^2} - \overline{X}^2$  is consistent for  $V[X]$ .

## BIAS OF THE PLUG-IN VARIANCE ESTIMATOR

Let  $(X_1, X_2, X_3, \dots)$  be a sequence of i.i.d. random variables with finite variance  $V[X]$ . Then

$$E[\tilde{V}(\mathbf{X})] = \frac{n-1}{n} V[X]$$

# A BETTER VARIANCE ESTIMATOR

## The Unbiased Variance Estimator

Given a sample  $\mathbf{X} = (X_1, X_2, \dots, X_n)$ , the *unbiased sample variance estimator* is,

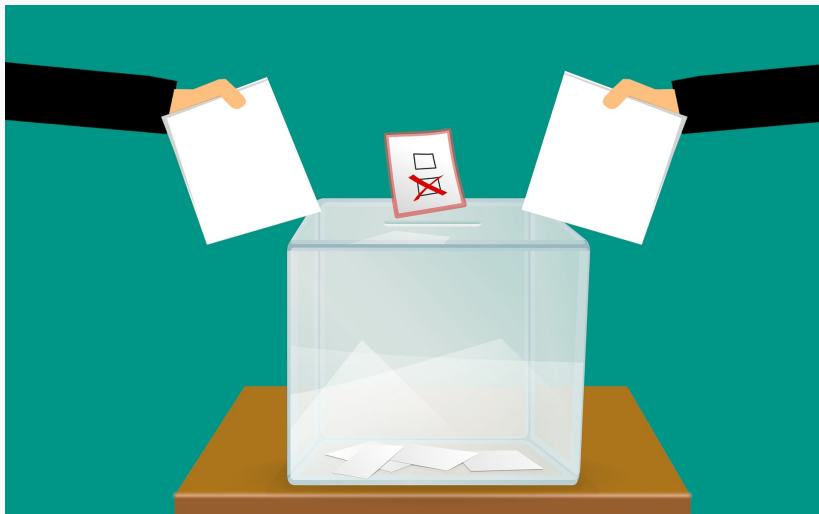
$$\hat{V}(\mathbf{X}) = \frac{n}{n-1} \left( \overline{X^2} - \overline{X}^2 \right)$$

# Standard Errors

---

**Point Estimate  $\Leftrightarrow$  Uncertainty**

# IMPORTANCE OF UNCERTAINTY



**Estimate:** Our candidate has 52% support.

# IMPORTANCE OF UNCERTAINTY



**Estimate:** The maximum dose of the medication is 26mg.

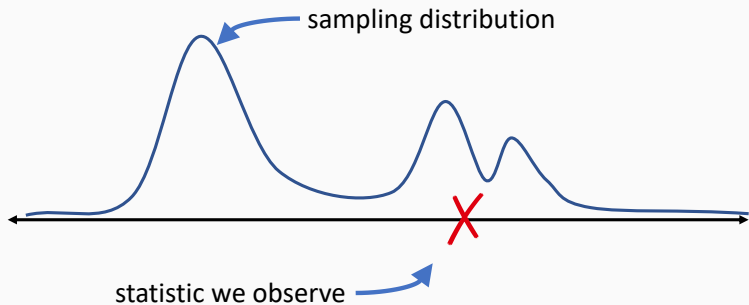
# IMPORTANCE OF UNCERTAINTY



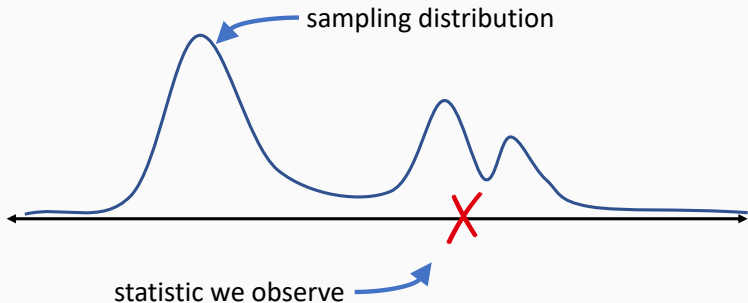
**Estimate:** The angle of the tower is 90 degrees.



