

# Economics 103 – Statistics for Economists

Francis J. DiTraglia

University of Pennsylvania

# Lecture #1 – Introduction

Overview – Population vs. Sample, Probability vs. Statistics

Polling – Sampling vs. Non-sampling Error, Random Sampling

Causality – Observational vs. Experimental Data, RCTs

# Racial Discrimination in the Labor Market

Source: Bureau of Labor Statistics

	Oct. 2018	Nov. 2018	Dec. 2018
White:	3.0	3.0	3.1
Black/African American:	6.2	5.8	6.2

**Table:** Unemployment rate in percentage points for men aged 20 and over in the last quarter of 2018.

The unemployment rate for African Americans has historically been much higher than for whites. What can this information by itself tell us about racial discrimination in the labor market?

# This Course: Use Sample to Learn About Population

## Population

Complete set of all items that interest investigator

## Sample

Observed subset, or portion, of a population

## Sample Size

# of items in the sample, typically denoted  $n$

Examples...

# In Particular: Use Statistic to Learn about Parameter

## Parameter

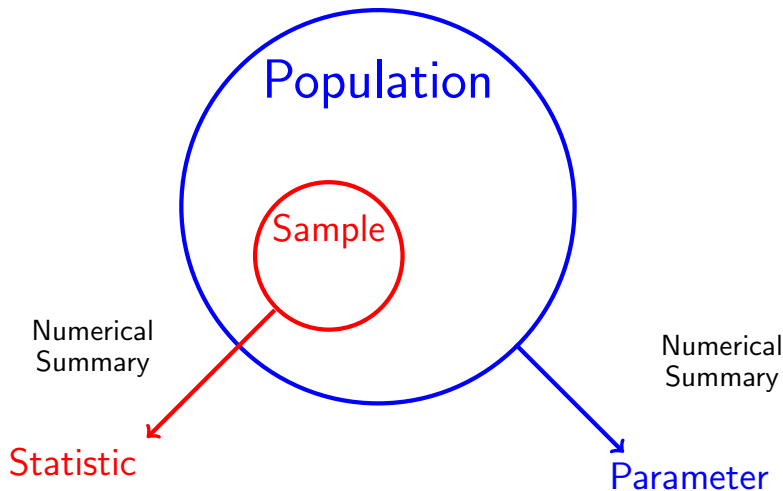
Numerical measure that describes specific characteristic of a population.

## Statistic

Numerical measure that describes specific characteristic of sample.

Examples...

## Essential Distinction You Must Remember!



# This Course

1. Descriptive Statistics: summarize data
  - ▶ Summary Statistics
  - ▶ Graphics
2. Probability: Population  $\rightarrow$  Sample
  - ▶ deductive: “safe” argument
    - ▶ All ravens are black. Mordecai is a raven, so Mordecai is black.
3. Inferential Statistics: Sample  $\rightarrow$  Population
  - ▶ inductive: “risky” argument
    - ▶ I’ve only every seen black ravens, so all ravens must be black.

# Sampling and Nonsampling Error

In statistics we use samples to learn about populations, but samples almost never be *exactly* like the population they are drawn from.

## 1. Sampling Error

- ▶ *Random* differences between sample and population
- ▶ Cancel out on average
- ▶ Decreases as sample size grows

## 2. Nonsampling Error

- ▶ *Systematic* differences between sample and population
- ▶ Does *not* cancel out on average
- ▶ Does *not* decrease as sample size grows



NEW COLORED MAP OF POLAND IN THIS ISSUE

Showing the Territorial Changes Wrought by the War

# The Literary Digest

(Title Reg. U.S. Pat. Off.)



New York FUNK & WAGNALLS COMPANY London

PUBLIC OPINION *New York* combined with *The LITERARY DIGEST*

Vol. 68, No. 8. Whole No. 1609

FEBRUARY 19, 1921

Price 15 CENTS

# Literary Digest – 1936 Presidential Election Poll



FDR versus Kansas Gov. Alf Landon

## Huge Sample

Sent out over 10 million ballots; 2.4 million replies! (Compared to less than 45 million votes cast in actual election)

## Prediction

Landslide for Landon: *Landonslide*, if you will.

# Spectacularly Mistaken!



FDR versus Kansas Gov. Alf Landon

	Roosevelt	Landon
Literary Digest Prediction:	41%	57%
Actual Result:	61%	37%

# What Went Wrong? *Non-sampling Error (aka Bias)*

Source: Squire (1988)

## Biased Sample

Some units more likely to be sampled than others.

- ▶ Ballots mailed those on auto reg. list and in phone books.

## Non-response Bias

Even if sample is unbiased, can't force people to reply.

- ▶ Among those who recieved a ballot, Landon supporters were more likely to reply.

In this case, neither effect *alone* was enough to throw off the result but together they did.

# Randomize to Get an Unbiased Sample

## Simple Random Sample

Each member of population is chosen strictly by chance, so that:  
(1) selection of one individual doesn't influence selection of any other, (2) each individual is just as likely to be chosen, (3) every possible sample of size  $n$  has the same chance of selection.

What about non-response bias? – we'll come back to this...

## “Negative Views of Trump’s Transition” (Jan, 2017)

Source: [Pew Research Center](#)

*Ahead of Donald Trump’s scheduled press conference in New York City on Wednesday, the public continues to give the president-elect low marks for how he is handling the transition process. . . The latest national survey by Pew Research Center, conducted Jan. 4-9 among 1,502 adults, finds that 39% approve of the job President-elect Trump has done so far explaining his policies and plans for the future to the American people, while a larger share (55%) say they disapprove.*

# Quantifying Sampling Error

## 95% Confidence Interval for Poll Based on Random Sample

### Margin of Error a.k.a. ME

We report  $P \pm \text{ME}$  where  $\text{ME} \approx 2\sqrt{P(1-P)/n}$

### Trump Transition Approval Rate

$P = 0.39$  and  $n = 1502$  so  $\text{ME} \approx 0.025$ . We'd report 39% plus or minus 2.5% if the poll were based on a simple random sample. . .

But Pew Reports an ME of 2.9% which doesn't agree with our calculation. What's going on here?!

# Non-response bias is a huge problem. . .

Source: Pew Research Center

---

## Surveys Face Growing Difficulty Reaching, Persuading Potential Respondents

	1997	2000	2003	2006	2009	2012
	%	%	%	%	%	%
<b>Contact rate</b> (percent of households in which an adult was reached)	90	77	79	73	72	62
<b>Cooperation rate</b> (percent of households contacted that yielded an interview)	43	40	34	31	21	14
<b>Response rate</b> (percent of households sampled that yielded an interview)	36	28	25	21	15	9

PEW RESEARCH CENTER 2012 Methodology Study. Rates computed according to American Association for Public Opinion Research (AAPOR) standard definitions for CON2, COOP3 and RR3. Rates are typical for surveys conducted in each year.

---



# Methodology – “Negative Views of Trump’s Transition”

Source: [Pew Research Center](#)

*The combined landline and cell phone sample are weighted using an iterative technique that matches gender, age, education, race, Hispanic origin and nativity and region to parameters from the 2015 Census Bureaus American Community Survey and population density to parameters from the Decennial Census. The sample also is weighted to match current patterns of telephone status (landline only, cell phone only, or both landline and cell phone), based on extrapolations from the 2016 National Health Interview Survey. The weighting procedure also accounts for the fact that respondents with both landline and cell phones have a greater probability of being included in the combined sample and adjusts for household size among respondents with a landline phone. The margins of error reported and statistical tests of significance are adjusted to account for the surveys design effect, a measure of how much efficiency is lost from the weighting procedures.*

# Simple Example of Weighting a Survey

## Post-stratification

- ▶ Women make up 49.6% of the population but suppose they are less likely to respond to your survey than men.
- ▶ If women have different opinions of Trump, this will skew the survey.
- ▶ Calculate Trump approval rate separately for men  $P_M$  vs. women  $P_W$ .
- ▶ Report  $0.496 \times P_W + 0.504 \times P_M$ , not the raw approval rate  $P$ .

## Caveats

- ▶ Post-stratification isn't a magic bullet: you have to figure out what factors could skew your poll to adjust for them.
- ▶ Calculating the ME is more complicated. Since this is an intro class we'll focus on simple random samples.



## Survey to find effect of Polio Vaccine

Ask random sample of parents if they vaccinated their kids or not and if the kids later developed polio. Compare those who were vaccinated to those who weren't.

Would this procedure:

- (a) Overstate effectiveness of vaccine
- (b) Correctly identify effectiveness of vaccine
- (c) Understate effectiveness of vaccine

# Confounding

Parents who vaccinate their kids may differ systematically from those who don't in *other ways* that impact child's chance of contracting polio!

Wealth is related to vaccination *and* whether child grows up in a hygienic environment.

## Confounder

Factor that influences both outcomes and whether subjects are treated or not. Masks true effect of treatment.

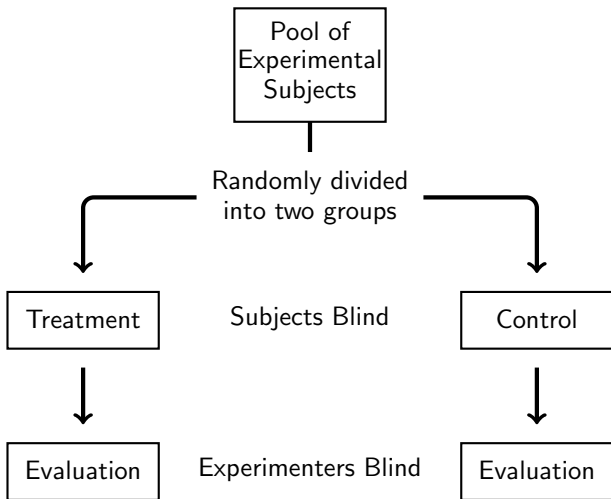
# Experiment Using Random Assignment: Randomized Experiment

Treatment Group Gets Vaccine, Control Group Doesn't

## Essential Point!

Random assignment *neutralizes* effect of all confounding factors: since groups are initially equal, on average, any difference that emerges must be the treatment effect.

## Placebo Effect and Randomized Double Blind Experiment



## Gold Standard: Randomized, Double-blind Experiment

*Randomized blind experiments ensure that on average the two groups are initially equal, and continue to be treated equally. Thus a fair comparison is possible.*

Randomized, double-blind experiments are considered the “gold standard” for untangling causation.

Sugar Doesn't Make Kids Hyper

<http://www.youtube.com/watch?v=mkr9YsmrPAI>

Randomization is not always possible, practical, or ethical.

## Observational Data

Data that do not come from a randomized experiment.

It much more challenging to untangle cause and effect using observational data because of confounders. But sometimes it's all we have.



# Racial Bias in the Labor Market

Bertrand & Mullainathan (2004, American Economic Review)

*When faced with observably similar African-American and White applicants, do they [employers] favor the White one? Some argue yes, citing either employer prejudice or employer perception that race signals lower productivity. Others argue that differential treatment by race is a relic of the past . . . Data limitations make it difficult to empirically test these views. Since researchers possess far less data than employers do, White and African-American workers that appear similar to researchers may look very different to employers. So any racial difference in labor market outcomes could just as easily be attributed to differences that are observable to employers but unobservable to researchers.*

## Racial Bias in the Labor Market: continued . . .

Bertrand & Mullainathan (2004, American Economic Review)

*To circumvent this difficulty, we conduct a field experiment . . . We send resumes in response to help-wanted ads in Chicago and Boston newspapers and measure call-back for interview for each sent resume. We experimentally manipulate the perception of race via the name of the fictitious job applicant. We randomly assign very White-sounding names (such as Emily Walsh or Greg Baker) to half the resumes and very African-American-sounding names (such as Lakisha Washington or Jamal Jones) to the other half.*

## Racial Bias in the Labor Market: continued . . .

Bertrand & Mullainathan (2004, American Economic Review)

Sample	White Names	African-American Names
All sent resumes	9.7	6.5
Females	9.9	6.6
Males	8.9	5.8

Table: % Callback by racial soundingness of names.

Later this semester: if there were no racial bias in callbacks, what is the chance that we would observe such large differences?

# Lecture #2 – Summary Statistics Part I

Class Survey

Types of Variables

Frequency, Relative Frequency, & Histograms

Measures of Central Tendency

Measures of Variability / Spread

# Class Survey

- ▶ Collect some data to analyze later in the semester.
- ▶ None of the questions are sensitive and your name will not be linked to your responses. I will post an anonymized version of the dataset on my website.
- ▶ The survey is *strictly voluntary* – if you don't want to participate, you don't have to.



## Multiple Choice Entry – What is your biological sex?

- (a) Male
- (b) Female



## Multiple Choice – What is Your Eye Color?

Please enter your eye color using your remote.

- (a) Black
- (b) Blue
- (c) Brown
- (d) Green
- (e) Gray
- (f) Hazel
- (g) Other



## How Right-Handed are You?

The sheet in front of you contains a handedness inventory. Please complete it and calculate your handedness score:

$$\frac{\text{Right} - \text{Left}}{\text{Right} + \text{Left}}$$

When finished, enter your score using your remote.





## What is your Height in Inches?

Using your remote, please enter your height in inches, rounded to the nearest inch:

$$4\text{ft} = 48\text{in}$$

$$5\text{ft} = 60\text{in}$$

$$6\text{ft} = 72\text{in}$$

$$7\text{ft} = 84\text{in}$$



## What is your Hand Span (in cm)?

On the sheet in front of you is a ruler. Please use it to measure the span of your right hand in centimeters, to the nearest  $1/2$  cm.

*Hand Span: the distance from thumb to little finger  
when your fingers are spread apart*

When ready, enter your measurement using your remote.



We chose (by computer) a random number between 0 and 100.  
The number selected and assigned to you is written on the slip of paper in front of you. Please do not show your number to anyone else or look at anyone else's number.

Please enter your number now using your remote.



Call your random number  $X$ . Do you think that the **percentage** of countries, among all those in the United Nations, that are in Africa is **higher** or **lower** than  $X$ ?

(a) Higher

(b) Lower

Please answer using your remote.



What is your best estimate of the **percentage** of countries, among all those that are in the United Nations, that are in Africa?

Please enter your answer using your remote.

# Types of Variables

Categorical = Qualitative

Numeric value either meaningless or indicates order only

**Nominal** unordered: eye color, sex

**Ordinal** ordered: course evaluations (0 = Poor, 1 = Fair)

Numerical = Quantitative

Numerical value is meaningful

**Discrete** # of credits you are taking this semester

**Continuous** height, handspan, handedness score

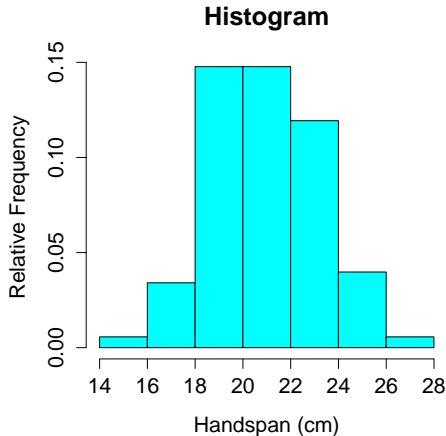
## Handspan - Frequency and Relative Frequency

cm	Freq.	Rel. Freq.
14.0	1	0.01
17.0	4	0.05
17.5	2	0.02
18.0	5	0.06
18.5	5	0.06
19.0	6	0.07
19.5	10	0.11
20.0	10	0.11
20.5	3	0.03
21.0	8	0.09
21.5	5	0.06
22.0	9	0.10
22.5	6	0.07
23.0	6	0.07
24.0	4	0.05
24.5	3	0.03
27.0	1	0.01
<hr/> $n = 88$		1.00



# Histogram – Density Estimate by Smoothing Barchart

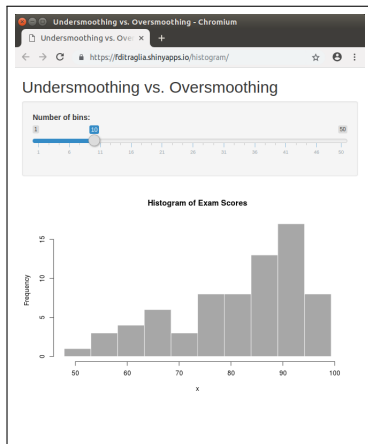
Bins	Freq.	Rel. Freq.
[14, 16)	1	0.01
[16, 18)	6	0.07
[18, 20)	26	0.30
[20, 22)	26	0.30
[22, 24)	21	0.24
[24, 26)	7	0.08
[26, 28)	1	0.01
$n = 88$		1.00



Group data into non-overlapping bins of equal width



<https://fditraglia.shinyapps.io/histogram/>



The number of histogram bins controls the degree of *smoothing*.

# Histogram - Density Estimate by Smoothing Barchart

## Why Histogram?

Summarize numerical data, especially continuous (few repeats)

## Too Many Bins – Undersmoothing

No longer a summary (lose the shape of distribution)

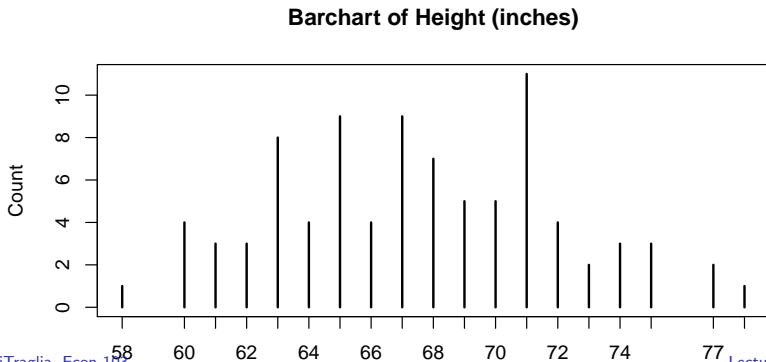
## Too Few Bins – Oversmoothing

Miss important detail

Don't confuse with barchart!

```
# Read data
data_url <- 'http://ditraglia.com/econ103/old_survey.csv'
survey <- read.csv(data_url)

#Make plot
plot(table(survey$height), main = 'Barchart of Height (inches)',
      xlab = '', ylab = 'Count')
```



```
hist(survey$height, freq = FALSE, main = 'Histogram of Height',  
     xlab = 'Height (in)', ylab = 'Relative Frequency')
```



# Summary Statistic = Numerical Summary of Sample

## Categories of Summary Statistic

1. Central Tendency: mean and median
2. Spread: range, interquartile range, variance, and std. dev.
3. Symmetry: skewness
4. Linear Dependence: covariance, correlation, and regression

## Questions ask yourself about each summary statistic

1. What does it measure?
2. What are its units compared to those of the data?
3. (How) do its units change if those of the data change?

# What is an Outlier?

## Outlier

A very unusual observation relative to the other observations in the dataset (i.e. very small or very big).

# Measures of Central Tendency

Suppose we have a dataset with observations  $x_1, x_2, \dots, x_n$

## Sample Mean

- ▶  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
- ▶ Only for numeric data
- ▶ Sensitive to asymmetry and outliers

## Sample Median

- ▶ Middle observation if  $n$  is odd, otherwise the mean of the two observations closest to the middle.
- ▶ Applicable to numerical or ordinal data
- ▶ Insensitive to outliers and skewness

## Mean is Sensitive to Outliers, Median Isn't

First Dataset: 1 2 3 4 5

Mean = 3, Median = 3

Second Dataset: 1 2 3 4 4990

Mean = 1000, Median = 3

When Does the Median Change?

Ranks would have to change so that 3 is no longer in the middle.



# Percentage of UN Countries that are in Africa

## You Were a Subject in a Randomized Experiment!

- ▶ There were only two numbers in the bag: 10 and 65
- ▶ Randomly assigned to Low group (10) or High group (65)

## Anchoring Heuristic (Kahneman and Tversky, 1974)

Subjects' estimates of an unknown quantity are influenced by an irrelevant previously supplied starting point.

Are Penn students subject to to this cognitive bias?

## Results from Anchoring Experiment (Previous Semester)

```
low <- subset(survey, rand.num == 10)$africa.percent
high <- subset(survey, rand.num == 65)$africa.percent
c(low = mean(low), high = mean(high))

##      low      high
## 17.09302 30.71739

c(low = median(low), high = median(high))

##  low high
##   17   30
```

# Percentiles (aka Quantiles) – Generalization of Median

## Percentiles (aka Quantiles)

Approx.  $P\%$  of the data are at or below the  $P^{\text{th}}$  percentile/quantile

## Quartiles

Q1 = 25th Percentile

Q2 = Median (i.e. 50th Percentile)

Q3 = 75th Percentile

There are some slightly tricky issues involved in actually *calculating* quantiles, but these only make a difference for very small datasets. We'll always use R to calculate quantiles. . .

```
quantile(survey$handspan, na.rm = TRUE)
```

```
##    0%   25%   50%   75%  100%
```

```
## 14.0 19.0 20.5 22.0 27.0
```

```
quantile(survey$handspan, 0.3, na.rm = TRUE)
```

```
##   30%
```

```
## 19.5
```

```
quantile(survey$handspan, c(0.1, 0.5, 0.9), na.rm = TRUE)
```

```
##   10%   50%   90%
```

```
## 18.0 20.5 23.0
```

## Boxplot: A Depiction of the “Five Number Summary”

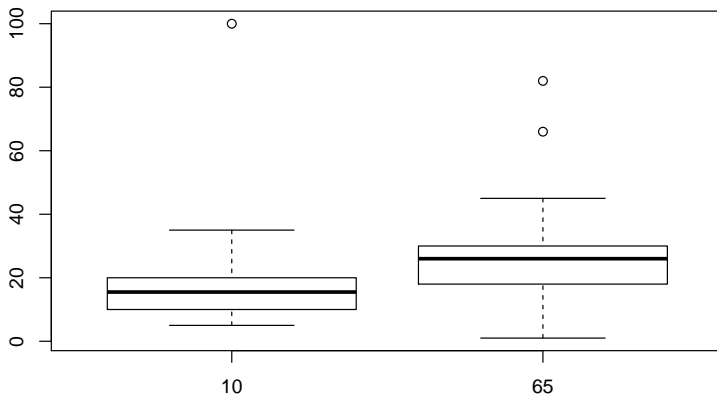


The `boxplot` command in R treats any observation more than 1.5 times the *width* of the box away from the box as an outlier.

```
boxplot(survey$handspan, main = 'Boxplot of Handspan',  
        ylab = 'Handspan (cm)')
```



```
boxplot(survey$africa.percent ~ survey$rand.num,  
        main = 'Boxplot for Anchoring Experiment',  
        ylab = 'Answer (% UN Countries from Africa)',  
        xlab = 'Random Number')
```



# Measures of Variability/Spread – 1

## Range

- ▶ Range = Maximum Observation - Minimum Observation
- ▶ Very sensitive to outliers.
- ▶ Displayed in boxplot.

