

This class is being conducted over Zoom. As the instructor, I will be **recording** this session. I have disabled the recording feature for others so that no one else will be able to record this session. I will be posting this session to the course's website.

If you have privacy concerns and **do not wish to appear in the recording**, you may turn video off (click "**stop video**") so that Zoom does not record you.

The chat box is always open for discussion and questions to the entire class. You may also send messages privately to the instructor or the TAs. Please note that Zoom saves all chat transcripts.

Introduction to Statistics

Stats 7

Mary Ryan

Aug. 4, 2020



<https://canvas.eee.uci.edu/courses/28451>

About the Teaching Team

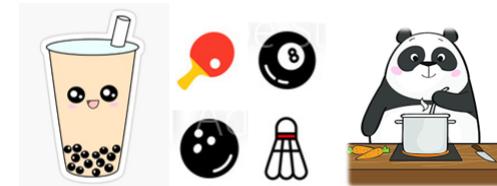
- Instructor: Mary Ryan
 - 5th year Statistics PhD student
 -  : marymr@uci.edu
 - OH:



- TA: Kyle Conniff
 - 5th year Statistics PhD student
 -  : krconnif@uci.edu
 - OH:



- TA: Jenifer Rim
 - 4th year Statistics PhD student
 -  : jsrim@uci.edu
 - OH:



Tentative Schedule

Week 1

- Descriptive statistics & data visualization
- Data sampling & bias

Week 2

- Regression
- Basic probability

Week 3

- **Midterm**
- Discrete probability distributions

Week 4

- Continuous probability distributions & sampling distributions
- One-sample inference

Week 5

- Two-sample inference
- Inference for categorical data

Final Exam

Textbook

Recommended text: **OpenIntro Statistics, 4th Edition**

- PDF free to download

The screenshot shows the Leanpub bookstore page for 'OpenIntro Statistics'. At the top, there's a navigation bar with links for Store, Books, Bundles, Courses, Featured, Newsletters, Podcast, Support, and Why Leanpub. Below the navigation, it says '81,495 READERS' and '422 PAGES'. The main content area features the book cover for 'OpenIntro Statistics, Fourth Edition'. The cover is dark blue with the title 'OpenIntro Statistics' and 'Fourth Edition' in white. A large number '4' is prominently displayed at the bottom right. Below the cover, it says 'This book is 100% complete'. To the left of the cover, there's a sidebar with author information: 'Includes 1st, 2nd, 3rd, and 4th Editions', 'David Diez, Mine Çetinkaya-Rundel, and Christopher Barr', and a 'Complete foundation for Statistics, also serving as a foundation for Data Science.' section. On the right side of the book cover, there's a price comparison table showing 'Free!', 'MINIMUM PRICE \$15.00', and 'SUGGESTED PRICE \$15.00'. It also shows 'YOU PAY \$0.00' and 'AUTHORS EARN \$0.00'. Below this, there are options for 'The Book' or 'The Book + Tablet-Friendly PDF of the Book'. A 'Read Free Sample' button is at the bottom left, and an 'Add Ebook to Cart' button is at the bottom right.

- Paperback copy \$20 from Amazon

The screenshot shows the Amazon product page for 'OpenIntro Statistics: Fourth Edition'. At the top, it says 'Books > Science & Math > Mathematics'. The main title is 'OpenIntro Statistics: Fourth Edition' by David Diez, Mine Çetinkaya-Rundel, and Christopher Barr. It has a rating of 4.5 stars from 16 ratings. A 'Look inside' button is available. The price is listed as '\$18.30' for Paperback, with a note that it's from \$18.30. There are tabs for 'Paperback' and 'Other Sellers'. Below the price, it says 'Buy new' and 'In Stock'. It also mentions 'Arrives: Jan 31 - Feb 5', 'Fastest delivery: Mon, Jan 27', and 'Order within 3 hrs 58 mins'. An 'Amazon Hub Locker+ (Irvine) - Irvine 92612' option is shown. On the right, there's a sidebar with shipping information: 'List Price: \$20.00', 'Save: \$1.70 (9%)', '1 New from \$18.30', and '& FREE Shipping on orders over \$25.00 shipped by Amazon. Details'. There are buttons for 'Add to Cart' and 'Buy Now'. At the bottom, it says '4 used & new from \$18.30' and 'See All Buying Options'.