# THE SAMPLING DISTRIBUTION OF THE STATISTIC



# THE SAMPLING DISTRIBUTION OF THE STATISTIC



**Standard Error:** (Estimated) standard deviation of the sampling distribution

## REPORTING STANDARD ERRORS

1. The mean number of mushrooms per pizza was 13.2 (SE 3.6).
2. The time for mice to navigate the maze was  $35 \pm 2$  seconds.

	Vitamin W	Vitamin X
3.	2.3 (0.3)	3.4 (0.9)

# **Standard Errors, The Sample Variance, and Standard Deviation**

---

# MANY MEASURES OF DISPERSION

**Sample Variance**

**Sampling Variance of the  
Sample Mean**

**Sample Standard  
Deviation**

**Standard Error of the  
Sample Mean**

# Motivating The Central Limit Theorem

---

# CAPTURING UNCERTAINTY

Need tools to capture how far off our estimate may be.

- Standard Error: the (estimated) standard deviation of the sampling distribution of the estimator.
- But a standard error is just one number...

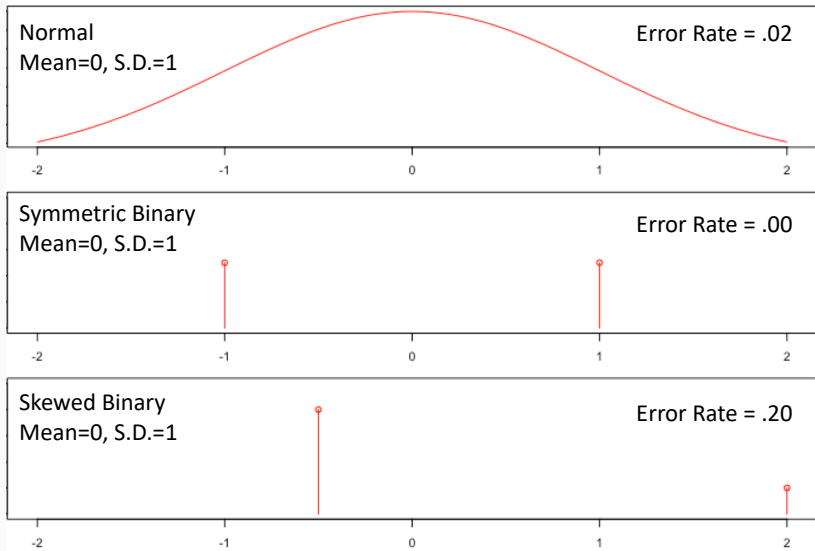
# UNCERTAINTY EXAMPLE



- Hidden fact: Population mean = 0
- Idea: Market if estimate is 2 standard errors above 0.



# THE IMPORTANCE OF SHAPE



# THE CENTRAL LIMIT THEOREM

## Central limit theorem (CLT): idea

For a **very broad** class of population distributions, the sampling distribution of the mean becomes approximately normal as the sample size grows large.

# **Reading: The Central Limit Theorem**

---

## READING: THE CENTRAL LIMIT THEOREM

**Note: This is a reading call, we're just placing it here for organization.** Read pages 108 and 109 of section 3.2.4. There's no *need* to read through the demonstration of Slutsky's theorem, but you can if you would like.

- Rather than proving the CLT, we are going to ask you to work through a short demonstration against data that we hope will convince you of the CLT's effectiveness.
- The book is terse in its presentation of when and how the CLT applies. We will fill that out when we come back together.

# **Apply the Central Limit Theorem**

---

# APPLY THE CENTRAL LIMIT THEOREM

**Note: This is a learnosity activity, just placing it here for organization.**

# Central Limit Theorem

---

## REMINDER OF CONTEXT

- The WLLN tells us what happens to the sample mean as  $n \rightarrow \infty$ :

$$\bar{X} \xrightarrow{p} E[X]$$

- We also know that, as  $n \rightarrow \infty$ , we can generate an increasingly good estimate for  $V[X]$  because

$$\overline{X^2} - \bar{X}^2 \xrightarrow{p} V[X]$$

- For more precise statements about uncertainty, we need the sampling distribution of the statistic.