## Interquartile Range (IQR)

- ▶  $IQR = Q_3 - Q_1$
- ▶ IQR = Range of middle 50% of the data.
- ▶ Insensitive to outliers.
- ▶ Displayed in boxplot.



# Measures of Variability/Spread – 2

## Variance

- ▶  $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$
- ▶ Essentially the average squared distance from the mean.
- ▶ (We'll talk about  $n - 1$  versus  $n$  later in the semester)
- ▶ Sensitive to both skewness and outliers.

## Standard Deviation

- ▶  $s = \sqrt{s^2}$
- ▶ Same information as variance but more convenient since it has the **same units as the data**

# Measures of Spread for Handspan

```
diff(range(survey$handspan, na.rm = TRUE))
```

```
## [1] 13
```

```
IQR(survey$handspan, na.rm = TRUE)
```

```
## [1] 3
```

```
var(survey$handspan, na.rm = TRUE)
```

```
## [1] 4.753788
```

```
sd(survey$handspan, na.rm = TRUE)
```

```
## [1] 2.180318
```

# Lecture #3 – Summary Statistics Part II

Why squares in the definition of variance?

Skewness & Symmetry

Sample versus Population, Empirical Rule

Centering, Standardizing, & Z-Scores

Relating Two Variables: Cross-tabs, Covariance, & Correlation

# Why Squares?

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

What's Wrong With This?

$$\begin{aligned} \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}) &= \frac{1}{n-1} \left[ \sum_{i=1}^n x_i - \sum_{i=1}^n \bar{x} \right] = \frac{1}{n-1} \left[ \sum_{i=1}^n x_i - n\bar{x} \right] \\ &= \frac{1}{n-1} \left[ \sum_{i=1}^n x_i - n \cdot \frac{1}{n} \sum_{i=1}^n x_i \right] \\ &= \frac{1}{n-1} \left[ \sum_{i=1}^n x_i - \sum_{i=1}^n x_i \right] = 0 \end{aligned}$$

# Skewness – A Measure of Symmetry

$$\text{Skewness} = \frac{1}{n} \frac{\sum_{i=1}^n (x_i - \bar{x})^3}{s^3}$$

What do the values indicate?

Zero  $\Rightarrow$  symmetry, positive right-skewed, negative left-skewed.

Why cubed?

To get the desired sign.

Why divide by  $s^3$ ?

So that skewness is unitless

Rule of Thumb

Typically (but not always), right-skewed  $\Rightarrow$  mean  $>$  median

left-skewed  $\Rightarrow$  mean  $<$  median

```
# Load Survey Data
```

```
data_url <- 'http://ditraglia.com/econ103/old_survey.csv'
```

```
survey <- read.csv(data_url)
```

```
# A Function to Calculate Skewness
```

```
get_skewness <- function(x) {  
  x <- na.omit(x)  
  n <- length(x)  
  xbar <- mean(x)  
  s <- sd(x)  
  skewness <- sum((x - xbar)^3) / (n * s^3)  
  return(skewness)  
}
```

```
# Handedness is left-skewed, handspan is symmetric
c(get_skewness(survey$handedness), get_skewness(survey$handspan))

## [1] -2.21905550  0.04331997

par(mfrow = c(1, 2))
hist(survey$handedness, main = 'Handedness', xlab = 'Handedness Score')
hist(survey$handspan, main = 'Handspan', xlab = 'Handspan (cm)')
```



# Sample vs. Population and Parameter vs. Statistic

## Sample vs. Population

For now, think of the **population** as a list of  $N$  objects  $(x_1, x_2, \dots, x_N)$  from which we draw a **sample** of  $n < N$  objects.

## Parameter vs. Statistic

Use a sample to calculate **statistics** (e.g.  $\bar{x}$ ,  $s^2$ ,  $s$ ) that estimate the corresponding population **parameters** (e.g.  $\mu$ ,  $\sigma^2$ ,  $\sigma$ ).

	Parameter (Population)	Statistic (Sample)
Mean	$\mu = \frac{1}{N} \sum_{i=1}^N x_i$	$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
Var.	$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$	$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$
S.D.	$\sigma = \sqrt{\sigma^2}$	$s = \sqrt{s^2}$



# Why Mean and Variance (and Std. Dev. )?

## Empirical Rule

For large populations that are approximately bell-shaped, std. dev. tells where most observations will be relative to the mean:

- ▶  $\approx 68\%$  of observations are in the interval  $\mu \pm \sigma$
- ▶  $\approx 95\%$  of observations are in the interval  $\mu \pm 2\sigma$
- ▶ Almost all of observations are in the interval  $\mu \pm 3\sigma$

This is a key reason why we will be interested in  $\bar{x}$  as an estimate of  $\mu$  and  $s$  as an estimate of  $\sigma$ .



Which is more “extreme?”

- (a) Handspan of 27cm
- (b) Height of 78in

# Centering: Subtract the Mean

Handspan	Height
$27\text{cm} - 20.6\text{cm} = 6.4\text{cm}$	$78\text{in} - 67.6\text{in} = 10.4\text{in}$

# Standardizing: Divide by S.D.

Handspan	Height
$27\text{cm} - 20.6\text{cm} = 6.4\text{cm}$	$78\text{in} - 67.6\text{in} = 10.4\text{in}$
$6.4\text{cm}/2.2\text{cm} \approx 2.9$	$10.4\text{in}/4.5\text{in} \approx 2.3$

The units have disappeared!

# Z-scores: How many standard deviations from the mean?

Best for Symmetric Distribution, No Outliers (Why?)

$$z_i = \frac{x_i - \bar{x}}{s}$$

## Unitless

Allows comparison of variables with different units.

## Detecting Outliers

Measures how “extreme” one observation is relative to the others.

## Linear Transformation

What is the sample mean of the z-scores?

$$\bar{z} = \frac{1}{n} \sum_{i=1}^n z_i = \frac{1}{n} \sum_{i=1}^n \frac{x_i - \bar{x}}{s} = \frac{1}{n \cdot s} \sum_{i=1}^n (x_i - \bar{x}) = 0$$

... using the same argument as on Slide 2 of this lecture!

## What is the variance of the z-scores?

$$\begin{aligned}s_z^2 &= \frac{1}{n-1} \sum_{i=1}^n (z_i - \bar{z})^2 = \frac{1}{n-1} \sum_{i=1}^n z_i^2 = \frac{1}{n-1} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s_x} \right)^2 \\&= \frac{1}{s_x^2} \left[ \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \right] = \frac{s_x^2}{s_x^2} = 1\end{aligned}$$

So what is the *standard deviation* of the z-scores?



## Population Z-scores and the Empirical Rule: $\mu \pm 2\sigma$

If  $\mu$  and  $\sigma$  were known, we could create a *population version* of a z-score. This lets us re-write the Empirical Rule as follows:

Bell-shaped population  $\Rightarrow$  approx. 95% of observations  $x_i$  satisfy

$$\mu - 2\sigma \leq x_i \leq \mu + 2\sigma$$

$$-2 \leq \frac{x_i - \mu}{\sigma} \leq 2$$



## Crosstabs – Show Relationship between Categorical Vars.

```
table(survey$eye.color, survey$sex)
```

```
##
```

```
##           Female Male
```

```
##   Black         2    5
```

```
##   Blue          4    6
```

```
##   Brown        32   26
```

```
##   Copper        0    1
```

```
##   Green         1    4
```

```
##   Hazel         2    2
```

```
##   Maroon        0    1
```

# Who Supported the Vietnam War?

In January 1971 the Gallup poll asked: “A proposal has been made in Congress to require the U.S. government to bring home all U.S. troops before the end of this year. Would you like to have your congressman vote for or against this proposal?”

Guess the results, for respondents in each education category, and fill out this table (the two numbers in each column should add up to 100%):

	Adults with:			
	Grade school education	High school education	College education	Total adults
% for withdrawal of U.S. troops (doves)				73%
% against withdrawal of U.S. troops (hawks)				27%
Total	100%	100%	100%	100%



## Who Were the Doves?

Which group do you think was most strongly **in favor of** the withdrawal of US troops from Vietnam?

- (a) Adults with only a Grade School Education
- (b) Adults with a High School Education
- (c) Adults with a College Education

Please respond with your remote.



## Who Were the Hawks?

Which group do you think was most strongly **opposed to** the withdrawal of US troops from Vietnam?

- (a) Adults with only a Grade School Education
- (b) Adults with a High School Education
- (c) Adults with a College Education

Please respond with your remote.

# Who *Really* Supported the Vietnam War

Gallup Poll, January 1971

	Adults with:			Total adults
	Grade school education	High school education	College education	
% for withdrawal of U.S. troops (doves)	80%	75%	60%	73%
% against withdrawal of U.S. troops (hawks)	20%	25%	40%	27%
Total	100%	100%	100%	100%

# Covariance and Correlation: Linear Dependence Measures

## Two Samples of Numeric Data

$x_1, \dots, x_n$  and  $y_1, \dots, y_n$  with means  $(\bar{x}, \bar{y})$  and std. devs.  $(s_x, s_y)$

## Dependence

Do  $x$  and  $y$  both tend to be large (or small) at the same time?

## Key Point

Use the idea of centering and standardizing to decide what “big” or “small” means in this context.

# Covariance

$$s_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

- ▶ Centers each observation around its mean and multiplies.
- ▶ Zero  $\Rightarrow$  no linear dependence
- ▶ Positive  $\Rightarrow$  positive linear dependence
- ▶ Negative  $\Rightarrow$  negative linear dependence
- ▶ Population parameter:  $\sigma_{xy}$
- ▶ Units?

# Correlation

$$r_{xy} = \frac{1}{n-1} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right) = \frac{s_{xy}}{s_x s_y}$$

- ▶ Centers *and* standardizes each observation
- ▶ Bounded between -1 and 1
- ▶ Zero  $\Rightarrow$  no linear dependence
- ▶ Positive  $\Rightarrow$  positive linear dependence
- ▶ Negative  $\Rightarrow$  negative linear dependence
- ▶ Population parameter:  $\rho_{xy}$
- ▶ Unitless



## Height and Handspan: Strongly Positively Associated

```
cov(survey$height, survey$handspan, use = 'complete.obs')  
  
## [1] 5.910786  
  
cor(survey$height, survey$handspan, use = 'complete.obs')  
  
## [1] 0.6042423
```

# Essential Distinction: Parameter vs. Statistic

## And Population vs. Sample

$N$  individuals in the Population,  $n$  individuals in the Sample:

	<b>Parameter</b> (Population)	<b>Statistic</b> (Sample)
Mean	$\mu_x = \frac{1}{N} \sum_{i=1}^N x_i$	$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
Var.	$\sigma_x^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$	$s_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$
S.D.	$\sigma_x = \sqrt{\sigma_x^2}$	$s_x = \sqrt{s^2}$
Cov.	$\sigma_{xy} = \frac{\sum_{i=1}^N (x_i - \mu_x)(y_i - \mu_y)}{N}$	$s_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$
Corr.	$\rho = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$	$r = \frac{s_{xy}}{s_x s_y}$

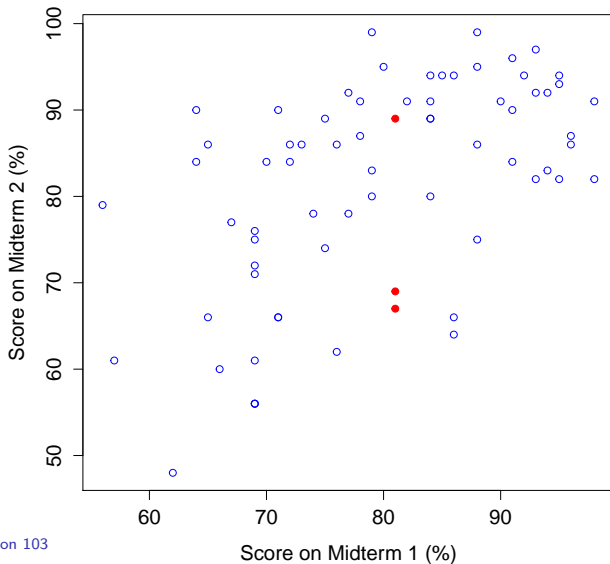
# Lecture #4 – Linear Regression I

Overview / Intuition for Linear Regression

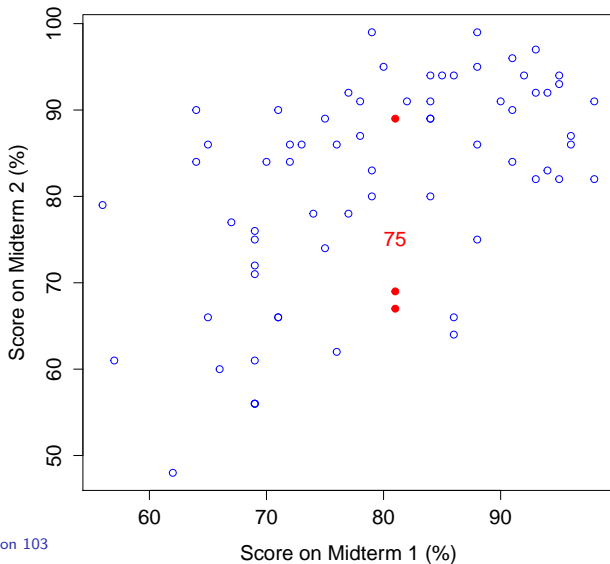
Deriving the Regression Equations

Relating Regression, Covariance and Correlation

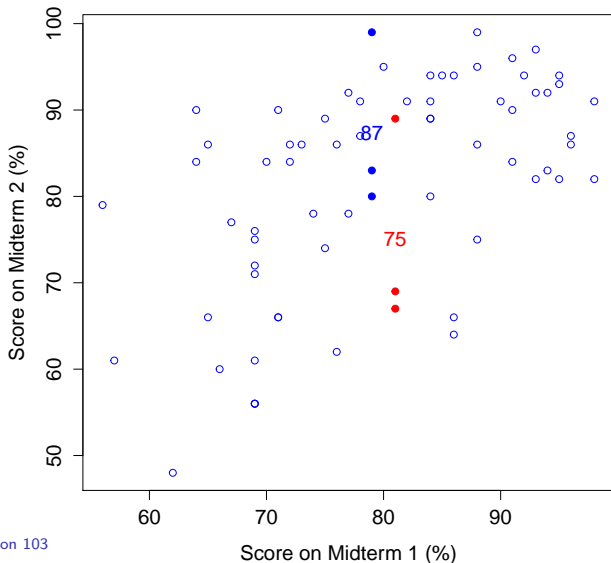
## Predict Second Midterm given 81 on First



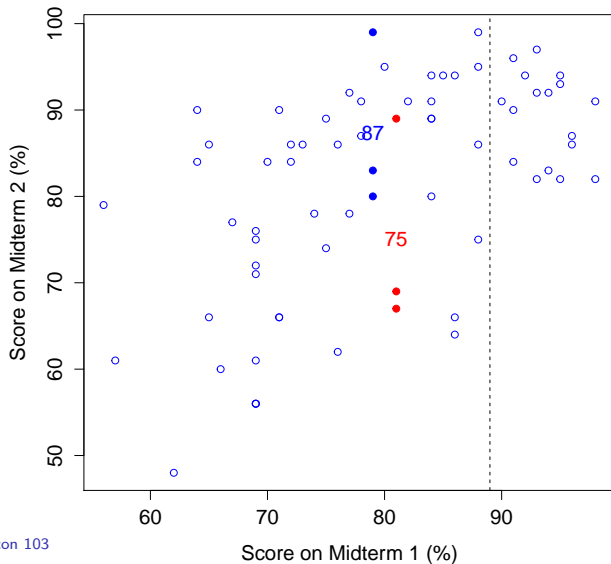
## Predict Second Midterm given 81 on First



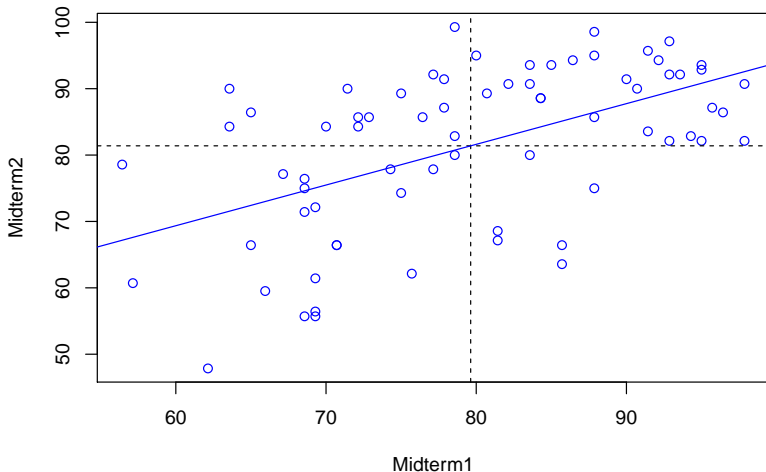
But if they'd only gotten 79 we'd predict higher?!



No one who took both exams got 89 on the first!

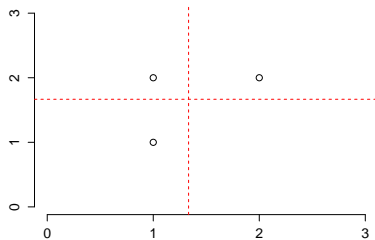
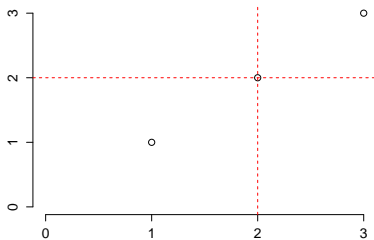


## Regression: “Best Fitting” Line Through Cloud of Points













# Least Squares Regression – Predict Using a Line

## The Prediction

Predict score  $\hat{y} = a + bx$  on 2nd midterm if you scored  $x$  on 1st

## How to choose $(a, b)$ ?

Linear regression chooses the slope ( $b$ ) and intercept ( $a$ ) that  
minimize the sum of squared vertical deviations

$$\sum_{i=1}^n d_i^2 = \sum_{i=1}^n (y_i - a - bx_i)^2$$

## Why Squared Deviations?

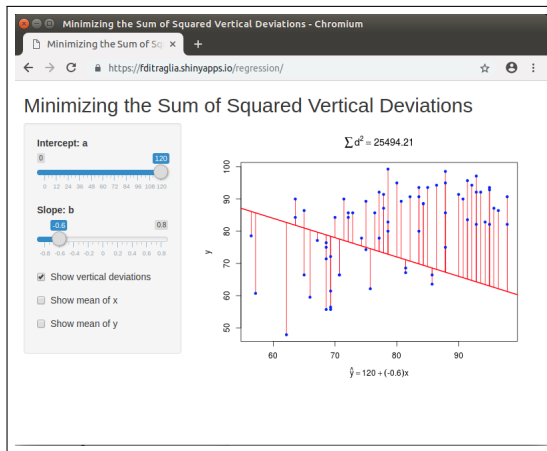
# Important Point About Notation

$$\underset{a,b}{\text{minimize}} \sum_{i=1}^n d_i^2 = \sum_{i=1}^n (y_i - a - bx_i)^2$$

$$\hat{y} = a + bx$$

- ▶  $(a, b)$  are our **choice variables**
- ▶  $(x_1, y_1), \dots, (x_n, y_n)$  are the **observed data**
- ▶  $\hat{y}$  is our **prediction** for a given value of  $x$
- ▶ Neither  $x$  nor  $\hat{y}$  needs to be in our dataset!

<https://fditraglia.shinyapps.io/regression/>



Try choosing  $(a, b)$  to minimize the sum of squared vertical deviations. . .