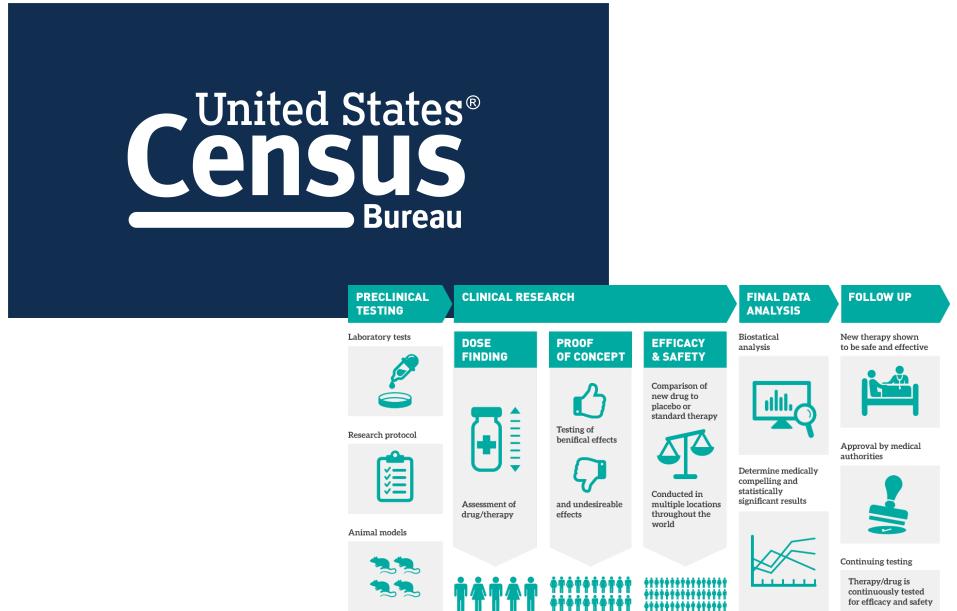
Assignments

- Video Quizzes & Surveys: 15%
 - Short topic videos with embedded quiz questions to be completed before appropriate lecture, found on [course website](#)
- Homework: 25%
 - Weekly assignments, found on the [course website](#)
 - Open Tuesdays @ 1p
 - Due [following Tuesday @ 12.59p](#) (before lecture) via Canvas
 - For [every day](#) late, homework grade will suffer a [20%](#) penalty
- Labs: 15%
 - Open Thursdays @ 4p, found on the [course website](#)
 - Due [following Tuesday @ 12.59p](#) (before lecture) via Canvas
- Midterm Exam: 25%
 - Date: [Tuesday, Aug. 18 @ 1p](#)
- Final Exam: 30%
 - Date: [Tuesday, Sept. 8 @ 1p](#)
 - Cumulative

If you have issues with either of the exam dates, let me know as soon as possible

What is Statistics?

What is Statistics?



FiveThirtyEight



BUREAU OF LABOR STATISTICS
U.S. DEPARTMENT OF LABOR



The Grammar of Science

Grammar takes a pile of words,



and helps turn it into a sentence that makes sense

*In this class, we are going to learn about
data, statistics, and how they impact our lives.*

The Grammar of Science

Grammar takes a pile of words,



and helps turn it into a sentence that makes sense

*In this class, we are going to learn about
data, statistics, and how they impact our lives.*

Statistics takes a pile of data,



and helps turn it into scientific conclusions that we can interpret

What Statistics Isn't

Just like all sentences that are grammatically correct aren't necessarily good at communicating their meaning...