# CONVERGENCE IN DISTRIBUTION

Statistic for  $n=1$   $T_{(1)}$



Statistic for  $n=2$   $T_{(2)}$



Statistic for  $n=3$   $T_{(3)}$



$\vdots$

$\vdots$

$\vdots$

Target Random Variable  $T$   
(Usually Normal)



# CONVERGENCE IN DISTRIBUTION

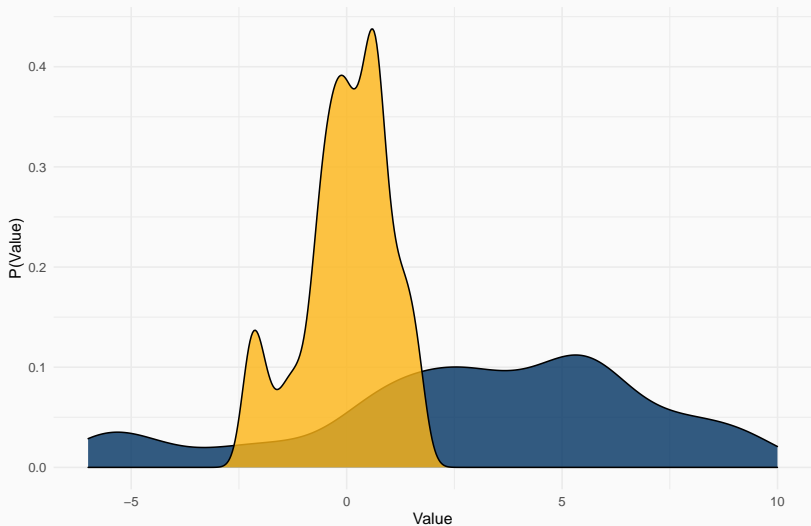
## Definition: Convergence in distribution

Let  $(T_{(1)}, T_{(2)}, T_{(3)}, \dots)$  be a sequence of random variables, with cdfs  $(F_{(1)}, F_{(2)}, F_{(3)}, \dots)$ , and let  $T$  be a random variable with cdf  $F$ . Then  $T_{(n)}$  *converges in distribution* to  $T$  if, for all  $t \in \mathbb{R}$  at which  $F$  is continuous,

$$\lim_{n \rightarrow \infty} F_{(n)}(t) = F(t).$$

We denote this as  $T_{(n)} \xrightarrow{d} T$ .

# THE NEED TO STANDARDIZE



# STANDARDIZING THE SAMPLE MEAN

## Definition: Standardized sample mean

For IID random variables  $(X_1, X_2, \dots, X_n)$  with finite  $E[x] = \mu$  and finite  $V[X] = \sigma^2$ , then **the standardized sample mean** is:

$$Z = \frac{(\bar{X} - E[\bar{X}])}{\sigma[\bar{X}]} = \frac{\sqrt{n}(\bar{X} - \mu)}{\sigma}.$$

- The standardized sample mean always has mean 0 and standard deviation 1.

# THE CENTRAL LIMIT THEOREM

## The Central Limit Theorem

Let  $(X_1, X_2, X_3, \dots)$  be a sequence of i.i.d. random variables with finite mean  $E[X] = \mu$  and finite variance  $V[X] = \sigma^2$ ,

$$Z = \frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \xrightarrow{d} N(0, 1)$$

# GENERALITY OF THE CLT

The CLT applies to:

- Continuous Random Variables
- Discrete Random Variables
- Symmetric Random Variables
- Asymmetric Random Variables

Only Requirements:

- Data points are IID.
- Population has finite variance.

Versions of the CLT exist for many other statistics.

# WHAT SAMPLE SIZE IS ENOUGH?

The CLT works in the limit as  $n \rightarrow \infty$ . What can we say for a finite  $n < \infty$ ?

- Rule of Thumb:  $n = 30$  for CLT to "kick in."
- Reality: Convergence depends on how non-normal population is.
  - A normal population requires  $n = 1$ .
  - Highly skewed distributions may require  $n = 100$ ,  $n = 1000$ , or more.