# Running the Regression in R

```
# Read data  
data_url <- 'http://ditraglia.com/econ103/midterms.csv'  
exams <- read.csv(data_url)  
  
# Drop students who missed an exam  
exams <- na.omit(exams)  
  
# Run the regression and display the slope and intercept  
reg <- lm(Midterm2 ~ Midterm1, data = exams)  
coef(reg)  
  
## (Intercept)      Midterm1  
## 32.5745441    0.6130357
```



# Predicting Midterm 2 Given 89 on Midterm 1

```
# By hand
```

```
32.5745441 + 0.6130357 * 89
```

```
## [1] 87.13472
```

```
# Using predict()
```

```
missing_student <- data.frame(Midterm1 = 89)
```

```
predict(reg, newdata = missing_student)
```

```
##          1
```

```
## 87.13472
```

# You Need to Know How To Derive This



Minimize the sum of squared vertical deviations from the line:

$$\min_{a,b} \sum_{i=1}^n (y_i - a - bx_i)^2$$

How should we proceed?

- (a) Differentiate with respect to  $x$
- (b) Differentiate with respect to  $y$
- (c) Differentiate with respect to  $x, y$
- (d) Differentiate with respect to  $a, b$
- (e) Can't solve this with calculus.

## Objective Function

$$\min_{a,b} \sum_{i=1}^n (y_i - a - bx_i)^2$$

FOC with respect to  $a$

$$-2 \sum_{i=1}^n (y_i - a - bx_i) = 0$$

$$\sum_{i=1}^n y_i - \sum_{i=1}^n a - b \sum_{i=1}^n x_i = 0$$

$$\frac{1}{n} \sum_{i=1}^n y_i - \frac{na}{n} - \frac{b}{n} \sum_{i=1}^n x_i = 0$$

$$\bar{y} - a - b\bar{x} = 0$$

## Regression Line Goes Through the Means!

$$\bar{y} = a + b\bar{x}$$

If your score equaled the class average on Midterm #1, we predict that your score will equal the class average on Midterm #2.

Substitute  $a = \bar{y} - b\bar{x}$

$$\begin{aligned}\sum_{i=1}^n (y_i - a - bx_i)^2 &= \sum_{i=1}^n (y_i - \bar{y} + b\bar{x} - bx_i)^2 \\ &= \sum_{i=1}^n [(y_i - \bar{y}) - b(x_i - \bar{x})]^2\end{aligned}$$

FOC wrt  $b$

$$-2 \sum_{i=1}^n [(y_i - \bar{y}) - b(x_i - \bar{x})] (x_i - \bar{x}) = 0$$

$$\sum_{i=1}^n (y_i - \bar{y}) (x_i - \bar{x}) - b \sum_{i=1}^n (x_i - \bar{x})^2 = 0$$

$$b = \frac{\sum_{i=1}^n (y_i - \bar{y}) (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

# Simple Linear Regression

## Problem

$$\min_{a,b} \sum_{i=1}^n (y_i - a - bx_i)^2$$

## Solution

$$b = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$a = \bar{y} - b\bar{x}$$

## Relating Regression to Covariance and Correlation

$$b = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} = \frac{s_{xy}}{s_x^2}$$

$$r = \frac{s_{xy}}{s_x s_y} = b \frac{s_x}{s_y}$$

# Comparing Regression, Correlation and Covariance

## Units

Correlation is unitless, covariance and regression coefficients ( $a$ ,  $b$ ) are not. (What are the units of these?)

## Symmetry

Correlation and covariance are symmetric, regression isn't.  
(Switching  $x$  and  $y$   $a$  and  $b$ : Review Exercise.)

## Extension Problem

Regression with z-scores rather than raw data gives  $a = 0$ ,  $b = r_{xy}$





$$s_{xy} = 6, \quad s_x = 5, \quad s_y = 2, \quad \bar{x} = 68, \quad \bar{y} = 21$$

What is the sample correlation between height ( $x$ ) and handspan ( $y$ )?



$$r = \frac{s_{xy}}{s_x s_y} = \frac{6}{5 \times 2} = 0.6$$

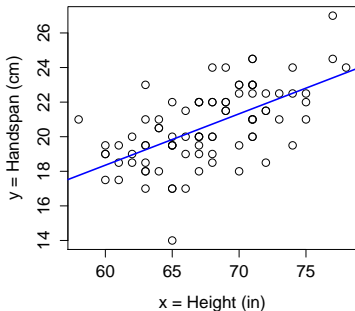


$$s_{xy} = 6, \quad s_x = 5, \quad s_y = 2, \quad \bar{x} = 68, \quad \bar{y} = 21$$

What is the value of  $b$  for the regression:

$$\hat{y} = a + bx$$

where  $x$  is height and  $y$  is handspan?



$$b = \frac{s_{xy}}{s_x^2} = \frac{6}{5^2} = 6/25 = 0.24$$



$$s_{xy} = 6, \quad s_x = 5, \quad s_y = 2, \quad \bar{x} = 68, \quad \bar{y} = 21$$

What is the value of  $a$  for the regression:

$$\hat{y} = a + bx$$

where  $x$  is height and  $y$  is handspan?  
(prev. slide  $b = 0.24$ )



$$a = \bar{y} - b\bar{x} = 21 - 0.24 \times 68 = 4.68$$

```
x <- seq(from = -1, to = 1, by = 0.1)
y <- x^2
cor(x,y)

## [1] 1.216307e-16

plot(x,y); abline(lm(y ~ x))
```



# Extremely Important Points to Remember!

- ▶ Regression, covariance, and correlation are all **measures of linear dependence**.
- ▶ Linear dependence **need not** imply a causal relationship.
- ▶ Dependence could be non-linear: always plot your data!

# Lecture #5 – Basic Probability I

Probability as Long-run Relative Frequency

Sets, Events and Axioms of Probability

“Classical” Probability

# Our Definition of Probability for this Course

Probability = Long-run Relative Frequency

That is, relative frequencies settle down to probabilities if we carry out an experiment over, and over, and over...

# Rolling a Fair, Six-Sided Die in R

```
# Function to plot relative frequencies
plot_freq <- function(x){
  n <- length(x)
  rel_freq <- prop.table(table(x))
  plot(rel_freq, ylab = 'Relative Frequency',
       xlab = bquote(n == .(n)))
}

# Roll a fair die 1 Million times
set.seed(1234567890)
dice <- sample(1:6, size = 1e6, replace = TRUE)
```



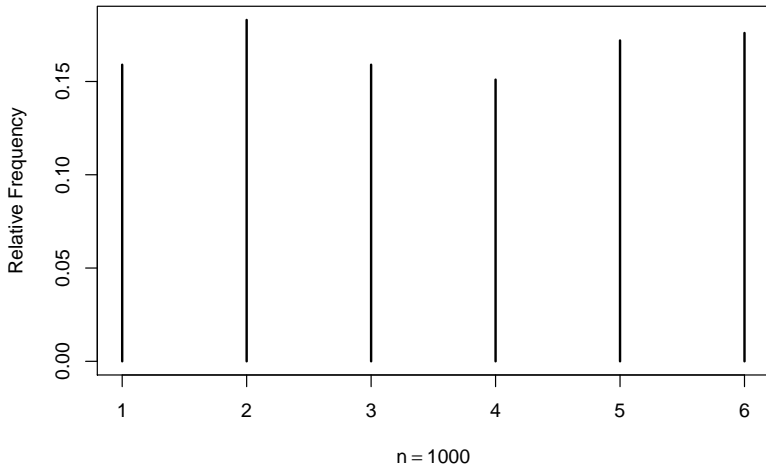
```
plot_freq(dice[1:10])
```



```
plot_freq(dice[1:50])
```



```
plot_freq(dice[1:1000])
```



```
plot_freq(dice)
```



# What do you think of this argument?



- ▶ The probability of flipping heads is  $1/2$ : if we flip a coin many times, about half of the time it will come up heads.
- ▶ The last ten throws in a row the coin has come up heads.
- ▶ The coin is bound to come up tails next time – it would be very rare to get 11 heads in a row.

(a) Agree

(b) Disagree

# The Gambler's Fallacy

Relative frequencies settle down to probabilities, but this does not mean that the trials are dependent.

Dependent = “Memory” of Prev. Trials

Independent = No “Memory” of Prev. Trials

# Terminology

## Random Experiment

An experiment whose outcomes are random.

## Basic Outcomes

Possible outcomes (mutually exclusive) of random experiment.

## Sample Space: $S$

Set of all basic outcomes of a random experiment.

## Event: $E$

A subset of the Sample Space (i.e. a collection of basic outcomes).

In set notation we write  $E \subseteq S$ .

# Example

## Random Experiment

Tossing a pair of dice.

## Basic Outcome

An ordered pair  $(a, b)$  where  $a, b \in \{1, 2, 3, 4, 5, 6\}$ , e.g.  $(2, 5)$

## Sample Space: $S$

All ordered pairs  $(a, b)$  where  $a, b \in \{1, 2, 3, 4, 5, 6\}$

Event:  $E = \{\text{Sum of two dice is less than 4}\}$

$\{(1, 1), (1, 2), (2, 1)\}$



## Visual Representation



The event  $E$  contains the basic outcomes  $O_3$  and  $O_2$  but not  $O_1$ .

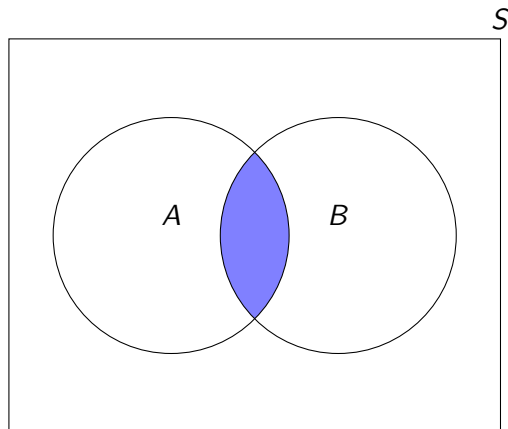
# Probability is Defined on *Sets*, and Events are Sets

## Complement of an Event: $A^c = \text{not } A$



**Figure:** The complement  $A^c$  of an event  $A \subseteq S$  is the collection of all basic outcomes from  $S$  not contained in  $A$ .

## Intersection of Events: $A \cap B = A \text{ and } B$



**Figure:** The intersection  $A \cap B$  of two events  $A, B \subseteq S$  is the collection of all basic outcomes from  $S$  contained in both  $A$  and  $B$

## Union of Events: $A \cup B = A \text{ or } B$



**Figure:** The union  $A \cup B$  of two events  $A, B \subseteq S$  is the collection of all basic outcomes from  $S$  contained in  $A$ ,  $B$  or both.

# Mutually Exclusive and Collectively Exhaustive

## Mutually Exclusive Events

A collection of events  $E_1, E_2, E_3, \dots$  is *mutually exclusive* if the intersection  $E_i \cap E_j$  of *any two different events* is empty.

## Collectively Exhaustive Events

A collection of events  $E_1, E_2, E_3, \dots$  is *collectively exhaustive* if, taken together, they contain *all of the basic outcomes in  $S$* .

Another way of saying this is that the union  $E_1 \cup E_2 \cup E_3 \cup \dots$  is  $S$ .

# Implications

## Mutually Exclusive Events

If one of the events occurs, then none of the others did.

## Collectively Exhaustive Events

One of these events *must* occur.

## Mutually Exclusive but *not Collectively Exhaustive*



Figure: Although  $A$  and  $B$  don't overlap, they also don't cover  $S$ .



## Collectively Exhaustive but *not Mutually Exclusive*



Figure: Together  $A$ ,  $B$ ,  $C$  and  $D$  cover  $S$ , but  $D$  overlaps with  $B$  and  $C$ .

## Collectively Exhaustive *and* Mutually Exclusive

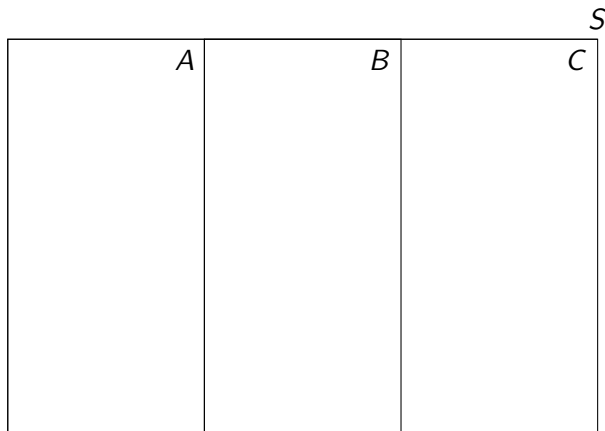


Figure:  $A$ ,  $B$ , and  $C$  cover  $S$  and don't overlap.

# Axioms of Probability

We assign every event  $A$  in the sample space  $S$  a real number  $P(A)$  called the **probability of  $A$**  such that:

Axiom 1  $0 \leq P(A) \leq 1$

Axiom 2  $P(S) = 1$

Axiom 3 If  $A_1, A_2, A_3, \dots$  are mutually exclusive events, then

$$P(A_1 \cup A_2 \cup A_3 \cup \dots) = P(A_1) + P(A_2) + P(A_3) + \dots$$

## “Classical” Probability

When all of the basic outcomes are equally likely, calculating the probability of an event is simply a matter of counting – count up all the basic outcomes that make up the event, and divide by the total number of basic outcomes.

# Recall from High School Math:

## Multiplication Rule for Counting

$n_1$  ways to make first decision,  $n_2$  ways to make second,  $\dots$ ,  $n_k$  ways to make  $k$ th  $\Rightarrow n_1 \times n_2 \times \dots \times n_k$  total ways to decide.

## Corollary – Number of Possible Orderings

$$k \times (k-1) \times (k-2) \times \dots \times 2 \times 1 = k!$$

## Permutations – Order $n$ people in $k$ slots

$$P_k^n = \frac{n!}{(n-k)!} \quad \text{(Order Matters)}$$

## Combinations – Choose committee of $k$ from group of $n$

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}, \text{ where } 0! = 1 \quad \text{(Order Doesn't Matter)}$$

# Poker – Deal 5 Cards, Order Doesn't Matter

## Basic Outcomes

$\binom{52}{5}$  possible hands

## How Many Hands have Four Aces?



48 (# of ways to choose the single card that is not an ace)

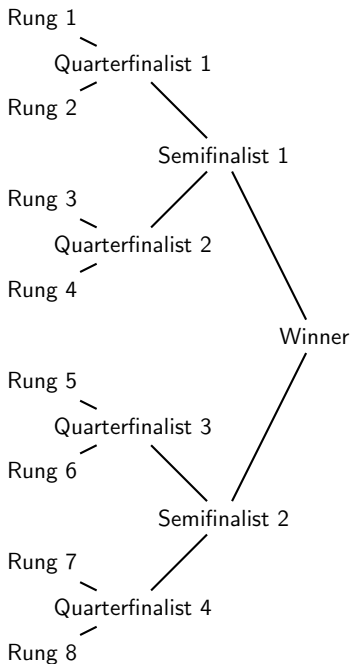
## Probability of Getting Four Aces

$$48 / \binom{52}{5} \approx 0.00002$$

## A Fairly Ridiculous Example



Roger Federer and Novak Djokovic have agreed to play in a tennis tournament against six Penn professors. Each player in the tournament is randomly allocated to one of the eight rungs in the ladder (next slide). Federer always beats Djokovic and, naturally, either of the two pros always beats any of the professors. What is the probability that Djokovic gets second place in the tournament?





# Solution: Order Matters!

## Denominator

8! basic outcomes – ways to arrange players on tournament ladder.

## Numerator

Sequence of three decisions:

1. Which rung to put Federer on? (8 possibilities)
2. Which rung to put Djokovic on?
  - For any given rung that Federer is on, only 4 rungs prevent Djokovic from meeting him until the final.
3. How to arrange the professors? (6! ways)

$$\frac{8 \times 4 \times 6!}{8!} = \frac{8 \times 4}{7 \times 8} = 4/7 \approx 0.57$$

# Lecture #6 – Basic Probability II

Complement Rule, Logical Consequence Rule, Addition Rule

Conditional Probability

Independence, Multiplication Rule

Law of Total Probability

## Recall: Axioms of Probability

Let  $S$  be the sample space. With each event  $A \subseteq S$  we associate a real number  $P(A)$  called the **probability of  $A$** , satisfying the following conditions:

**Axiom 1**  $0 \leq P(A) \leq 1$

**Axiom 2**  $P(S) = 1$

**Axiom 3** If  $A_1, A_2, A_3, \dots$  are mutually exclusive events, then

$$P(A_1 \cup A_2 \cup A_3 \cup \dots) = P(A_1) + P(A_2) + P(A_3) + \dots$$

## The Complement Rule: $P(A^c) = 1 - P(A)$

Since  $A, A^c$  are mutually exclusive and collectively exhaustive:

$$P(A \cup A^c) = P(A) + P(A^c) = P(S) = 1$$

Rearranging:

$$P(A^c) = 1 - P(A)$$



Figure:  $A \cap A^c = \emptyset$ ,  
 $A \cup A^c = S$

## Another Important Rule – Equivalent Events

If A and B are Logically Equivalent, then  $P(A) = P(B)$ .

In other words, if A and B contain exactly the same basic outcomes, then  $P(A) = P(B)$ .

Although this seems obvious it's important to keep in mind. . .

# The Logical Consequence Rule

If  $B$  Logically Entails  $A$ , then  $P(B) \leq P(A)$

For example, the probability that someone comes from Texas cannot exceed the probability that she comes from the USA.

In Set Notation

$$B \subseteq A \Rightarrow P(B) \leq P(A)$$

Why is this so?

If  $B \subseteq A$ , then all the basic outcomes in  $B$  are also in  $A$ .

## Proof of Logical Consequence Rule

Since  $B \subseteq A$ , we have  $B = A \cap B$  and  $A = B \cup (A \cap B^c)$ . Combining these,

$$A = (A \cap B) \cup (A \cap B^c)$$

Now since  $(A \cap B) \cap (A \cap B^c) = \emptyset$ ,

$$\begin{aligned} P(A) &= P(A \cap B) + P(A \cap B^c) \\ &= P(B) + P(A \cap B^c) \\ &\geq P(B) \end{aligned}$$

because  $0 \leq P(A \cap B^c) \leq 1$ .

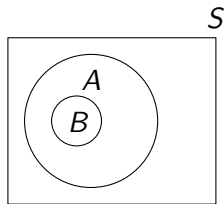


Figure:

$B = A \cap B$ , and  
 $A = B \cup (A \cap B^c)$

## “Odd Question” # 2

Pia is thirty-one years old, single, outspoken, and smart. She was a philosophy major. When a student, she was an ardent supporter of Native American rights, and she picketed a department store that had no facilities for nursing mothers. Rank the following statements in order from most probable to least probable.

- (A) Pia is an active feminist.
- (B) Pia is a bank teller.
- (C) Pia works in a small bookstore.
- (D) Pia is a bank teller and an active feminist.
- (E) Pia is a bank teller and an active feminist who takes yoga classes.
- (F) Pia works in a small bookstore and is an active feminist who takes yoga classes.



## Using the Logical Consequence Rule...

- (A) Pia is an active feminist.
- (B) Pia is a bank teller.
- (C) Pia works in a small bookstore.
- (D) Pia is a bank teller and an active feminist.
- (E) Pia is a bank teller and an active feminist who takes yoga classes.
- (F) Pia works in a small bookstore and is an active feminist who takes yoga classes.

Any Correct Ranking Must Satisfy:

$$P(A) \geq P(D) \geq P(E)$$

$$P(B) \geq P(D) \geq P(E)$$

$$P(A) \geq P(F)$$

$$P(C) \geq P(F)$$

# Throw a Fair Die Once

$E$  = roll an even number

What are the basic outcomes?

$\{1, 2, 3, 4, 5, 6\}$

What is  $P(E)$ ?



$E = \{2, 4, 6\}$  and the basic outcomes are equally likely (and mutually exclusive), so

$$P(E) = 1/6 + 1/6 + 1/6 = 3/6 = 1/2$$

# Throw a Fair Die Once

$E$  = roll an even number

$M$  = roll a 1 or a prime number

What is  $P(E \cup M)$ ?



Key point:  $E$  and  $M$  are not mutually exclusive!

$$P(E \cup M) = P(\{1, 2, 3, 4, 5, 6\}) = 1$$

$$P(E) = P(\{2, 4, 6\}) = 1/2$$

$$P(M) = P(\{1, 2, 3, 5\}) = 4/6 = 2/3$$

$$P(E) + P(M) = 1/2 + 2/3 = 7/6 \neq P(E \cup M) = 1$$

## The Addition Rule – Don't Double-Count!

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$



Construct a formal proof as an optional homework problem.

Who's on the other side?

# Three Cards, Each with a Face on the Front and Back



1. Gaga/Gaga
2. Obama/Gaga
3. Obama/Obama

I draw a card at random and look at one side: it's Obama.  
What is the probability that the other side is also Obama?



# Let's Try The Method of Monte Carlo...

When you don't know how to calculate, simulate.