# **Learnosity: When Does the CLT Apply?**

---



## WHERE AND WHEN DOES THE CLT APPLY? PART III

**Which of these random variables has an approximately normal distribution because of the CLT?**

- $X$  = hours spent by one individual on a 203 homework assignment
- $X$  = average hours spent by randomly assigned study groups of size 4 on the same 203 homework
- $X$  = number of barks by a dog named Rex for a one-hour period at night

## WHERE AND WHEN DOES THE CLT APPLY? PART III

**Which of these random variables has an approximately normal distribution because of the CLT?**

- $X$  = total number of barks by a neighborhood of dogs for a one-hour period at night
- $X$  = age of a randomly selected MIDS student
- $X$  = average of sample of 10 randomly selected MIDS students
- $X$  = VADER sentiment of a sample of 100 SMS messages

# The Plug-In Principle

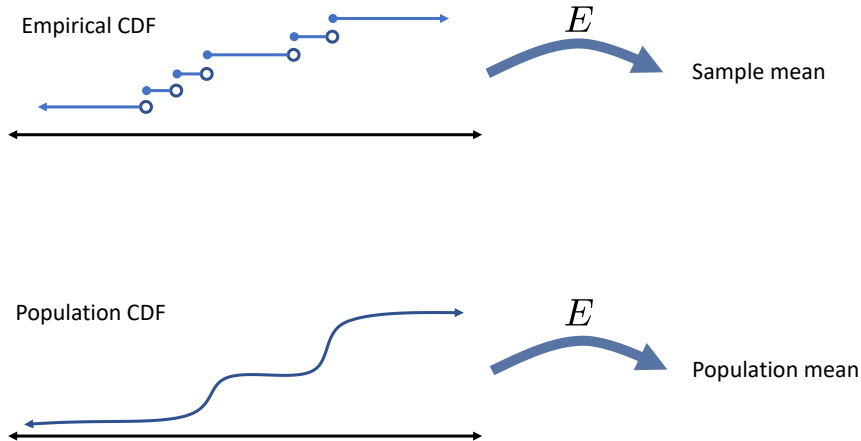
---

# THE PLUG-IN PRINCIPLE

- Want mean of population  $\rightarrow$  Mean of sample.
- Want variance of population  $\rightarrow$  Variance of sample.
- Want  $f(E[X], V[X])$ ?  $\rightarrow f(\bar{X}, \hat{V}(X))$

See the Pattern?

# THE PLUG-IN PRINCIPLE



# Reading Assignment

---

## READING ASSIGNMENT

**Only if you're interested,** read section 3.3 - 3.3.1 to learn more about the Plug-In Principle.

# Asymptotic Theory

---



# **Asymptotic Theory**

---

**Asymptotics Rescue Data Science**

**Across a range of applications, small samples are problems.**

- If we've got a *lot* of data, though, we can rely on weaker requirements
- Asymptotic properties of estimators ask the question, **What happens when we have a lot of data?**
- In general, as  $n \rightarrow \infty$ , does  $\hat{\theta}$  converge in distribution or probability to something that is useful or desirable?

# Asymptotic Theory

---

**Reading: Asymptotics**

## READING: ASYMPTOTICS

**Note: This is a READING CALL. We're placing it here for organization.**

- The interested student can read pages 111-114.
- But, we might recommend skipping it if you're constrained.
- The take home is that there is an asymptotic statement of:
  - Being Normally Distributed
  - SE, MSE and Efficiency
  - Sampling Variance and Sampling Standard Error

# **Asymptotic Theory**

---

**Reading: The Plug-In Principle**

## READING: THE PLUG-IN PRINCIPLE

- $F(x)$  is the cumulative distribution function, the *CDF*
- $f(x)$  is the probability distribution function, the *PDF*
- $\hat{F}(x) \neq F(x)$  and  $\hat{f}(x) \neq f(x)$
- As  $n \rightarrow \infty$  we hope that  $\hat{F}(x) \xrightarrow{d} F(x)$  and  $\hat{f}(x) \xrightarrow{d} f(x)$ .

# **Asymptotic Theory**

---

## **The Plug-In Principle**

# THE PLUG-IN PRINCIPLE

## Expectation and Variance Functionals

You *know* what the processes for calculating an expectation and variance are.

$$E[X] = T_E(F) = \int x \cdot dF(x)$$

$$V[E] = T_V(F) = \int (x - E[X])^2 \cdot dF(x)$$

These are the *statistical functionals* for expectation and variance.



# THE PLUG-IN PRINCIPLE (CONT.)

## Plug-In Estimators

For i.i.d. random variables  $X_1, X_2, \dots, X_n$ , with a common CDF  $F$ , the plug in estimator for  $\theta = T(F)$  is just  $\hat{\theta} = T(\hat{F})$ .

$$\begin{aligned}\hat{E}(X) &= T_E(\hat{F}) \\ &= \sum x \cdot \hat{f}(x) \\ &= \sum x \cdot \frac{I(X_i = x)}{n} \\ &= \frac{1}{n} \sum x \\ &= \bar{X}\end{aligned}$$

# EXAMPLE OF THE PLUG-IN PRINCIPLE

## Example with Discrete Data

Suppose you there is some discrete process that places values into the random variable  $X$ .

- $F_X(1) = 1/3$
- $F_X(2) = 1/6$
- $F_X(3) = 0$
- $F_X(4) = 1/2$

But, you do get to produce i.i.d. random draws from the process. Of 1000 draws:

<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>
331	156	0	513

## EXAMPLE OF THE PLUG-IN PRINCIPLE, CONT'D

### Example with Discrete Data

$$\hat{E}(\mathbf{X}) = T_E(\hat{F})$$

# Asymptotic Theory

---

**Reading: Kernel Methods Estimate the PDF**

## READING: KERNEL METHODS ESTTIMATE THE PDF

Read pages 121 - 124.

# **Asymptotic Theory**

---

**Kernel Methods Estimate the PDF**

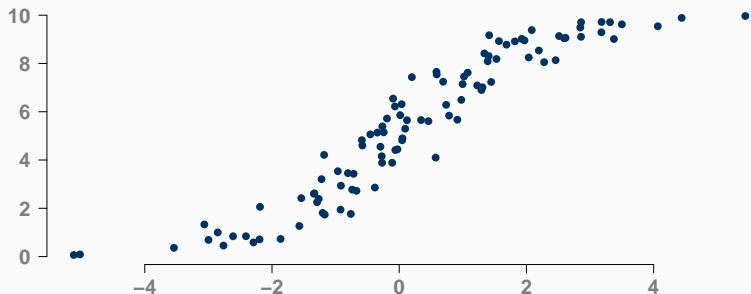
## Kernel Density Estimator of PDF

- Cannot *directly* observe the joint PDF.
- For discrete RV, approximations come through frequency tables.
- For continuous RV, approximations come through kernel density estimates:

$$\hat{f}_K(x) = \frac{1}{n} \sum_{i=1}^n K_b(x - X_i), \forall x \in \mathbb{R}$$

## KERNEL METHODS, PART II

- Kernel methods are smoothing methods
- The *kernel* is some weighting function





## KERNEL METHODS, PART III

As smoothing function, these permit plug in estimates that are directly analogous to the estimating functionals of the CDF.

### Kernel Plug-In Estimator

The *Kernel plug-in estimator* of  $\theta = T(F)$  is

$$\hat{\theta}_k = T(\hat{F}_K),$$

where,  $\hat{F}_K = \int_{-\infty}^x \hat{f}(u) du$ .

# KERNEL METHODS, PART IV

## Kernel Plug-In Estimator for Expected Value

We know the functional for the expected value has a specific “shape”

$$E[Y] = \int y \cdot df(y).$$

So, the feasible kernel method is

$$\hat{E}_k(\mathbf{Y}) = \int_{-\infty}^{\infty} y \cdot d\hat{f}_k(y)$$

## KERNEL METHODS, PART V

- $E[Y]$  and  $E[Y|X]$  are lowest MSE estimates of  $Y$

$$\hat{E}_K(\mathbf{Y}) = \int y \hat{f}(y)$$
$$\hat{E}_K(\mathbf{Y}|X = x] = \int y \hat{f}_{Y|X}(y|x)$$

- If you can produce i.i.d. samples from  $f_{Y|X}$ , you can produce estimates,  $\hat{f}_K(y|x)$  that get ever closer to the true value

# **Asymptotic Theory**

---

**Apply Kernel Methods**

# APPLY KERNEL METHODS

**Note: This is a LEARNOSITY activity. We're just placing it here for organization.**