## Procedure

1. Close your eyes and thoroughly shuffle your cards.
2. Keeping eyes closed, draw a card and place it on your desk.
3. Stand if Obama is face-up on your chosen card.
4. We'll count those standing and call the total  $N$
5. Of those standing, sit down if Obama is *not* on the back of your chosen card.
6. We'll count those *still* standing and call the total  $m$ .

$$\text{Monte Carlo Approximation of Desired Probability} = \frac{m}{N}$$

```
draw_simulation <- function() {  
  cards <- c('GG', 'OG', 'OO')  
  random_card <- sample(cards, size = 1)  
  if(random_card == 'GG') {  
    faces <- c('G', 'G')  
  } else if (random_card == 'OO') {  
    faces <- c('O', 'O')  
  } else {  
    faces <- c('O', 'G')  
  }  
  out <- sample(faces)  
  names(out) <- c('front', 'back')  
  return(out)  
}
```



```
set.seed(54321)
simulations <- replicate(n = 1000, draw_simulation())
simulations <- data.frame(t(simulations))
head(simulations)
```

```
##      front back
## 1         0    G
## 2         G    G
## 3         G    G
## 4         0    G
## 5         G    G
## 6         0    0
```

```
Obama_on_front <- subset(simulations, front == '0')
mean(Obama_on_front$back == '0')
```

```
## [1] 0.6633065
```



# Conditional Probability – Reduced Sample Space

Set of relevant outcomes restricted by condition

$$P(A|B) = \frac{P(A \cap B)}{P(B)}, \text{ provided } P(B) > 0$$



**Figure:**  $B$  becomes the “new sample space” so we need to re-scale by  $P(B)$  to keep probabilities between zero and one.

## Who's on the other side?

Let  $F$  be the event that Obama is on the front of the card of the card we draw and  $B$  be the event that he is on the back.

$$P(B|F) = \frac{P(B \cap F)}{P(F)} = \frac{1/3}{1/2} = 2/3$$

## Conditional Versions of Probability Axioms

1.  $0 \leq P(A|B) \leq 1$
2.  $P(B|B) = 1$
3. If  $A_1, A_2, A_3, \dots$  are mutually exclusive given  $B$ , then
$$P(A_1 \cup A_2 \cup A_3 \cup \dots | B) = P(A_1|B) + P(A_2|B) + P(A_3|B) \dots$$

## Conditional Versions of Other Probability Rules

- ▶  $P(A|B) = 1 - P(A^c|B)$
- ▶  $A_1$  logically equivalent to  $A_2 \iff P(A_1|B) = P(A_2|B)$
- ▶  $A_1 \subseteq A_2 \implies P(A_1|B) \leq P(A_2|B)$
- ▶  $P(A_1 \cup A_2|B) = P(A_1|B) + P(A_2|B) - P(A_1 \cap A_2|B)$

However:  $P(A|B) \neq P(B|A)$  and  $P(A|B^c) \neq 1 - P(A|B)$ !

# Independence and The Multiplication Rule

## The Multiplication Rule

Rearrange the definition of conditional probability:

$$P(A \cap B) = P(A|B)P(B)$$

## Statistical Independence

$$P(A \cap B) = P(A)P(B)$$

## By the Multiplication Rule

$$\text{Independence} \iff P(A|B) = P(A)$$

## Interpreting Independence

Knowledge that  $B$  has occurred tells nothing about whether  $A$  will.

## Will Having 5 Children Guarantee a Boy?



A couple plans to have five children. Assuming that each birth is independent and male and female children are equally likely, what is the probability that they have at least one boy?

By Independence and the Complement Rule,

$$\begin{aligned}P(\text{no boys}) &= P(5 \text{ girls}) \\&= 1/2 \times 1/2 \times 1/2 \times 1/2 \times 1/2 \\&= 1/32\end{aligned}$$

$$\begin{aligned}P(\text{at least 1 boy}) &= 1 - P(\text{no boys}) \\&= 1 - 1/32 = 31/32 = 0.97\end{aligned}$$

# The Law of Total Probability

If  $E_1, E_2, \dots, E_k$  are mutually exclusive, collectively exhaustive events and  $A$  is another event, then

$$P(A) = P(A|E_1)P(E_1) + P(A|E_2)P(E_2) + \dots + P(A|E_k)P(E_k)$$



## Example of Law of Total Probability

Define the following events:

$F$  = Obama on front of card

$A$  = Draw card with two Gagas

$B$  = Draw card with two Obamas

$C$  = Draw card with BOTH Obama and Gaga

$$\begin{aligned}P(F) &= P(F|A)P(A) + P(F|B)P(B) + P(F|C)P(C) \\&= 0 \times 1/3 + 1 \times 1/3 + 1/2 \times 1/3 \\&= 1/2\end{aligned}$$

## Deriving the Law of Total Probability For $k = 2$

Since  $A \cap B$  and  $A \cap B^c$  are mutually exclusive and their union equals  $A$ ,

$$P(A) = P(A \cap B) + P(A \cap B^c)$$

But by the multiplication rule:

$$P(A \cap B) = P(A|B)P(B)$$

$$P(A \cap B^c) = P(A|B^c)P(B^c)$$

Combining,

$$P(A) = P(A|B)P(B) + P(A|B^c)P(B^c)$$

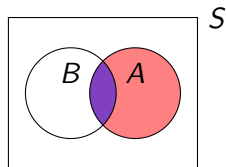


Figure:

$$A = (A \cap B) \cup (A \cap B^c), \\ (A \cap B) \cap (A \cap B^c) = \emptyset$$

# Lecture #7 – Basic Probability III / Discrete RVs I

Bayes' Rule and the Base Rate Fallacy

Overview of Random Variables

Probability Mass Functions

# Four Volunteers Please!

# The Lie Detector Problem

From accounting records, we know that 10% of employees in the store are stealing merchandise.

The managers want to fire the thieves, but their only tool in distinguishing is a lie detector test that is 80% accurate:

Innocent  $\Rightarrow$  Pass test with 80% Probability

Thief  $\Rightarrow$  Fail test with 80% Probability

What is the probability that someone is a thief *given* that she has failed the lie detector test?



# Monte Carlo Simulation – Roll a 10-sided Die Twice

Managers will split up and visit employees. Employees roll the die twice **but keep the results secret!**

## First Roll – Thief or not?

0  $\Rightarrow$  Thief, 1 – 9  $\Rightarrow$  Innocent

## Second Roll – Lie Detector Test

0, 1  $\Rightarrow$  Incorrect Test Result, 2 – 9 Correct Test Result

	0 or 1	2–9
Thief	Pass	<b>Fail</b>
Innocent	<b>Fail</b>	Pass

What percentage of those who failed the test are guilty?

# Who Failed Lie Detector Test:

# Of Thieves Among Those Who Failed:

```
draw_simulation <- function() {  
  guilty <- FALSE  
  fail <- FALSE  
  die1 <- sample(0:9, size = 1)  
  die2 <- sample(0:9, size = 1)  
  if(die1 == 0){ # Thief  
    guilty <- TRUE  
    if(die2 >= 2) fail <- TRUE  
  } else { # Innocent  
    if(die2 < 2) fail <- TRUE  
  }  
  return(c(guilty = guilty, fail = fail))  
}
```



```
set.seed(123456)
simulations <- replicate(n = 1000, draw_simulation())
simulations <- data.frame(t(simulations))
head(simulations)

##    guilty  fail
## 1  FALSE FALSE
## 2  FALSE FALSE
## 3  FALSE  TRUE
## 4  FALSE  TRUE
## 5  FALSE  TRUE
## 6  FALSE FALSE

failed_test <- subset(simulations, fail)
mean(failed_test$guilty)

## [1] 0.311828
```

# Base Rate Fallacy – Failure to Consider Prior Information

## Base Rate – Prior Information

Before the test we know that 10% of Employees are stealing.

People tend to focus on the fact that the test is 80% accurate and ignore the fact that only 10% of the employees are thieves.

# Thief (Y/N), Lie Detector (P/F)

	0	1	2	3	4	5	6	7	8	9
0	YP	YP	YF	YF	YF	YF	YF	YF	YF	YF
1	NF	NF	NP	NP	NP	NP	NP	NP	NP	NP
2	NF	NF	NP	NP	NP	NP	NP	NP	NP	NP
3	NF	NF	NP	NP	NP	NP	NP	NP	NP	NP
4	NF	NF	NP	NP	NP	NP	NP	NP	NP	NP
5	NF	NF	NP	NP	NP	NP	NP	NP	NP	NP
6	NF	NF	NP	NP	NP	NP	NP	NP	NP	NP
7	NF	NF	NP	NP	NP	NP	NP	NP	NP	NP
8	NF	NF	NP	NP	NP	NP	NP	NP	NP	NP
9	NF	NF	NP	NP	NP	NP	NP	NP	NP	NP

**Table:** Each outcome in the table is equally likely. The 26 given in red correspond to failing the test, but only 8 of these (YF) correspond to being a thief.

## Base Rate of Thievery is 10%



**Figure:** Although  $\frac{9}{50} + \frac{4}{50} = \frac{13}{50}$  fail the test, only  $\frac{4/50}{13/50} = \frac{4}{13} \approx 0.31$  are actually thieves!

## Deriving Bayes' Rule

Intersection is symmetric:  $A \cap B = B \cap A$  so  $P(A \cap B) = P(B \cap A)$

By the definition of conditional probability,

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

And by the multiplication rule:

$$P(B \cap A) = P(B|A)P(A)$$

Finally, combining these

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

# Understanding Bayes' Rule

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

## Reversing the Conditioning

Express  $P(A|B)$  in terms of  $P(B|A)$ . *Relative magnitudes* of the two conditional probabilities determined by the ratio  $P(A)/P(B)$ .

## Base Rate

$P(A)$  is called the “base rate” or the “prior probability.”

## Denominator

Typically, we calculate  $P(B)$  using the law of total probability

In General  $P(A|B) \neq P(B|A)$



### Question

Most college students are Democrats. Does it follow that most Democrats are college students? (A = YES, B = NO)

### Answer

There are many more Democrats than college students:

$$P(\text{Dem}) > P(\text{Student})$$

so  $P(\text{Student}|\text{Dem})$  is small even though  $P(\text{Dem}|\text{Student})$  is large.

## Solving the Lie Detector Problem with Bayes' Rule

$T$  = Employee is a Thief,  $F$  = Employee Fails Lie Detector Test

$$P(T|F) = \frac{P(F|T)P(T)}{P(F)}$$

$$\begin{aligned}P(F) &= P(F|T)P(T) + P(F|T^c)P(T^c) \\&= 0.8 \times 0.1 + 0.2 \times 0.9 \\&= 0.08 + 0.18 = 0.26\end{aligned}$$

$$P(T|F) = \frac{0.08}{0.26} = \frac{8}{26} = \frac{4}{13} \approx 0.31$$



# Random Variables

# Random Variables

*A random variable is neither random nor a variable.*

## Random Variable (RV): $X$

A *fixed* function that assigns a *number* to each basic outcome of a random experiment.

## Realization: $x$

A particular numeric value that an RV could take on. We write  $\{X = x\}$  to refer to the *event* that the RV  $X$  took on the value  $x$ .

## Support Set (aka Support)

The set of all possible realizations of a RV.

# Random Variables (continued)

## Notation

Capital latin letters for RVs, e.g.  $X$ ,  $Y$ ,  $Z$ , and the corresponding lowercase letters for their realizations, e.g.  $x$ ,  $y$ ,  $z$ .

## Intuition

A RV is machine that spits out random numbers. The machine is deterministic: outputs are random because *inputs* are random.

## Why Random Variables?

Different random experiments can have the same structure: e.g. flipping a fair coin vs. drawing a ball from an urn with 5 red and 5 blue. RVs abstract from coin vs. urn and let us study both at once.

## Example: Coin Flip Random Variable



**Figure:** This random variable assigns numeric values to the random experiment of flipping a fair coin once: Heads is assigned 1 and Tails 0.

Which of these is a realization of the Coin Flip RV?



- (a) Tails
- (b) 2
- (c) 0
- (d) Heads
- (e)  $1/2$

## What is the support set of the Coin Flip RV?



- (a) {Heads, Tails}
- (b)  $1/2$
- (c) 0
- (d)  $\{0, 1\}$
- (e) 1

Let  $X$  denote the Coin Flip RV



What is  $P(X = 1)$ ?

- (a) 0
- (b) 1
- (c)  $1/2$
- (d) Not enough information to determine

# Two Kinds of RVs: Discrete and Continuous

**Discrete** support set is discrete, e.g.  $\{0, 1, 2\}$ ,  
 $\{\dots, -2, -1, 0, 1, 2, \dots\}$

**Continuous** support set is continuous, e.g.  $[-1, 1]$ ,  $\mathbb{R}$ .

Start with the discrete case since it's easier, but most of the ideas we learn will carry over to the continuous case.



# Discrete Random Variables I

# Probability Mass Function (pmf)

A function that gives  $P(X = x)$  for any realization  $x$  in the support set of a discrete RV  $X$ . We use the following notation for the pmf:

$$p(x) = P(X = x)$$

Plug in a realization  $x$ , get out a probability  $p(x)$ .

# Probability Mass Function for Coin Flip RV

$$X = \begin{cases} 0, \text{Tails} \\ 1, \text{Heads} \end{cases}$$

$$p(0) = 1/2$$

$$p(1) = 1/2$$



Figure: Plot of pmf for Coin Flip Random Variable

## Important Note about Support Sets

Whenever you write down the pmf of a RV, it is **crucial** to also write down its Support Set. Recall that this is the set of *all possible realizations for a RV*. Outside of the support set, all probabilities are zero. In other words, the pmf is **only defined** on the support.

# Properties of Probability Mass Functions

If  $p(x)$  is the pmf of a random variable  $X$ , then

(i)  $0 \leq p(x) \leq 1$  for all  $x$

(ii)  $\sum_{\text{all } x} p(x) = 1$

where “all  $x$ ” is shorthand for “all  $x$  in the support of  $X$ .”

# Lecture #8 – Discrete RVs II

Cumulative Distribution Functions (CDFs)

The Bernoulli Random Variable

Definition of Expected Value

Expected Value of a Function

Linearity of Expectation

## Recall: Properties of Probability Mass Functions

If  $p(x)$  is the pmf of a random variable  $X$ , then

(i)  $0 \leq p(x) \leq 1$  for all  $x$

(ii)  $\sum_{\text{all } x} p(x) = 1$

where “all  $x$ ” is shorthand for “all  $x$  in the support of  $X$ .”

# Cumulative Distribution Function (CDF)

This Def. is **the same** for continuous RVs.

The CDF gives the probability that a RV  $X$  **does not exceed** a specified threshold  $x_0$ , as a function of  $x_0$

$$F(x_0) = P(X \leq x_0)$$

**Important!**

The threshold  $x_0$  is allowed to be *any real number*. In particular, it doesn't have to be in the support of  $X$ !



## Discrete RVs: Sum the pmf to get the CDF

$$F(x_0) = \sum_{x \leq x_0} p(x)$$

Why?

The events  $\{X = x\}$  are mutually exclusive, so we sum to get the probability of their union for all  $x \leq x_0$ :

$$F(x_0) = P(X \leq x_0) = P\left(\bigcup_{x \leq x_0} \{X = x\}\right) = \sum_{x \leq x_0} P(X = x) = \sum_{x \leq x_0} p(x)$$

## Probability Mass Function



$$p(0) = 1/2$$

$$p(1) = 1/2$$

## Cumulative Dist. Function



$$F(x_0) = \begin{cases} 0, & x_0 < 0 \\ \frac{1}{2}, & 0 \leq x_0 < 1 \\ 1, & x_0 \geq 1 \end{cases}$$

# Properties of CDFs

These are also true for continuous RVs.

1.  $\lim_{x_0 \rightarrow \infty} F(x_0) = 1$
2.  $\lim_{x_0 \rightarrow -\infty} F(x_0) = 0$
3. Non-decreasing:  $x_0 < x_1 \Rightarrow F(x_0) \leq F(x_1)$
4. Right-continuous (“open” versus “closed” on prev. slide)

Since  $F(x_0) = P(X \leq x_0)$ , we have  $0 \leq F(x_0) \leq 1$  for all  $x_0$

# Bernoulli Random Variable – Generalization of Coin Flip

## Support Set

$\{0, 1\}$  – 1 traditionally called “success,” 0 “failure”

## Probability Mass Function

$$p(0) = 1 - p$$

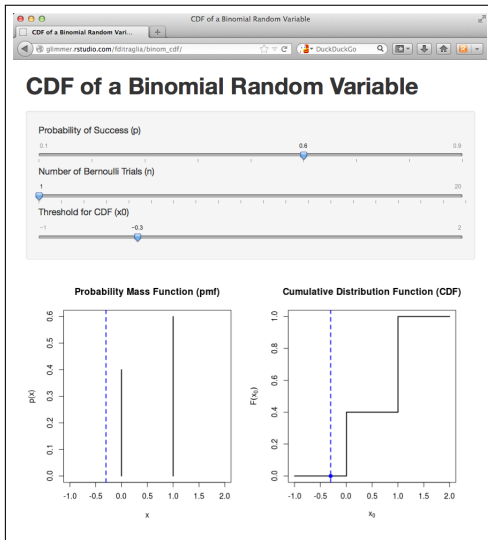
$$p(1) = p$$

## Cumulative Distribution Function

$$F(x_0) = \begin{cases} 0, & x_0 < 0 \\ 1 - p, & 0 \leq x_0 < 1 \\ 1, & x_0 \geq 1 \end{cases}$$

[http://fditraglia.shinyapps.io/binom\\_cdf/](http://fditraglia.shinyapps.io/binom_cdf/)

Set the second slider to 1 and play around with the others.



## Average Winnings Per Trial



If the realizations of the coin-flip RV were **payoffs**, how much would you expect to win per play *on average* in a long sequence of plays?

$$X = \begin{cases} \$0, \text{Tails} \\ \$1, \text{Heads} \end{cases}$$

## Expected Value (aka Expectation)

The expected value of a discrete RV  $X$  is given by

$$E[X] = \sum_{\text{all } x} x \cdot p(x)$$

In other words, the expected value of a discrete RV is the *probability-weighted average of its realizations*.

### Notation

We sometimes write  $\mu$  as shorthand for  $E[X]$ .

## Expected Value of Bernoulli RV

$$X = \begin{cases} 0, \text{Failure: } 1 - p \\ 1, \text{Success: } p \end{cases}$$

$$\sum_{\text{all } x} x \cdot p(x) = 0 \cdot (1 - p) + 1 \cdot p = p$$



## Your Turn to Calculate an Expected Value



Let  $X$  be a random variable with support set  $\{1, 2, 3\}$  where  $p(1) = p(2) = 1/3$ . Calculate  $E[X]$ .

$$E[X] = \sum_{\text{all } x} x \cdot p(x) = 1 \times 1/3 + 2 \times 1/3 + 3 \times 1/3 = 2$$

# Random Variables and Parameters

Notation:  $X \sim \text{Bernoulli}(p)$

Means  $X$  is a Bernoulli RV with  $P(X = 1) = p$  and  $P(X = 0) = 1 - p$ . The tilde is read “distributes as.”

Parameter

Any constant that appears in the definition of a RV, here  $p$ .

# Constants Versus Random Variables

This is a crucial distinction that students sometimes miss:

## Random Variables

- ▶ Suppose  $X$  is a RV – the values it takes on are random
- ▶ A function  $g(X)$  of a RV is itself a RV as we'll learn today.

## Constants

- ▶  $E[X]$  is a constant (you should convince yourself of this)
- ▶ Realizations  $x$  are constants. What is random is *which* realization the RV takes on.
- ▶ Parameters are constants (e.g.  $p$  for Bernoulli RV)
- ▶ Sample size  $n$  is a constant

# The St. Petersburg Game

## How Much Would You Pay?



How much would you be willing to pay for the right to play the following game?

*Imagine a fair coin. The coin is tossed once. If it falls heads, you receive a prize of \$2 and the game stops. If not, it is tossed again. If it falls heads on the second toss, you get \$4 and the game stops. If not, it is tossed again. If it falls heads on the third toss, you get \$8 and the game stops, and so on. The game stops after the first head is thrown. If the first head is thrown on the  $x^{\text{th}}$  toss, the prize is  $\$2^x$*

$X =$  Trial Number of First Head

$x$	$2^x$	$p(x)$	$2^x \cdot p(x)$
1	2	$1/2$	1
2	4	$1/4$	1
3	8	$1/8$	1
$\vdots$	$\vdots$	$\vdots$	$\vdots$
$n$	$2^n$	$1/2^n$	1
$\vdots$	$\vdots$	$\vdots$	$\vdots$

$$E[Y] = \sum_{\text{all } x} 2^x \cdot p(x) = 1 + 1 + 1 + \dots = \infty$$

# Functions of Random Variables are Themselves Random Variables

Example:  $X \sim \text{Bernoulli}(p)$ ,  $Y = (X + 1)^2$

Support Set for  $Y$

$$\{(0 + 1)^2, (1 + 1)^2\} = \{1, 4\}$$

Probability Mass Function for  $Y$

$$p_Y(y) = \begin{cases} 1 - p & y = 1 \\ p & y = 4 \\ 0 & \text{otherwise} \end{cases}$$

Expected Value of  $Y$

$$\sum_{y \in \{1, 4\}} y \times p_Y(y) = 1 \times (1 - p) + 4 \times p = 1 + 3p$$



Example:  $X \sim \text{Bernoulli}(p)$ ,  $Y = (X + 1)^2$

$$E[g(X)] = E[(X + 1)^2]$$

$$\sum_{y \in \{1,4\}} y \times p_Y(y) = 1 \times (1 - p) + 4 \times p = 1 + 3p$$

$$g(E[X]) = (E[X] + 1)^2$$

$$(E[X] + 1)^2 = (p + 1)^2 = 1 + 2p + p^2$$

In general:  $1 + 3p \neq 1 + 2p + p^2$ !

$$E[g(X)] \neq g(E[X])$$

(Expected value of Function  $\neq$  Function of Expected Value)

## Expectation of a Function of a Discrete RV

Let  $X$  be a random variable and  $g$  be a function. Then:

$$E[g(X)] = \sum_{\text{all } x} g(x)p(x)$$

This is how we proceeded in the St. Petersburg Game Example

## Your Turn: Calculate $E[X^2]$



$X$  has support  $\{-1, 0, 1\}$ ,  $p(-1) = p(0) = p(1) = 1/3$ .

$$\begin{aligned} E[X^2] &= \sum_{\text{all } x} x^2 p(x) = \sum_{x \in \{-1, 0, 1\}} x^2 p(x) \\ &= (-1)^2 \cdot (1/3) + (0)^2 \cdot (1/3) + (1)^2 \cdot (1/3) \\ &= 1/3 + 1/3 \\ &= 2/3 \approx 0.67 \end{aligned}$$

```
set.seed(794729)
sims <- sample(c(-1, 0, 1), size = 1e6, replace = TRUE,
               prob = c(1/3, 1/3, 1/3))
head(sims)

## [1]  1 -1  0  0  1  1

mean(sims)

## [1] -0.001182

mean(sims^2)

## [1] 0.66682
```

# Linearity of Expectation

Holds for Continuous RVs as well, but proof is different.

Let  $X$  be a RV and  $a, b$  be constants. Then:

$$E[a + bX] = a + bE[X]$$

## This is a Crucial Exception

In general  $E[g(X)]$  does not equal  $g(E[X])$ . But in the special case where  $g$  is a **linear function**,  $g(X) = a + bX$ , the two **are equal**.

## Example: Linearity of Expectation



Let  $X \sim \text{Bernoulli}(1/3)$  and define  $Y = 3X + 2$

1. What is  $E[X]$ ?  $E[X] = 0 \times 2/3 + 1 \times 1/3 = 1/3$
2. What is  $E[Y]$ ?  $E[Y] = E[3X + 2] = 3E[X] + 2 = 3$

## Proof: Linearity of Expectation For Discrete RV

$$\begin{aligned}E[a + bX] &= \sum_{\text{all } x} (a + bx)p(x) \\&= \sum_{\text{all } x} p(x) \cdot a + \sum_{\text{all } x} p(x) \cdot bx \\&= a \sum_{\text{all } x} p(x) + b \sum_{\text{all } x} x \cdot p(x) \\&= a + bE[X]\end{aligned}$$



# Lecture #9 – Discrete RVs III

Variance and Standard Deviation of a Random Variable

Binomial Random Variable

# Variance and Standard Deviation of a RV

The Defs are the same for continuous RVs, but the method of calculating will differ.

## Variance (Var)

$$\sigma^2 = \text{Var}(X) = E[(X - \mu)^2] = E[(X - E[X])^2]$$

## Standard Deviation (SD)

$$\sigma = \sqrt{\sigma^2} = \text{SD}(X)$$

# Key Point

Variance and std. dev. are *expectations of functions of a RV*

It follows that:

1. Variance and SD are constants
2. To derive facts about them you can use the facts you know about expected value

# How To Calculate Variance for Discrete RV?

Remember: it's just a function of  $X$ !

$$\text{Recall that } \mu = E[X] = \sum_{\text{all } x} xp(x)$$

$$\text{Var}(X) = E[(X - \mu)^2] = \sum_{\text{all } x} (x - \mu)^2 p(x)$$

## Shortcut Formula For Variance

This is *not* the definition, it's a shortcut for doing calculations:

$$\text{Var}(X) = E[(X - \mu)^2] = E[X^2] - (E[X])^2$$

We'll prove this in an upcoming lecture.

## Example: The Shortcut Formula



Let  $X \sim \text{Bernoulli}(1/2)$ . Calculate  $\text{Var}(X)$ .

$$E[X] = 0 \times 1/2 + 1 \times 1/2 = 1/2$$

$$E[X^2] = 0^2 \times 1/2 + 1^2 \times 1/2 = 1/2$$

$$E[X^2] - (E[X])^2 = 1/2 - (1/2)^2 = 1/4$$

## Variance of Bernoulli RV – via the Shortcut Formula

Step 1 –  $E[X]$

$$\mu = E[X] = \sum_{x \in \{0,1\}} p(x) \cdot x = (1-p) \cdot 0 + p \cdot 1 = p$$

Step 2 –  $E[X^2]$

$$E[X^2] = \sum_{x \in \{0,1\}} x^2 p(x) = 0^2(1-p) + 1^2 p = p$$

Step 3 – Combine with Shortcut Formula

$$\sigma^2 = \text{Var}[X] = E[X^2] - (E[X])^2 = p - p^2 = p(1-p)$$

## Variance of a Linear Transformation

$$\begin{aligned}\text{Var}(a + bX) &= E \left[ \{(a + bX) - E(a + bX)\}^2 \right] \\&= E \left[ \{(a + bX) - (a + bE[X])\}^2 \right] \\&= E \left[ (bX - bE[X])^2 \right] \\&= E[b^2(X - E[X])^2] \\&= b^2 E[(X - E[X])^2] \\&= b^2 \text{Var}(X) = b^2 \sigma^2\end{aligned}$$

The key point here is that variance is defined in terms of expectation and expectation is linear.



## Variance and SD are *NOT* Linear

$$\text{Var}(a + bX) = b^2\sigma^2$$

$$\text{SD}(a + bX) = |b|\sigma$$

These should look familiar from the related results for sample variance and std. dev. that you worked out on an earlier problem set.

# Binomial Random Variable

Let  $X$  = the sum of  $n$  independent Bernoulli trials, each with probability of success  $p$ . Then we say that:  $X \sim \text{Binomial}(n, p)$

## Parameters

$p$  = probability of “success,”  $n$  = # of trials

## Support

$\{0, 1, 2, \dots, n\}$

## Probability Mass Function (pmf)

$$p(x) = \binom{n}{x} p^x (1 - p)^{n-x}$$

[http://fditraglia.shinyapps.io/binom\\_cdf/](http://fditraglia.shinyapps.io/binom_cdf/)

Try playing around with all three sliders. If you set the second to 1 you get a Bernoulli.



# Where does the Binomial pmf come from?

## Question

Suppose we flip a fair coin 3 times. What is the probability that we get exactly 2 heads?

## Answer

Three basic outcomes make up this event:  $\{HHT, HTH, THH\}$ , each has probability  $1/8 = 1/2 \times 1/2 \times 1/2$ . Basic outcomes are mutually exclusive, so sum to get  $3/8 = 0.375$

# Where does the Binomial pmf come from?

## Question

Suppose we flip an *unfair* coin 3 times, where the probability of heads is  $1/3$ . What is the probability that we get exactly 2 heads?

## Answer

No longer true that *all* basic outcomes are equally likely, but those with exactly two heads *still are*

$$P(HHT) = (1/3)^2(1 - 1/3) = 2/27$$

$$P(THH) = 2/27$$

$$P(HTH) = 2/27$$

Summing gives  $2/9 \approx 0.22$

# Where does the Binomial pmf come from?

Starting to see a pattern?

Suppose we flip an unfair coin 4 times, where the probability of heads is  $1/3$ . What is the probability that we get exactly 2 heads?

HHTT    TTHH

HTHT    THTH

HTTH    THTT

Six equally likely, mutually exclusive  
basic outcomes make up this event:

$$\binom{4}{2} (1/3)^2 (2/3)^2$$

# R Commands for Binomial( $n, p$ ) RV

## Probability Mass Function

`dbinom(x, size, prob)`, where `size` is  $n$  and `prob` is  $p$

## Cumulative Distribution Function

`pbinom(q, size, prob)`, where `q` is  $x_0$ , `size` is  $n$  and `prob` is  $p$

## Make Random Draws

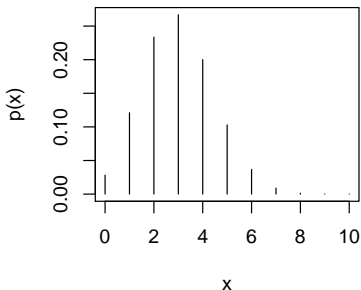
`rbinom(n, size, prob)`, where `n` is the number of draws, `size` is  $n$  and `prob` is  $p$

```

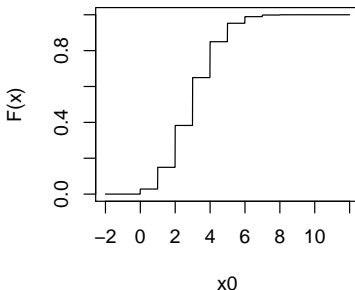
x <- 0:10
px <- dbinom(x, size = 10, prob = 0.3)
x0 <- seq(from = -2, to = 12, by = 0.01)
Fx <- pbinom(x0, size = 10, prob = 0.3)
par(mfrow = c(1, 2))
plot(x, px, type = 'h', ylab = 'p(x)', main = 'Binom(10, 0.3) pmf')
plot(x0, Fx, type = 'l', ylab = 'F(x)', main = 'Binom(10, 0.3) CDF')

```

**Binom(10, 0.3) pmf**



**Binom(10, 0.3) CDF**

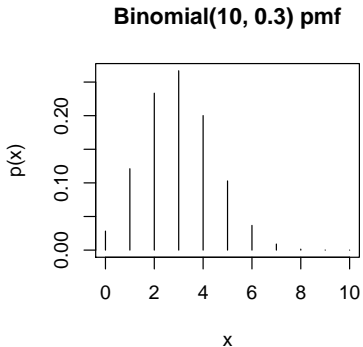
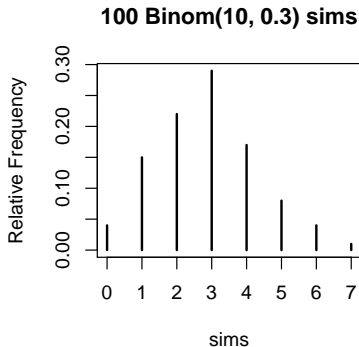




```

set.seed(5545)
sims <- rbinom(100, size = 10, prob = 0.3)
par(mfrow = c(1, 2))
rel_freq <- prop.table(table(sims))
plot(rel_freq, main = '100 Binom(10, 0.3) sims',
     ylab = 'Relative Frequency')
plot(x, px, type = 'h', ylab = 'p(x)', main = 'Binomial(10, 0.3) pmf')

```



# Lecture #10 – Discrete RVs IV

Joint vs. Marginal Probability Mass Functions

Conditional Probability Mass Function & Independence

Expectation of a Function of Two Discrete RVs, Covariance

Linearity of Expectation Reprise, Properties of Binomial RV

## Multiple RVs *at once* - Definition of Joint PMF

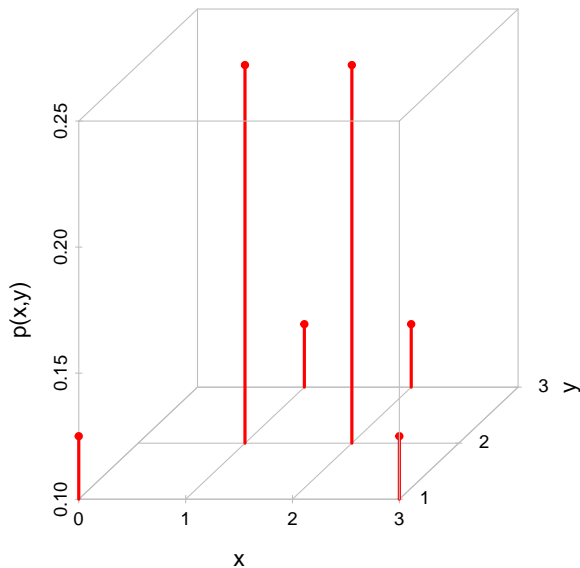
Let  $X$  and  $Y$  be discrete random variables. The joint probability mass function  $p_{XY}(x, y)$  gives the probability of each pair of realizations  $(x, y)$  in the support:

$$p_{XY}(x, y) = P(X = x \cap Y = y)$$

## Example: Joint PMF in Tabular Form

		Y		
		1	2	3
X	0	1/8	0	0
	1	0	1/4	1/8
	2	0	1/4	1/8
	3	1/8	0	0

# Plot of Joint PMF



What is  $p_{XY}(1, 2)$ ?



		Y		
		1	2	3
X	0	1/8	0	0
	1	0	1/4	1/8
	2	0	1/4	1/8
	3	1/8	0	0

$$p_{XY}(1, 2) = P(X = 1 \cap Y = 2) = 1/4$$

$$p_{XY}(2, 1) = P(X = 2 \cap Y = 1) = 0$$

# Properties of Joint PMF

1.  $0 \leq p_{XY}(x, y) \leq 1$  for any pair  $(x, y)$
2. The sum of  $p_{XY}(x, y)$  over all pairs  $(x, y)$  in the support is 1:

$$\sum_x \sum_y p(x, y) = 1$$

# Joint versus Marginal PMFs

## Joint PMF

$$p_{XY}(x, y) = P(X = x \cap Y = y)$$

## Marginal PMFs

$$p_X(x) = P(X = x)$$

$$p_Y(y) = P(Y = y)$$

You can't calculate a joint pmf from marginals alone but you *can* calculate marginals from the joint!



## Marginals from Joint

$$p_X(x) = \sum_{\text{all } y} p_{XY}(x, y)$$

$$p_Y(y) = \sum_{\text{all } x} p_{XY}(x, y)$$

Why?

$$\begin{aligned} p_Y(y) &= P(Y = y) = P\left(\bigcup_{\text{all } x} \{X = x \cap Y = y\}\right) \\ &= \sum_{\text{all } x} P(X = x \cap Y = y) = \sum_{\text{all } x} p_{XY}(x, y) \end{aligned}$$

To get the marginals sum “into the margins” of the table.

		Y			
		1	2	3	
X	0	1/8	0	0	1/8
	1	0	1/4	1/8	3/8
	2	0	1/4	1/8	3/8
	3	1/8	0	0	1/8
					1

$$p_X(0) = 1/8 + 0 + 0 = 1/8$$

$$p_X(1) = 0 + 1/4 + 1/8 = 3/8$$

$$p_X(2) = 0 + 1/4 + 1/8 = 3/8$$

$$p_X(3) = 1/8 + 0 + 0 = 1/8$$

What is  $p_Y(2)$ ?



		Y			
		1	2	3	
X	0	1/8	0	0	
	1	0	1/4	1/8	
	2	0	1/4	1/8	
	3	1/8	0	0	
		1/4	1/2	1/4	1

$$p_Y(1) = 1/8 + 0 + 0 + 1/8 = 1/4$$

$$p_Y(2) = 0 + 1/4 + 1/4 + 0 = 1/2$$

$$p_Y(3) = 0 + 1/8 + 1/8 + 0 = 1/4$$

# Definition of Conditional PMF

How does the distribution of  $y$  change with  $x$ ?

$$p_{Y|X}(y|x) = P(Y = y|X = x) = \frac{P(Y = y \cap X = x)}{P(X = x)} = \frac{p_{XY}(x, y)}{p_X(x)}$$

## Conditional PMF of $Y$ given $X = 2$

		$Y$			
		1	2	3	
$X$	0	1/8	0	0	1/8
	1	0	1/4	1/8	3/8
	2	0	1/4	1/8	3/8
	3	1/8	0	0	1/8

$$p_{Y|X}(1|2) = \frac{p_{XY}(2,1)}{p_X(2)} = \frac{0}{3/8} = 0$$

$$p_{Y|X}(2|2) = \frac{p_{XY}(2,2)}{p_X(2)} = \frac{1/4}{3/8} = 2/3$$

$$p_{Y|X}(3|2) = \frac{p_{XY}(2,3)}{p_X(2)} = \frac{1/8}{3/8} = 1/3$$

What is  $p_{X|Y}(1|2)$ ?



		Y			
		1	2	3	
X	0	1/8	0	0	
	1	0	1/4	1/8	
	2	0	1/4	1/8	
	3	1/8	0	0	
		1/4	1/2	1/4	

$$p_{X|Y}(1|2) = \frac{p_{XY}(1,2)}{p_Y(2)} = \frac{1/4}{1/2} = 1/2$$

Similarly:

$$p_{X|Y}(0|2) = 0, \quad p_{X|Y}(2|2) = 1/2, \quad p_{X|Y}(3|2) = 0$$

# Independent RVs: Joint Equals Product of Marginals

## Definition

Two discrete RVs are **independent** if and only if

$$p_{XY}(x, y) = p_X(x)p_Y(y)$$

for all pairs  $(x, y)$  in the support.

## Equivalent Definition

$$p_{Y|X}(y|x) = p_Y(y) \text{ and } p_{X|Y}(x|y) = p_X(x)$$

for all pairs  $(x, y)$  in the support.

# Are $X$ and $Y$ Independent?



(A = YES, B = NO)

		Y			
		1	2	3	
X	0	1/8	0	0	1/8
	1	0	1/4	1/8	3/8
	2	0	1/4	1/8	3/8
	3	1/8	0	0	1/8
		1/4	1/2	1/4	

$$p_{XY}(2, 1) = 0$$

$$p_X(2) \times p_Y(1) = (3/8) \times (1/4) \neq 0$$

Therefore  $X$  and  $Y$  are *not* independent.



## Expectation of Function of Two Discrete RVs

$$E[g(X, Y)] = \sum_x \sum_y g(x, y) p_{XY}(x, y)$$

# Some Extremely Important Examples

Same For Continuous Random Variables

Let  $\mu_X = E[X], \mu_Y = E[Y]$

Covariance

$$\sigma_{XY} = \text{Cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]$$

Correlation

$$\rho_{XY} = \text{Corr}(X, Y) = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

## Shortcut Formula for Covariance

Much easier for calculating:

$$\text{Cov}(X, Y) = E[XY] - E[X]E[Y]$$

I'll mention this again in a few slides...

## Calculating $\text{Cov}(X, Y)$

		Y			
		1	2	3	
X	0	1/8	0	0	1/8
	1	0	1/4	1/8	3/8
	2	0	1/4	1/8	3/8
	3	1/8	0	0	1/8
		1/4	1/2	1/4	

$$E[X] = 3/8 + 2 \times 3/8 + 3 \times 1/8 = 3/2$$

$$E[Y] = 1/4 + 2 \times 1/2 + 3 \times 1/4 = 2$$

$$\begin{aligned} E[XY] &= 1/4 \times (2 + 4) + 1/8 \times (3 + 6 + 3) \\ &= 3 \end{aligned}$$

$$\begin{aligned} \text{Cov}(X, Y) &= E[XY] - E[X]E[Y] \\ &= 3 - 3/2 \times 2 = 0 \end{aligned}$$

$$\text{Corr}(X, Y) = \text{Cov}(X, Y) / [SD(X)SD(Y)] = 0$$

Hence, zero covariance (correlation) does *not* imply independence!

# Zero Covariance versus Independence

While zero covariance (correlation) *does not* imply independence, independence *does* imply zero covariance (correlation).

You will prove this in an extension problem. . .

# Linearity of Expectation, Again

Holds for Continuous RVs as well, but different proof.

In general  $E[g(X, Y)] \neq g(E[X], E[Y])$ . But if  $g$  is linear, then:

$$E[aX + bY + c] = aE[X] + bE[Y] + c$$

where  $X, Y$  are random variables and  $a, b, c$  are constants.

There's an optional proof on the course website.

## Application: Proof of Shortcut Formula for Variance

By the Linearity of Expectation,

$$\begin{aligned} \text{Var}(X) &= E[(X - \mu)^2] = E[X^2 - 2\mu X + \mu^2] \\ &= E[X^2] - 2\mu E[X] + \mu^2 \\ &= E[X^2] - 2\mu^2 + \mu^2 \\ &= E[X^2] - \mu^2 \end{aligned}$$

## Expected Value of Sum = Sum of Expected Values

Repeatedly applying the linearity of expectation,

$$E[X_1 + X_2 + \dots + X_n] = E[X_1] + E[X_2] + \dots + E[X_n]$$

regardless of how the RVs  $X_1, \dots, X_n$  are related to each other. In particular it **doesn't matter if they're dependent or independent**.



# Independent and Identically Distributed (iid) RVs

## Example

$$X_1, X_2, \dots, X_n \sim \text{iid Bernoulli}(p)$$

## Independent

Realization of one of the RVs gives no information about the others.

## Identically Distributed

Each  $X_i$  is the same kind of RV, with the same values for any parameters. (Hence same pmf, cdf, mean, variance, etc.)

## Recall: Binomial( $n, p$ ) Random Variable

### Definition

Sum of  $n$  independent Bernoulli RVs, each with probability of “success,” i.e. 1, equal to  $p$

### Using Our New Notation

Let  $X_1, X_2, \dots, X_n \sim \text{iid Bernoulli}(p)$ ,  $Y = X_1 + X_2 + \dots + X_n$ .

Then  $Y \sim \text{Binomial}(n, p)$ .

## Expected Value of Binomial RV

Use the fact that a Binomial( $n, p$ ) RV is defined as the sum of  $n$  iid Bernoulli( $p$ ) Random Variables and the Linearity of Expectation:

$$\begin{aligned} E[Y] &= E[X_1 + X_2 + \dots + X_n] = E[X_1] + E[X_2] + \dots + E[X_n] \\ &= p + p + \dots + p \\ &= np \end{aligned}$$

## Variance of a Sum $\neq$ Sum of Variances!

$$\begin{aligned} \text{Var}(aX + bY) &= E \left[ \{(aX + bY) - E[aX + bY]\}^2 \right] \\ &\vdots \\ &= a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab \text{Cov}(X, Y) \end{aligned}$$

You'll fill in the missing steps as an extension problem. . .

Since  $\sigma_{XY} = \rho\sigma_X\sigma_Y$ , this is sometimes written as:

$$\text{Var}(aX + bY) = a^2\sigma_X^2 + b^2\sigma_Y^2 + 2ab\rho\sigma_X\sigma_Y$$

$$\text{Independence} \Rightarrow \text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$$

$X$  and  $Y$  independent  $\Rightarrow \text{Cov}(X, Y) = 0$ . Hence:

$$\begin{aligned}\text{Var}(X + Y) &= \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y) \\ &= \text{Var}(X) + \text{Var}(Y)\end{aligned}$$

Also true for three or more RVs

If  $X_1, X_2, \dots, X_n$  are independent, then

$$\text{Var}(X_1 + X_2 + \dots + X_n) = \text{Var}(X_1) + \text{Var}(X_2) + \dots + \text{Var}(X_n)$$

# Crucial Distinction

## Expected Value

Always true that

$$E[X_1 + X_2 + \dots + X_n] = E[X_1] + E[X_2] + \dots + E[X_n]$$

## Variance

Not true in general that

$$\text{Var}[X_1 + X_2 + \dots + X_n] = \text{Var}[X_1] + \text{Var}[X_2] + \dots + \text{Var}[X_n]$$

except in the special case where  $X_1, \dots, X_n$  are independent (or at least uncorrelated).

# Variance of Binomial Random Variable

## Definition from Sequence of Bernoulli Trials

If  $X_1, X_2, \dots, X_n \sim \text{iid Bernoulli}(p)$  then

$$Y = X_1 + X_2 + \dots + X_n \sim \text{Binomial}(n, p)$$

## Using Independence

$$\begin{aligned} \text{Var}[Y] &= \text{Var}[X_1 + X_2 + \dots + X_n] \\ &= \text{Var}[X_1] + \text{Var}[X_2] + \dots + \text{Var}[X_n] \\ &= p(1 - p) + p(1 - p) + \dots + p(1 - p) \\ &= np(1 - p) \end{aligned}$$

# Lecture #11 – Continuous RVs I

Introduction: Probability as Area

Probability Density Function (PDF)

Relating the PDF to the CDF

Calculating the Probability of an Interval

Calculating Expected Value for Continuous RVs



# Continuous RVs – What Changes?

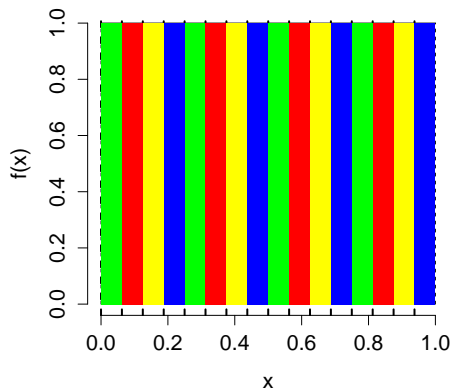
1. Probability Density Functions replace Probability Mass Functions
2. Integrals Replace Sums

Everything Else is Essentially Unchanged!

What is the probability of “Yellow?”



## From Twister to Density – Probability as *Area*



For continuous RVs, probability is defined as *area under a curve*.

Zero area means zero probability!

# Probability Density Function (PDF)

For a continuous random variable  $X$ ,

$$P(a \leq X \leq b) = \int_a^b f(x) dx$$

where  $f(x)$  is the *probability density function* for  $X$ .

## Extremely Important

For any realization  $x$ ,  $P(X = x) = 0$  since  $\int_a^a f(x)dx = 0$ . In other words, zero area means zero probability!

## For a Continuous RV, Zero Probability $\neq$ Impossible

It is *crucial* to specify the support set of a continuous RV:

- ▶ Any  $x$  outside the support set of  $X$  is *impossible*.
- ▶ Any  $x$  in the support set of  $X$  is a *possible outcome* even though  $P(X = x) = 0$  for all  $x$ .

There is no way around this slightly awkward situation: it is a consequence of defining probability as the *area under a curve*.

# Properties of PDFs

1.  $f(x) \geq 0$  for all  $x$  in the support of  $X$  and zero otherwise.
2.  $\int_{-\infty}^{\infty} f(x) dx = 1$

*Warning:  $f(x)$  is not a probability*

Can have  $f(x) > 1$  for some  $x$  as long as  $\int_{-\infty}^{\infty} f(x) dx = 1$ .

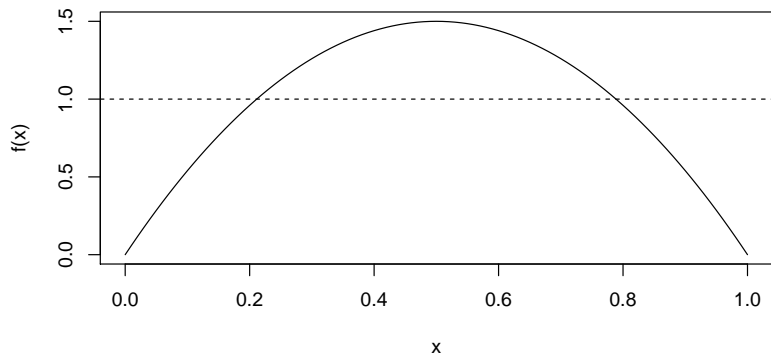
## Relating the CDF to the PDF

$$F(x_0) \equiv P(X \leq x_0) = \int_{-\infty}^{x_0} f(x) dx$$

## Example: Suppose $X$ has Support Set $[0, 1]$

Let  $f(x) = 6x(1 - x)$  for  $x \in [0, 1]$  and zero otherwise.

```
curve(6 * x * (1 - x), from = 0, to = 1, ylab = 'f(x)')  
abline(h = 1, lty = 2)
```



## Example: Suppose $X$ has Support Set $[0, 1]$

Let  $f(x) = 6x(1 - x)$  for  $x \in [0, 1]$  and zero otherwise.

Is  $f$  a valid PDF?

1. Is  $f(x) \geq 0$  for  $x \in [0, 1]$  and zero otherwise?
2. Does the total area under  $f$  equal one?

$$\begin{aligned}\int_{-\infty}^{\infty} f(x) dx &= \int_0^1 6x(1 - x) dx = 6 \int_0^1 (x - x^2) dx \\ &= 6 \left( \frac{x^2}{2} - \frac{x^3}{3} \right) \Big|_0^1 = 1\end{aligned}$$

So yes,  $f$  is a valid PDF ✓



# Integrating a Function in R

```
pdf <- function(x) {  
  6 * x * (1 - x)  
}  
  
integrate(pdf, lower = 0, upper = 1)  
  
## 1 with absolute error < 1.1e-14
```

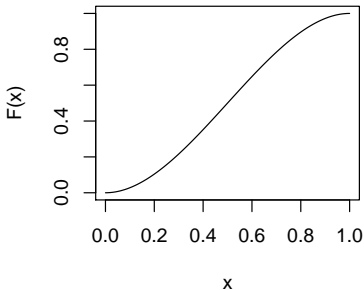
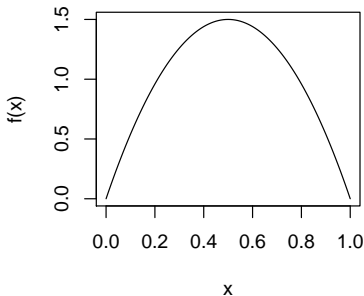
You can use this to check your work!

Example:  $f(x) = 6x(1 - x)$  for  $x \in [0, 1]$ , zero otherwise.

What is the CDF of  $X$ ?

$$\begin{aligned} F(x_0) &\equiv P(X \leq x_0) = \int_{-\infty}^{x_0} f(x) dx = \int_0^{x_0} 6x(1 - x) dx \\ &= 6 \left( \frac{x^2}{2} - \frac{x^3}{3} \right) \Big|_0^{x_0} = 3x_0^2 - 2x_0^3 \\ F(x_0) &= \begin{cases} 0, & x_0 < 0 \\ 3x_0^2 - 2x_0^3, & 0 \leq x_0 \leq 1 \\ 1, & x_0 > 1 \end{cases} \end{aligned}$$

```
par(mfrow = c(1,2))  
curve(6 * x * (1 - x), from = 0, to = 1, ylab = 'f(x)')  
curve(3 * x^2 - 2 * x^3, from = 0, to = 1, ylab = 'F(x)')
```



```
par(mfrow = c(1,1))
```

# Relationship between PDF and CDF

Integrate PDF to get CDF

$$F(x_0) = P(X \leq x_0) = \int_{-\infty}^{x_0} f(x) dx$$

Differentiate CDF to get PDF

$$f(x) = \frac{d}{dx} F(x)$$

This is just the First Fundamental Theorem of Calculus.

Example:  $f(x) = 6x(1 - x)$  for  $x \in [0, 1]$ , zero otherwise.

Differentiate CDF to get PDF

$$\begin{aligned} f(x) &= \frac{d}{dx} F(x) = \frac{d}{dx} (3x^2 - 2x^3) \\ &= 6x - 6x^2 \\ &= 6x(1 - x) \end{aligned}$$

## Key Idea: Probability of an Interval for a Continuous RV

$$P(a \leq X \leq b) = \int_a^b f(x) dx = F(b) - F(a)$$

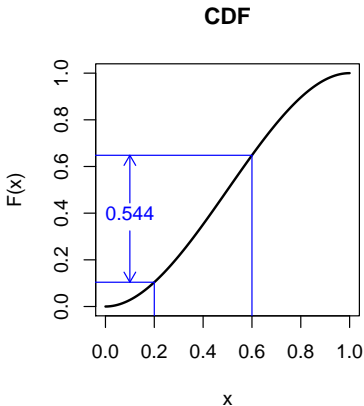
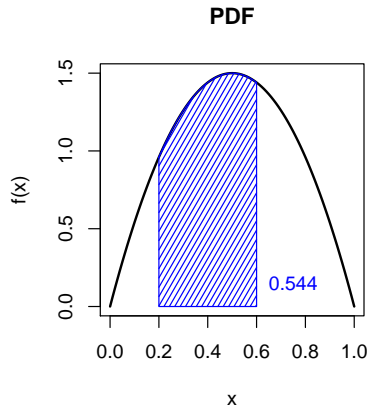
This is just the Second Fundamental Theorem of Calculus.

Example:  $f(x) = 6x(1 - x)$  for  $x \in [0, 1]$ , zero otherwise.

Two equivalent ways of calculating  $P(0.2 \leq X \leq 0.6)$

```
cdf <- function(x0) {  
  3 * x0^2 - 2 * x0^3  
}  
cdf(0.6) - cdf(0.2)  
  
## [1] 0.544  
  
integrate(pdf, lower = 0.2, upper = 0.6)  
  
## 0.544 with absolute error < 6e-15
```

Example:  $f(x) = 6x(1 - x)$  for  $x \in [0, 1]$ , zero otherwise.



$$P(0.2 \leq X \leq 0.6) = 0.544$$



## Expected Value for Continuous RVs

$$E[X] = \int_{-\infty}^{\infty} xf(x) dx$$

$$E[g(X)] = \int_{-\infty}^{\infty} g(x)f(x) dx$$

Integrals Replace Sums!

## What about all those rules for expected value?

- ▶ The only difference between expectation for continuous versus discrete is how we do the *calculation*.
- ▶ Sum for discrete; integral for continuous.
- ▶ All *properties* of expected value **continue to hold!**
- ▶ Includes linearity, shortcut for variance, etc.

## Variance of Continuous RV

$$\text{Var}(X) = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx$$

where

$$\mu = E[X] = \int_{-\infty}^{\infty} xf(x) dx$$

Shortcut formula still holds for continuous RVs!

$$\text{Var}(X) = E[X^2] - (E[X])^2$$

Example:  $f(x) = 6x(1 - x)$  for  $x \in [0, 1]$ , zero otherwise.

$$E[X] = \int_{-\infty}^{\infty} xf(x) dx = \int_0^1 x \cdot 6x(1 - x) dx = 6 \left( \frac{x^3}{3} - \frac{x^4}{4} \right) \Big|_0^1 = \frac{1}{2}$$

$$E[X^2] = \int_{-\infty}^{\infty} x^2 f(x) dx = \int_0^1 x^2 \cdot 6x(1 - x) dx = 6 \left( \frac{x^4}{4} dx - \frac{x^5}{5} \right) \Big|_0^1 = \frac{3}{10}$$

$$\text{Var}(X) = E[X^2] - (E[X])^2 = \frac{3}{10} - \left( \frac{1}{2} \right)^2 = 1/20$$

Complete the algebra at home and check using `integrate` in R.

## Simulating a Beta(2, 2) Random Variable

Our example from above is a special case of the *Beta distribution*. The command `rbeta(n, 2, 2)` makes `n` draws for this RV. These simulations agree with our calculations from above:

```
set.seed(12345)
sims <- rbeta(10000, 2, 2)
mean(sims)

## [1] 0.5007002

var(sims)

## [1] 0.05012776
```

# Simulating a Beta(2, 2) Random Variable

```
mean(sims^2)

## [1] 0.3008234

hist(sims, freq = FALSE)
```



# The Uniform Random Variable

Several of your review questions along with one of your extension questions will involve the so-called *Uniform Random Variable*:

## Uniform(0,1) Random Variable

$f(x) = 1$  for  $x \in [0, 1]$ , zero otherwise.

## Uniform(a,b) Random Variable

$f(x) = 1/(b - a)$  for  $x \in [a, b]$ , zero otherwise.

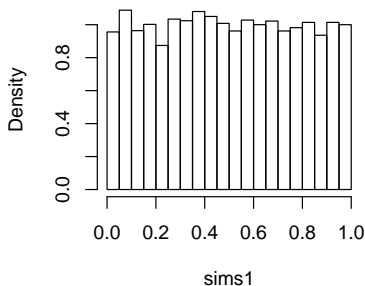
## Simulating from a Uniform RV

`runif(n, a, b)` makes `n` draws from a  $\text{Uniform}(a, b)$  RV.

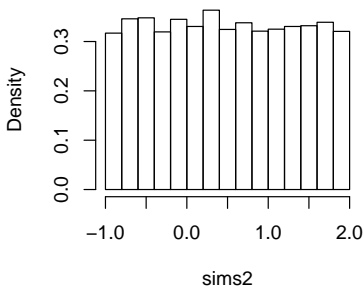
# Simulating Uniform Random Variables

```
sims1 <- runif(10000, 0, 1)
sims2 <- runif(10000, -1, 2)
par(mfrow = c(1, 2))
hist(sims1, freq = FALSE)
hist(sims2, freq = FALSE)
```

**Histogram of sims1**



**Histogram of sims2**





We don't have time to cover these in Econ 103:

### Joint Density

$$P(a \leq X \leq b \cap c \leq Y \leq d) = \int_c^d \int_a^b f(x, y) \, dx dy$$

### Marginal Densities

$$f_X(x) = \int_{-\infty}^{\infty} f(x, y) \, dy, \quad f_Y(y) = \int_{-\infty}^{\infty} f(x, y) \, dx$$

### Independence in Terms of Joint and Marginal Densities

$$f_{XY}(x, y) = f_X(x)f_Y(y)$$

### Conditional Density

$$f_{Y|X} = f_{XY}(x, y)/f_X(x)$$

# So where does that leave us?

## What We've Accomplished

We've covered all the basic properties of RVs on this [Handout](#).

## Where are we headed next?

Next up is the most important RV of all: the normal RV. After that it's time to do some statistics!

## How should you be studying?

If you *master* the material on RVs (both continuous and discrete) and in particular the normal RV the rest of the semester will seem easy. If you don't, you're in for a rough time. . .

# Lecture #12 – Continuous RVs II: The Normal RV

The Standard Normal RV

Linear Combinations and the  $N(\mu, \sigma^2)$  RV

Transforming to a Standard Normal

Percentiles/Quantiles for Continuous RVs

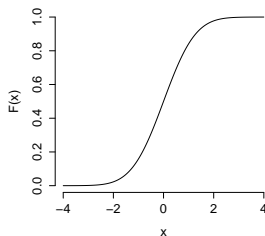
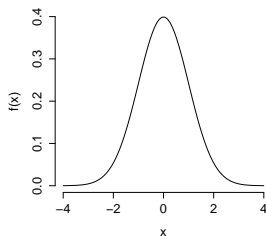
Symmetric Intervals for the  $N(0, 1)$  RV

Available on Etsy, Made using R!



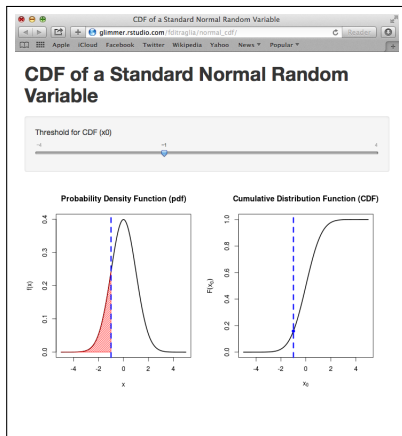
Figure: Standard Normal RV (PDF)

## Standard Normal RV: PDF at left, CDF at right



- ▶ Notation:  $X \sim N(0, 1)$
- ▶ Support Set =  $(-\infty, \infty)$
- ▶ PDF symmetric about 0, bell-shaped
- ▶  $E[X] = 0$ ,  $Var[X] = 1$
- ▶ For Econ 103, don't need formula for PDF.
- ▶ No closed-form expression for CDF.

[https://fditraglia.shinyapps.io/normal\\_cdf/](https://fditraglia.shinyapps.io/normal_cdf/)



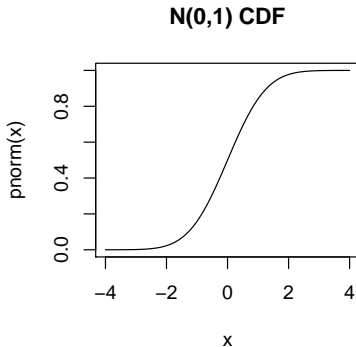
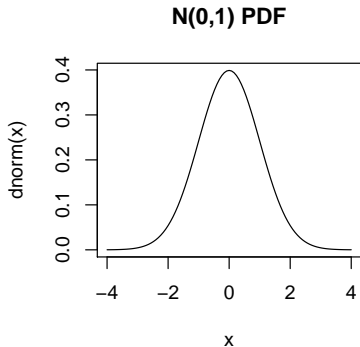
# R Commands for the Standard Normal RV

PDF $f(x)$	<code>dnorm(x)</code>
CDF $F(x)$	<code>pnorm(x)</code>
Make $n$ Random Draws	<code>rnorm(n)</code>

## Mnemonic

- ▶ `norm` = “Normal”
- ▶ `d` = “density”
- ▶ `p` = “probability”
- ▶ `r` = “random”

```
par(mfrow = c(1, 2))  
curve(dnorm(x), -4, 4, main = 'N(0,1) PDF')  
curve(pnorm(x), -4, 4, main = 'N(0,1) CDF')
```



```
par(mfrow = c(1, 1))
```



```
set.seed(1234)
normal_sims <- rnorm(10000)

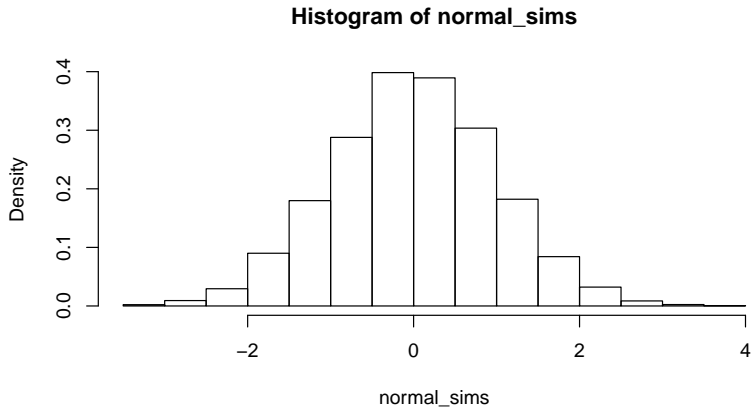
mean(normal_sims)

## [1] 0.006115893

var(normal_sims)

## [1] 0.9752143
```

```
hist(normal_sims, freq = FALSE)
```



## $Y \sim N(\mu, \sigma^2)$ Random Variable

### Linear Function of $N(0, 1)$

Let  $X \sim N(0, 1)$  and define  $Y = \mu + \sigma X$  where  $\mu, \sigma$  are constants.

### Properties of $N(\mu, \sigma^2)$

- ▶ Parameters:  $\mu, \sigma^2$ .
- ▶ Support Set =  $(-\infty, \infty)$
- ▶ PDF symmetric about  $\mu$ , bell-shaped.
- ▶ Special case:  $N(0, 1)$  has  $\mu = 0$  and  $\sigma^2 = 1$ .

What are the mean and variance of a  $N(\mu, \sigma^2)$ ? How do we know?

# Expected Value: $\mu$ shifts PDF

all of these have  $\sigma = 1$

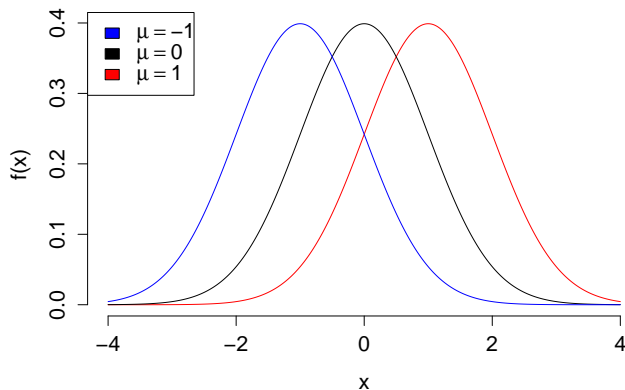


Figure: Blue  $\mu = -1$ , Black  $\mu = 0$ , Red  $\mu = 1$

# Standard Deviation: $\sigma$ scales PDF

all of these have  $\mu = 0$

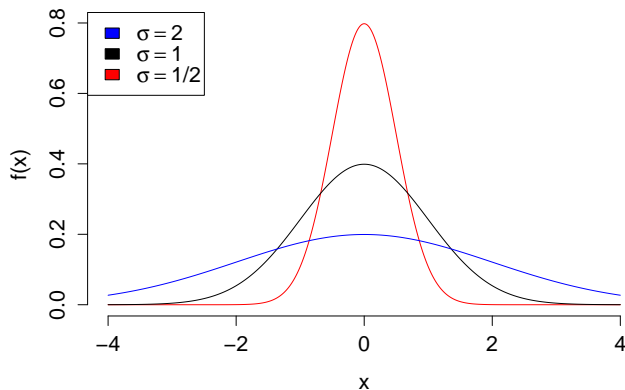


Figure: Blue  $\sigma^2 = 4$ , Black  $\sigma^2 = 1$ , Red  $\sigma^2 = 1/4$

# Linear Function of Normal RV is a Normal RV

Let  $a, b$  be constants with  $b \neq 0$

$$X \sim N(\mu, \sigma^2) \implies (a + bX) \sim N(a + b\mu, b^2\sigma^2)$$

## Key Point

Linear transformation of a normal RV is *also* a normal RV!

## Example



Suppose  $X \sim N(\mu, \sigma^2)$  and let  $Z = (X - \mu)/\sigma$ . What is the distribution of  $Z$ ?

- (a)  $N(\mu, \sigma^2)$
- (b)  $N(\mu, \sigma)$
- (c)  $N(0, \sigma^2)$
- (d)  $N(0, \sigma)$
- (e)  $N(0, 1)$

# Linear Combinations of *Multiple Independent* Normals

Let  $a, b, c$  be constants and at least one of  $a, b$  nonzero.

$X \sim N(\mu_x, \sigma_x^2)$  is independent of  $Y \sim N(\mu_y, \sigma_y^2)$  then

$$aX + bY + c \sim N(a\mu_x + b\mu_y + c, a^2\sigma_x^2 + b^2\sigma_y^2)$$

## Key Points

- ▶ Result assumes independence
- ▶ Extends to more than two Normal RVs



Suppose  $X_1, X_2, \sim \text{iid } N(\mu, \sigma^2)$



Let  $\bar{X} = (X_1 + X_2)/2$ . What is the distribution of  $\bar{X}$ ?

- (a)  $N(\mu, \sigma^2/2)$
- (b)  $N(0, 1)$
- (c)  $N(\mu, \sigma^2)$
- (d)  $N(\mu, 2\sigma^2)$
- (e)  $N(2\mu, 2\sigma^2)$

# The “Empirical Rule” Gives Probabilities for a Normal RV!

## Empirical Rule

Approximately 68% of observations within  $\mu \pm \sigma$

Approximately 95% of observations within  $\mu \pm 2\sigma$

Nearly all observations within  $\mu \pm 3\sigma$

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

$$P(\mu - \sigma \leq X \leq \mu + \sigma) \approx 0.683$$

$$P(\mu - 2\sigma \leq X \leq \mu + 2\sigma) \approx 0.954$$

$$P(\mu - 3\sigma \leq X \leq \mu + 3\sigma) \approx 0.997$$

For a continuous RV,  $P(a \leq X \leq b) = \int_a^b f(x) dx = F(b) - F(a)$

```
pnorm(1) - pnorm(-1)  # Approx. 68% Prob. in (-1,1)
```

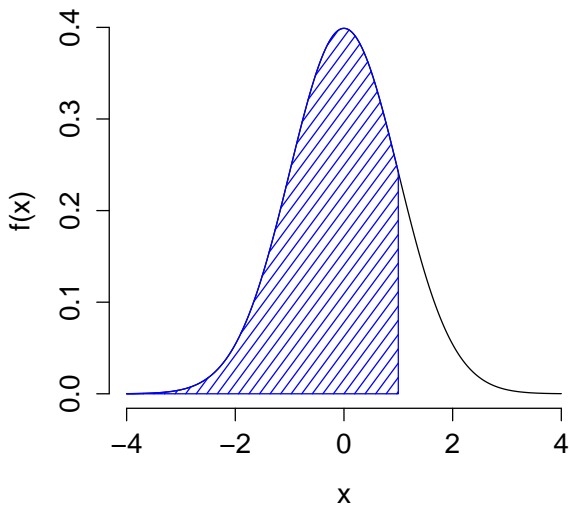
```
## [1] 0.6826895
```

```
pnorm(2) - pnorm(-2)  # Approx. 95% Prob. in (-2,2)
```

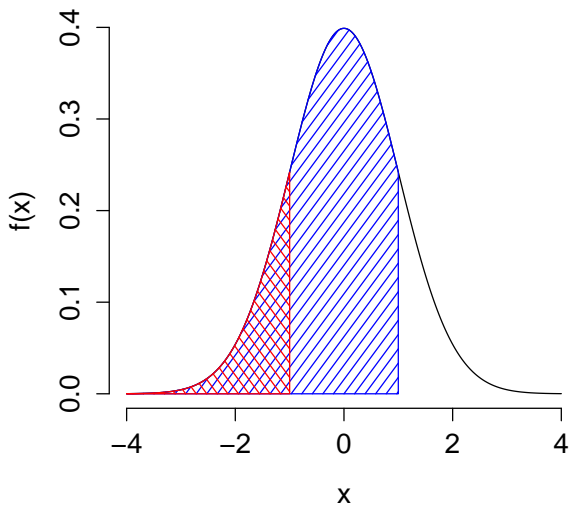
```
## [1] 0.9544997
```

```
pnorm(3) - pnorm(-3)  # > 99% Prob. in (-3,3)
```

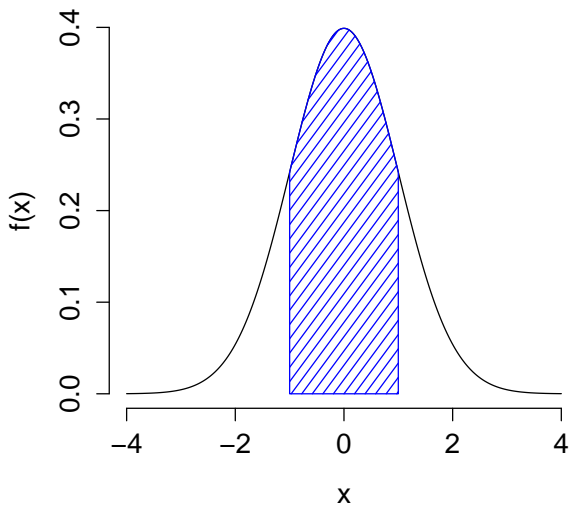
```
## [1] 0.9973002
```



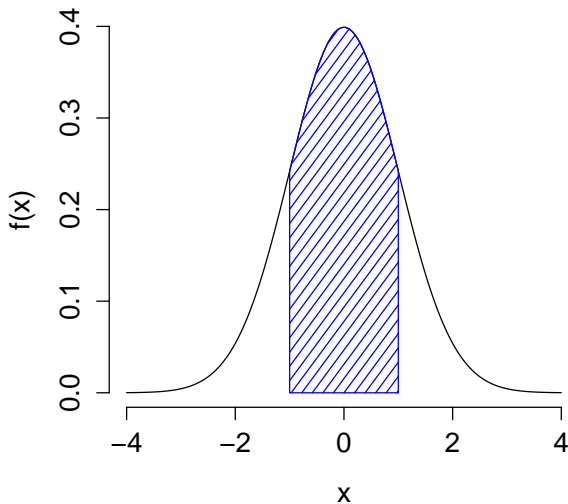
$$\text{pnorm}(1) \approx 0.84$$



$$\text{pnorm}(1) - \text{pnorm}(-1) \approx 0.84 - 0.16$$



$$\text{pnorm}(1) - \text{pnorm}(-1) \approx 0.68$$



Middle 68% of  $N(0, 1) \Rightarrow \text{approx. } (-1, 1)$

## Transforming to a Standard Normal: Example #1

Suppose  $X \sim N(\mu = 1, \sigma^2 = 4)$ . What is  $P(-1 \leq X \leq 3)$ ?

### Key Point

If  $X \sim N(\mu, \sigma^2)$  then  $\frac{X-\mu}{\sigma} \sim N(0, 1)$ .

$$\begin{aligned} P(-1 \leq X \leq 3) &= P(-2 \leq X - 1 \leq 2) \\ &= P\left(-1 \leq \frac{X - 1}{2} \leq 1\right) \\ &= \text{pnorm}(1) - \text{pnorm}(-1) \\ &\approx 0.68 \end{aligned}$$



## Transforming to a Standard Normal: Example #2

Suppose  $X \sim N(3, 16)$ . What is  $P(X \geq 10)$ ?

### Key Point

If  $X \sim N(\mu, \sigma^2)$  then  $\frac{X-\mu}{\sigma} \sim N(0, 1)$ .

$$\begin{aligned}P(X \geq 10) &= 1 - P(X \leq 10) \\&= 1 - P(X - 3 \leq 7) \\&= 1 - P\left(\frac{X - 3}{4} \leq \frac{7}{4}\right) \\&= 1 - \text{pnorm}(7/4) \approx 0.04\end{aligned}$$

# Quantile Function of a Continuous RV

Quantiles are also known as Percentiles

CDF  $F(x_0)$

- ▶  $F(x_0) \equiv P(X \leq x_0) = \int_{-\infty}^{x_0} f(x) dx$
- ▶ Input threshold  $x_0$ , get probability that  $X \leq x_0$ .

Quantile Function  $Q(p)$

- ▶  $Q(p) = F^{-1}(p)$
- ▶ Input probability  $p$ , get threshold  $x_0$  such that  $P(X \leq x_0) = p$ .
- ▶ In other words:  $p = \int_{-\infty}^{x_0} f(x) dx$

# The Median of a Continuous RV

$$\text{Median} = Q(0.5)$$

Median is the threshold  $x_0$   
such that  $P(X \leq x_0) = 0.5$ .

Median of  $N(\mu, \sigma^2)$  RV

Normal RV is symmetric  
about  $\mu$  so its median is  $\mu$ .



Figure: Median of  $N(0, 1)$  is zero.

## R Commands for the Standard Normal RV

PDF $f(x)$	<code>dnorm(x)</code>
CDF $F(x)$	<code>pnorm(x)</code>
Quantile Function $Q(p)$	<code>qnorm(p)</code>
Make $n$ Random Draws	<code>rnorm(n)</code>

### Mnemonic

- ▶ `norm` = “Normal”
- ▶ `d` = “density”
- ▶ `p` = “probability”
- ▶ `r` = “random.”
- ▶ `q` = “quantile”



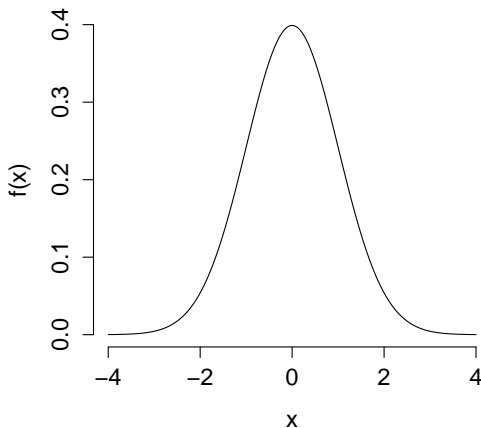
```
qnorm(0.9)  # 90th Percentile of Standard Normal
```

```
## [1] 1.281552
```

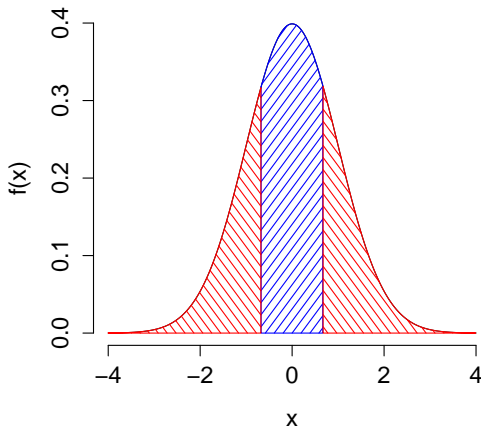
```
pnorm(1.281552)  # Check our answer using the CDF
```

```
## [1] 0.9000001
```

If  $X \sim N(0, 1)$ , for what  $c$  is  $P(-c \leq X \leq c) = 0.5$ ?



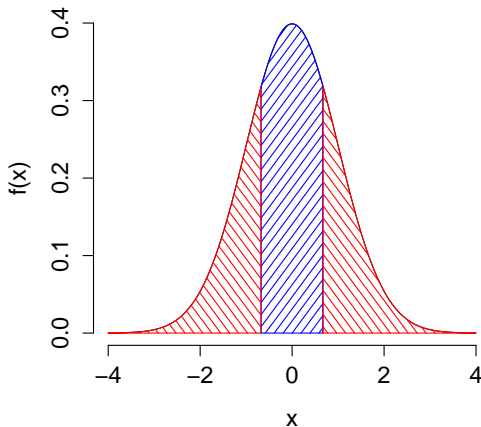
If  $X \sim N(0, 1)$ , for what  $c$  is  $P(-c \leq X \leq c) = 0.5$ ?



50% Probability in Blue; 50% Probability in Red

Boundaries of blue region are  $(-c, c)$

If  $X \sim N(0, 1)$ , for what  $c$  is  $P(-c \leq X \leq c) = 0.5$ ?

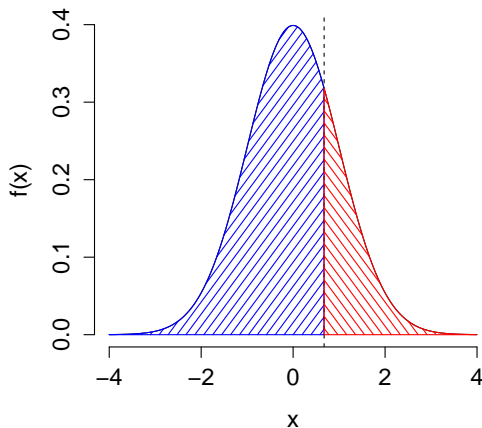


Symmetric Interval: **each red region has 25% probability**

Boundaries of blue region are  $(-c, c)$

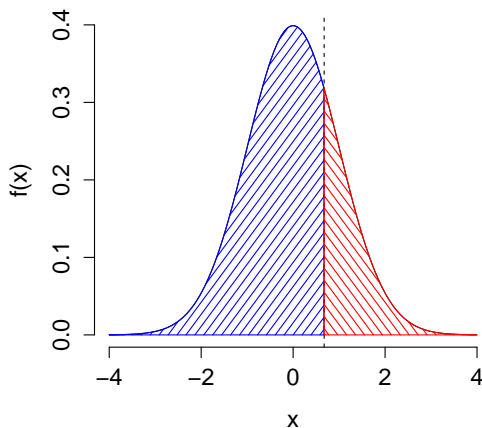


If  $X \sim N(0, 1)$ , for what  $c$  is  $P(-c \leq X \leq c) = 0.5$ ?



Let's find the right-hand boundary:  $c$

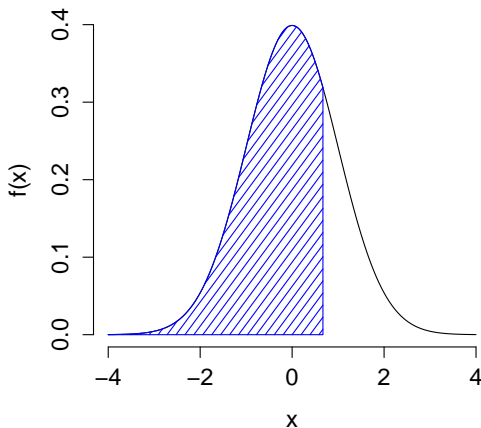
If  $X \sim N(0, 1)$ , for what  $c$  is  $P(-c \leq X \leq c) = 0.5$ ?



25% Probability to the right of  $c$

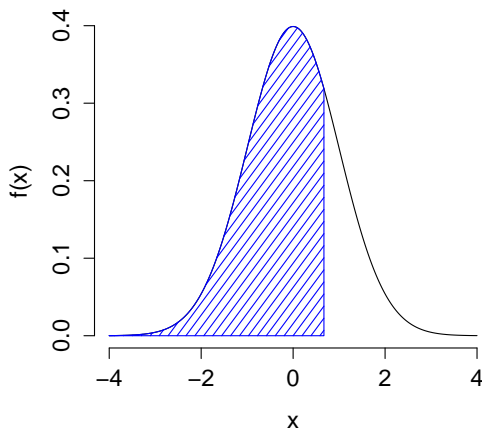
Hence, 75% to the left of  $c$

If  $X \sim N(0, 1)$ , for what  $c$  is  $P(-c \leq X \leq c) = 0.5$ ?



For what  $c$  is 75% of the probability to the left of  $c$ ?

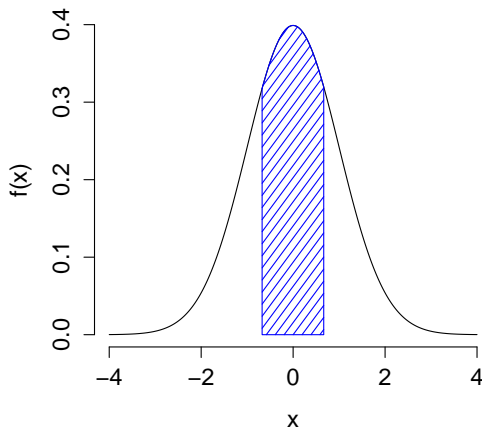
If  $X \sim N(0, 1)$ , for what  $c$  is  $P(-c \leq X \leq c) = 0.5$ ?



$$\text{qnorm}(0.75) \approx 0.67$$

Therefore  $c = 0.67$ !

If  $X \sim N(0, 1)$ , for what  $c$  is  $P(-c \leq X \leq c) = 0.5$ ?



Checking our work:

$$\text{pnorm}(0.67) - \text{pnorm}(-0.67) \approx 0.5$$



# Lecture #13 – Sampling Distributions and Estimation I

Candy Weighing Experiment

Random Sampling Redux

Unbiasedness of Sample Mean

Standard Error of the Mean

Some More Intuition for Sampling Distributions

Estimator versus Estimate

# Weighing a Random Sample

## Bag Contains 100 Candies

Estimate total weight of candies by weighing a random sample of size 5 and multiplying the result by 20.

## Your Chance to Win

The bag of candies and a digital scale will make their way around the room **during the lecture**. Each student gets a chance to draw 5 candies and weigh them.

**Student with closest estimate wins the bag of candy!**

# Weighing a Random Sample

## Procedure

When the bag and scale reach you, do the following:

1. Fold the top of the bag over and shake to randomize.
2. Randomly draw 5 candies **without replacement**.
3. Weigh your sample and record the result **in grams** along with your name on the sign-up sheet.
4. Replace your sample and shake again to re-randomize.
5. Pass bag and scale to next person.



# Sampling and Estimation

## Questions to Answer

1. How accurately do sample statistics estimate population parameters?
2. How can we quantify the uncertainty in our estimates?
3. What's so good about random sampling?

# Random Sample

## Verbal Definition from Lecture #1

Each member of population is chosen strictly by chance, so that:  
(1) selection of one individual doesn't influence selection of any other, (2) each individual is just as likely to be chosen, (3) every possible sample of size  $n$  has the same chance of selection.

## Mathematical Definition

$X_1, X_2, \dots, X_n \sim \text{iid } f(x)$  if continuous

$X_1, X_2, \dots, X_n \sim \text{iid } p(x)$  if discrete

## Random Sample Means *Sample With Replacement*

- ▶ Sampling *without replacement* creates dependence between samples (Extension Problem #11).
- ▶ But if the population is large relative to the sample, this dependence is negligible: candy experiment isn't bogus!

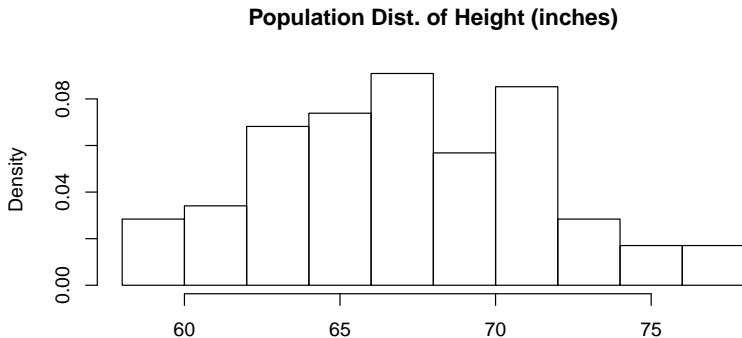
## Example: Sampling from Econ 103 Class List

- ▶ Pretend the students in this class are a population of interest.
- ▶ What is the population mean height?
- ▶ In reality I know this since I know all of your heights!
- ▶ Suppose I didn't: I could take a random sample of  $n$  students and use the sample mean to estimate the population mean.
- ▶ I know all of your heights, so I can simulate this in R.

Use this idea to explore the properties of random sampling. . .

## Example: Sampling from the Econ 103 Class List

```
survey <- read.csv('http://ditraglia.com/econ103/old_survey.csv')  
height <- na.omit(survey$height)  
hist(height, freq = FALSE, xlab = '',  
      main = 'Population Dist. of Height (inches)')
```



```
# What is the population mean?
```

```
mean(height)
```

```
## [1] 67.54545
```

```
# Draw a random sample of  $n = 5$  and compute the sample mean
```

```
set.seed(3827)
```

```
random_sample <- sample(height, 5, replace = TRUE)
```

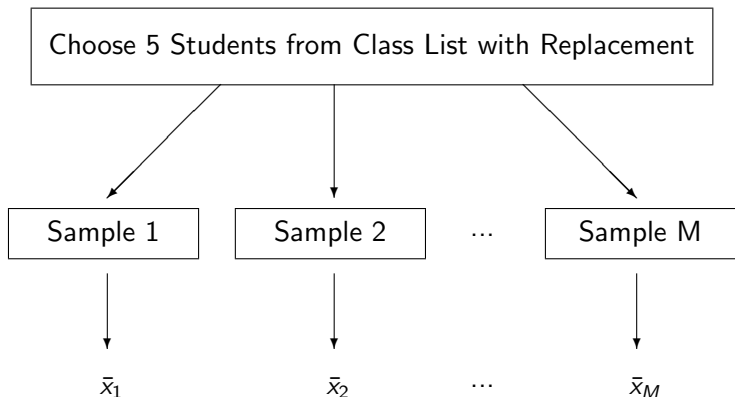
```
random_sample
```

```
## [1] 65 75 69 67 69
```

```
mean(random_sample)
```

```
## [1] 69
```

## Sampling Distribution of $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$



Repeat  $M$  times  $\rightarrow$  get  $M$  different sample means

Sampling Dist: relative frequencies of the  $\bar{x}_i$  when  $M = \infty$

```

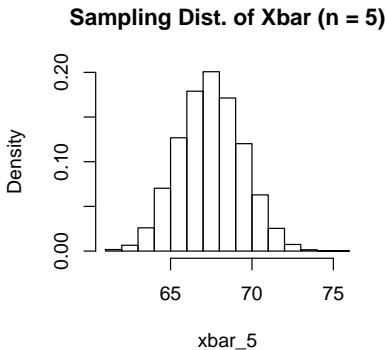
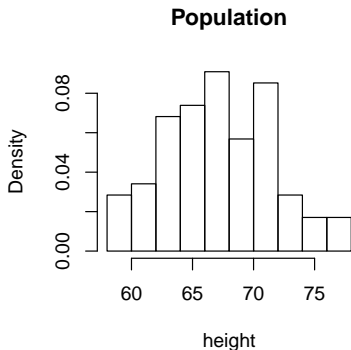
set.seed(2985)
# Function: take a random sample of size n, compute sample mean
draw_xbar <- function(n) {
  random_sample <- sample(height, size = n, replace = TRUE)
  mean(random_sample)
}
# Calculate the mean of 10000 random samples with n = 5
M <- 10000
xbar_5 <- replicate(M, draw_xbar(5))
# Compare simulated sample means to population mean: 67.5454 in.
head(xbar_5)

## [1] 65.0 64.6 69.6 68.6 64.6 65.8

```



```
# Compare popn. dist. of height to histogram of the simulated x-bars  
par(mfrow = c(1,2))  
hist(height, freq = FALSE, main = 'Population')  
hist(xbar_5, freq = FALSE, main = 'Sampling Dist. of Xbar (n = 5)')
```



```
par(mfrow = c(1,1))
```

```
# Population mean height
```

```
mean(height)
```

```
## [1] 67.54545
```

```
# Mean of sampling dist. of x-bar (n = 5)
```

```
mean(xbar_5)
```

```
## [1] 67.55678
```

```
# Population variance
```

```
var(height)
```

```
## [1] 19.74504
```

```
# Variance of sampling dist of x-bar (n = 5)
```

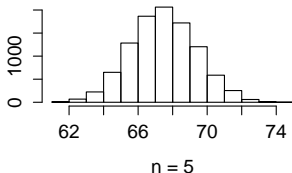
```
var(xbar_5)
```

```
## [1] 3.780202
```

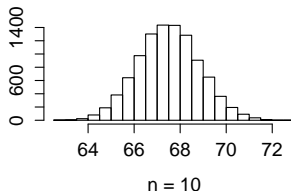
# Histograms of sampling distribution of sample mean $\bar{X}_n$

Random Sampling With Replacement, 10000 Reps. Each

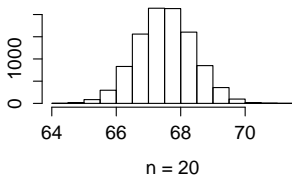
**Mean = 67.6, Var = 3.6**



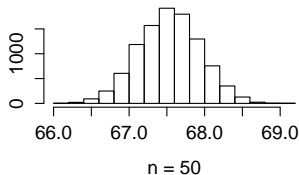
**Mean = 67.5, Var = 1.8**



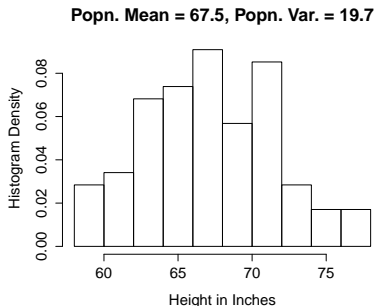
**Mean = 67.5, Var = 0.8**



**Mean = 67.5, Var = 0.2**



# Population Distribution vs. Sampling Distribution of $\bar{X}_n$



Sampling Dist. of $\bar{X}_n$		
$n$	Mean	Variance
5	67.6	3.6
10	67.5	1.8
20	67.5	0.8
50	67.5	0.2

## Things to Notice:

1. Sampling dist. “correct on average”
2. Sampling variability decreases with  $n$
3. Sampling dist. bell-shaped even though population isn't!

## Mean of Sampling Distribution of $\bar{X}_n$

$X_1, \dots, X_n \sim \text{iid}$  with mean  $\mu$

$$E[\bar{X}_n] = E\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} \sum_{i=1}^n E[X_i] = \frac{1}{n} \sum_{i=1}^n \mu = \frac{n\mu}{n} = \mu$$

Hence, sample mean is “correct on average.” The formal term for this is *unbiased*.

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

$X_1, \dots, X_n \sim \text{iid}$  with mean  $\mu$  and variance  $\sigma^2$

$$\begin{aligned}\text{Var}[\bar{X}_n] &= \text{Var}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) \\ &= \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n}\end{aligned}$$

The sampling variance of  $\bar{X}_n$  *decreases linearly with sample size.*

# Standard Error

Std. Dev. of a sampling distribution is called a **standard error**.

## Standard Error of the Sample Mean

$$SE(\bar{X}_n) = \sqrt{\text{Var}(\bar{X}_n)} = \sqrt{\sigma^2/n} = \sigma/\sqrt{n}$$

## Step 1: Population as RV rather than List of Objects

---

### Old Way

In the 2016 election, 65,853,625 out of 137,100,229 voters voted for Hillary Clinton

### New Way

Bernoulli( $p = 0.48$ ) RV

---

### Old Way

List of heights for 97 million US adult males with mean 69 in and std. dev. 6 in

### New Way

$N(\mu = 69, \sigma^2 = 36)$  RV

---

Second example assumes distribution of height is bell-shaped.



## Step 2: iid RVs Represent Random Sampling from Popn.

### Hillary Voters Example

Poll random sample of 1000 people who voted in 2016:

$$X_1, \dots, X_{1000} \sim \text{iid Bernoulli}(p = 0.48)$$

### Height Example

Measure the heights of random sample of 50 US males:

$$Y_1, \dots, Y_{50} \sim \text{iid } N(\mu = 69, \sigma^2 = 36)$$

### Key Question

What do the properties of the population imply about the properties of the sample?

## The rest of the probabilities. . .

Suppose that exactly half of US voters plan to vote for Hillary Clinton and we poll a random sample of 4 voters.

$$P(\text{Exactly 0 Hillary Voters in the Sample}) = 0.0625$$

$$P(\text{Exactly 1 Hillary Voters in the Sample}) = 0.25$$

$$P(\text{Exactly 2 Hillary Voters in the Sample}) = 0.375$$

$$P(\text{Exactly 3 Hillary Voters in the Sample}) = 0.25$$

$$P(\text{Exactly 4 Hillary Voters in the Sample}) = 0.0625$$

You should be able to work these out yourself. If not, review the lecture slides on the Binomial RV.

# Population Size is Irrelevant Under Random Sampling

## Crucial Point

*None* of the preceding calculations involved the population size: I didn't even tell you what it was! We'll never talk about population size again in this course.

## Why?

Draw with replacement  $\implies$  only the sample size and the *proportion* of Hillary supporters in the population matter.

## (Sample) Statistic

Any function of the data *alone*, e.g. sample mean  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ .

Used to estimate a population parameter: e.g.  $\bar{x}$  estimates of  $\mu$ .

## Step 3: Random Sampling $\Rightarrow$ *Sample Statistics* are RVs

This is *the crucial point of the course*: if we draw a random sample, the dataset we get is random. Since a statistic is a function of the data, it is a random variable!

# Sampling Distribution

Under random sampling, a statistic is a RV so it has a PDF if continuous or PMF if discrete: this is its **sampling distribution**.

## Sampling Dist. of Sample Mean in Polling Example

$$p(0) = 0.0625$$

$$p(0.25) = 0.25$$

$$p(0.5) = 0.375$$

$$p(0.75) = 0.25$$

$$p(1) = 0.0625$$

## Contradiction? No, but we need better terminology. . .

- ▶ Under random sampling, a statistic is a RV
- ▶ Given dataset is *fixed* so statistic is a *constant number*
- ▶ Distinguish between: **Estimator** vs. **Estimate**

### **Estimator**

Description of a general procedure.

### **Estimate**

Particular result obtained from applying the procedure.

$\bar{X}_n$  is an Estimator = Procedure = Random Variable

1. Take a random sample:  $X_1, \dots, X_n$
2. Average what you get:  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$

$\bar{x}$  is an Estimate = Result of Procedure = Constant

- ▶ Result of taking a random sample was the dataset:  $x_1, \dots, x_n$
- ▶ Result of averaging the observed data was  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

Sampling Distribution of  $\bar{X}_n$

**Thought experiment:** suppose I were to repeat the procedure of taking the mean of a random sample over and over **forever**. What **relative frequencies** would I get for the sample means?



# Lecture #14 – Sampling Distributions and Estimation II

Bias of an Estimator

Why divide by  $n - 1$  in sample variance?

Biased Sampling and the Candy-Weighing Experiment

Efficiency: Choosing between Unbiased Estimators

Mean-Squared Error: Choosing Between Biased Estimators

Consistency and the Law of Large Numbers

# Unbiased means “Right on Average”

## Bias of an Estimator

Let  $\hat{\theta}_n$  be a sample estimator of a population parameter  $\theta_0$ . The *bias* of  $\hat{\theta}_n$  is  $E[\hat{\theta}_n] - \theta_0$ .

## Unbiased Estimator

A sample estimator  $\hat{\theta}_n$  of a population parameter  $\theta_0$  is called *unbiased* if  $E[\hat{\theta}_n] = \theta_0$

## Why $(n - 1)$ for sample variance?

We will show that having  $n - 1$  in the denominator ensures:

$$E[S^2] = E \left[ \frac{1}{n - 1} \sum_{i=1}^n (X_i - \bar{X})^2 \right] = \sigma^2$$

under random sampling.

## Why $(n - 1)$ for sample variance?

Step #1 – Steps similar to Extension Problem #3 give:

$$\sum_{i=1}^n (X_i - \bar{X})^2 = \left[ \sum_{i=1}^n (X_i - \mu)^2 \right] - n(\bar{X} - \mu)^2$$

## Why $(n - 1)$ for sample variance?

Step # 2 – Take Expectations of Step # 1:

$$\begin{aligned} E \left[ \sum_{i=1}^n (X_i - \bar{X})^2 \right] &= E \left[ \left\{ \sum_{i=1}^n (X_i - \mu)^2 \right\} - n(\bar{X} - \mu)^2 \right] \\ &= E \left[ \sum_{i=1}^n (X_i - \mu)^2 \right] - E [n(\bar{X} - \mu)^2] \\ &= \sum_{i=1}^n E [(X_i - \mu)^2] - n E [(\bar{X} - \mu)^2] \end{aligned}$$

Where we have used the linearity of expectation.

## Why $(n - 1)$ for sample variance?

Step # 3 – Use assumption of random sampling:

$X_1, \dots, X_n \sim$  iid with mean  $\mu$  and variance  $\sigma^2$

$$\begin{aligned} E \left[ \sum_{i=1}^n (X_i - \bar{X})^2 \right] &= \sum_{i=1}^n E \left[ (X_i - \mu)^2 \right] - n E \left[ (\bar{X} - \mu)^2 \right] \\ &= \sum_{i=1}^n \text{Var}(X_i) - n E \left[ (\bar{X} - E[\bar{X}])^2 \right] \\ &= \sum_{i=1}^n \text{Var}(X_i) - n \text{Var}(\bar{X}) = n\sigma^2 - \sigma^2 \\ &= (n - 1)\sigma^2 \end{aligned}$$

Since  $E[\bar{X}] = \mu$  and  $\text{Var}(\bar{X}) = \sigma^2/n$  under random sampling.

## Why $(n - 1)$ for sample variance?

Finally – Divide Step # 3 by  $(n - 1)$ :

$$E[S^2] = E\left[\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2\right] = \frac{(n-1)\sigma^2}{n-1} = \sigma^2$$

Hence, having  $(n - 1)$  in the denominator ensures that the sample variance is “correct on average,” that is *unbiased*.

## A Different Estimator of the Population Variance

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

$$E[\hat{\sigma}^2] = E \left[ \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \right] = \frac{1}{n} E \left[ \sum_{i=1}^n (X_i - \bar{X})^2 \right] = \frac{(n-1)\sigma^2}{n}$$

Bias of  $\hat{\sigma}^2$

$$E[\hat{\sigma}^2] - \sigma^2 = \frac{(n-1)\sigma^2}{n} - \sigma^2 = \frac{(n-1)\sigma^2}{n} - \frac{n\sigma^2}{n} = -\sigma^2/n$$



# How Large is the Average Family?



How many brothers and sisters are in your family, including yourself?

# What's Going On Here?

Twenty years ago the average number of children per family was about 2.0. But our average was much higher!

## Biased Sample!

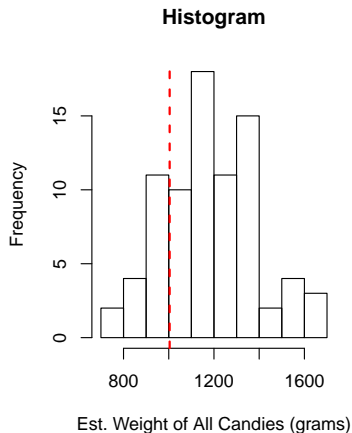
- ▶ Zero children  $\Rightarrow$  didn't send any to college
- ▶ Sampling by *children* so large families **oversampled**

## Candy Weighing: 80 Estimates, Each With $n = 5$

$$\hat{\theta} = 20 \times (X_1 + \dots + X_5)$$

Summary of Sampling Dist.	
Overestimates	63
Exactly Correct	0
Underestimates	17
$E[\hat{\theta}]$	1194 grams
$SD(\hat{\theta})$	206 grams

Actual Mass:  $\theta_0 = 1004$  grams



# What was in the bag?

100 Candies Total:

- ▶ 20 Fun Size Snickers Bars (large)
- ▶ 30 Reese's Miniatures (medium)
- ▶ 50 Tootsie Roll "Midgees" (small)

## So What Happened?

Not a random sample! The Snickers bars were *oversampled*.

Could we have avoided this? How?



Let  $X_1, X_2, \dots, X_n \sim iid$  mean  $\mu$ , variance  $\sigma^2$ . True or False:

*$X_1$  is an unbiased estimator of  $\mu$*

(a) True

(b) False

TRUE!

## How to choose between two unbiased estimators?

Suppose  $X_1, X_2, \dots, X_n \sim iid$  with mean  $\mu$  and variance  $\sigma^2$

From Last Lecture:

$$E[\bar{X}_n] = \mu, \quad \text{Var}(\bar{X}_n) = \sigma^2/n$$

Compared To:

$$E[X_1] = \mu, \quad \text{Var}(X_1) = \sigma^2$$

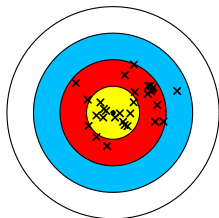
Both  $\bar{X}_n$  and  $X_1$  are unbiased estimators of  $\mu$ , but  $\bar{X}_n$  has a lower variance!

## Efficiency - Compare Unbiased Estimators by Variance

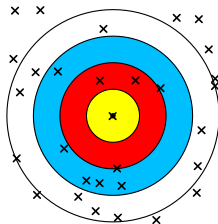
Let  $\hat{\theta}_1$  and  $\hat{\theta}_2$  be unbiased estimators of  $\theta_0$ . We say that  $\hat{\theta}_1$  is *more efficient* than  $\hat{\theta}_2$  if  $\text{Var}(\hat{\theta}_1) < \text{Var}(\hat{\theta}_2)$ .

# Bias and Variance are Both Bad Things

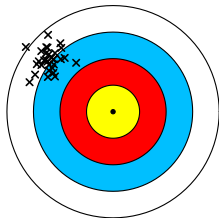
**Low Bias, Low Variance**



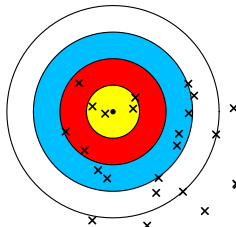
**Low Bias, High Variance**



**High Bias, Low Variance**



**High Bias, High Variance**





## Mean-Squared Error: Trading Bias Against Variance

- ▶ Unbiased estimator with a huge variance is bad.
- ▶ Highly biased estimator with a low variance is bad.
- ▶ Often there is a “tradeoff” between bias and variance:
  - ▶ Low bias estimators often have high variance.
  - ▶ Low variance estimators often have high bias.

### Mean-Squared Error (MSE):

Compare estimators accounting for **both** bias and variance:

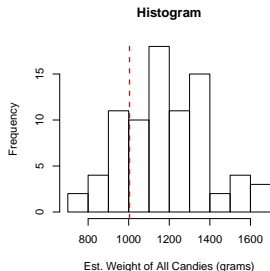
$$MSE(\hat{\theta}) = \text{Bias}(\hat{\theta})^2 + \text{Var}(\hat{\theta})$$

**Root Mean-Squared Error (RMSE):**  $\sqrt{\text{MSE}}$

# Calculate MSE for Candy Experiment



$E[\hat{\theta}]$	1194 grams
$\theta_0$	1004 grams
$SD(\hat{\theta})$	206 grams



$$\begin{aligned}\text{Bias} &= 1194 \text{ grams} - 1004 \text{ grams} \\ &= 190 \text{ grams}\end{aligned}$$

$$\begin{aligned}\text{MSE} &= \text{Bias}^2 + \text{Variance} \\ &= (190^2 + 206^2) \text{ grams}^2 \\ &= 7.8536 \times 10^4 \text{ grams}^2\end{aligned}$$

$$\text{RMSE} = \sqrt{\text{MSE}} = 280 \text{ grams}$$

# Finite Sample versus Asymptotic Properties of Estimators

## Finite Sample Properties

For *fixed sample size*  $n$  what are the properties of the sampling distribution of  $\hat{\theta}_n$ ? (E.g. bias and variance.)

## Asymptotic Properties

What happens to the sampling distribution of  $\hat{\theta}_n$  *as the sample size  $n$  gets larger and larger?*

1. Law of Large Numbers (today)
2. Central Limit Theorem (Lecture 16)

# Consistency

## Definition

We say that an estimator  $\hat{\theta}_n$  is *consistent for* a parameter  $\theta_0$  if  $\lim_{n \rightarrow \infty} \text{MSE}(\hat{\theta}_n) = 0$ , in other words, if both the bias and variance of  $\hat{\theta}_n$  disappear as the sample size grows.

Intuitively, this means  $\hat{\theta}_n$  becomes “less random” as the sample size increases, eventually converging to a constant:  $\theta_0$ .

# Law of Large Numbers

Let  $X_1, X_2, \dots, X_n \sim iid$  mean  $\mu$ , variance  $\sigma^2$ . Then the sample mean  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  is consistent for the population mean  $\mu$ .

How do we know this?

From our last lecture:

$$E[\bar{X}_n] = \mu, \quad \text{Var}(\bar{X}_n) = \sigma^2/n$$

and hence:

$$\begin{aligned} \text{MSE}(\bar{X}_n) &= \text{Bias}(\bar{X}_n)^2 + \text{Var}(\bar{X}_n) \\ &= (E[\bar{X}_n] - \mu)^2 + \text{Var}(\bar{X}_n) \\ &= 0 + \sigma^2/n \rightarrow 0 \end{aligned}$$

```
set.seed(12345)
n <- 10000
x <- rnorm(n, mean = 0, sd = 10)
xbar_n <- cumsum(x) / (1:n)
plot(xbar_n, type = 'l', xlab = 'n', ylab = 'Sample Mean')
```