*Buffalo buffalo Buffalo buffalo
buffalo buffalo Buffalo buffalo*

Statistics is not simply a machine where data goes in and truth pops out the other side

We need to use Statistics **consciously** and understand its **limitations**

Types of Data

- A **variable** is a measured characteristic of your data
 - A dataset is often made up of many different variables, of many different types

Types of Data

- A **variable** is a measured characteristic of your data
 - A dataset is often made up of many different variables, of many different types
- **Categorical** (qualitative) variables
 - Ordinal variable: categorical variable whose categories have a **specific ordering**
 - Nominal variable: categorical variable **without any ordering**
 - Binary variable: categorical variable with **only 2 categories**

Types of Data

- A **variable** is a measured characteristic of your data
 - A dataset is often made up of many different variables, of many different types
- **Categorical** (qualitative) variables
 - Ordinal variable: categorical variable whose categories have a **specific ordering**
 - Nominal variable: categorical variable **without any ordering**
 - Binary variable: categorical variable with **only 2 categories**
- **Quantitative** variables
 - Discrete variable: takes on numerical values in **jumps**
 - Continuous variable: takes on numerical values that can go out to **infinite decimal points**

Types of Data

- 1) Number of cats at the animal shelter
- 2) Species of pet
- 3) Weight
- 4) Mood
- 5) Diagnosis
- 6) Population of Orange County
- 7) Airplane ticket class

Types of Data

- 1) Number of cats at the animal shelter Discrete
- 2) Species of pet
- 3) Weight
- 4) Mood
- 5) Diagnosis
- 6) Population of Orange County
- 7) Airplane ticket class

Types of Data

- | | |
|---|----------|
| 1) Number of cats at the animal shelter | Discrete |
| 2) Species of pet | Nominal |
| 3) Weight | |
| 4) Mood | |
| 5) Diagnosis | |
| 6) Population of Orange County | |
| 7) Airplane ticket class | |

Types of Data

- | | |
|---|------------|
| 1) Number of cats at the animal shelter | Discrete |
| 2) Species of pet | Nominal |
| 3) Weight | Continuous |
| 4) Mood | |
| 5) Diagnosis | |
| 6) Population of Orange County | |
| 7) Airplane ticket class | |

Types of Data

- | | |
|---|------------|
| 1) Number of cats at the animal shelter | Discrete |
| 2) Species of pet | Nominal |
| 3) Weight | Continuous |
| 4) Mood | Nominal |
| 5) Diagnosis | |
| 6) Population of Orange County | |
| 7) Airplane ticket class | |

Types of Data

1) Number of cats at the animal shelter	Discrete
2) Species of pet	Nominal
3) Weight	Continuous
4) Mood	Nominal
5) Diagnosis	Nominal
6) Population of Orange County	
7) Airplane ticket class	

Types of Data

1) Number of cats at the animal shelter	Discrete
2) Species of pet	Nominal
3) Weight	Continuous
4) Mood	Nominal
5) Diagnosis	Nominal
6) Population of Orange County	Discrete
7) Airplane ticket class	

Types of Data

1) Number of cats at the animal shelter	Discrete
2) Species of pet	Nominal
3) Weight	Continuous
4) Mood	Nominal
5) Diagnosis	Nominal
6) Population of Orange County	Discrete
7) Airplane ticket class	Ordinal

Types of Data

1) Number of cats at the animal shelter	Discrete
2) Species of pet	Nominal
3) Weight	Continuous
4) Mood	Nominal
5) Diagnosis	Nominal
6) Population of Orange County	Discrete
7) Airplane ticket class	Ordinal
8) Money you spend to wash 1 load of laundry	

Types of Data

1) Number of cats at the animal shelter	Discrete
2) Species of pet	Nominal
3) Weight	Continuous
4) Mood	Nominal
5) Diagnosis	Nominal
6) Population of Orange County	Discrete
7) Airplane ticket class	Ordinal
8) Money you spend to wash 1 load of laundry	Continuous

Types of Data

8) Money you spend to wash 1 load of laundry

- I spend \$1.75 per wash cycle

Types of Data

8) Money you spend to wash 1 load of laundry

- I spend \$1.75 per wash cycle
 - Looks discrete, right?

Types of Data

8) Money you spend to wash 1 load of laundry

- I spend \$1.75 per wash cycle
 - Looks discrete, right?
 - **Converted to Iraqi Dinar: 2080.75 د.إ**

Types of Data

8) Money you spend to wash 1 load of laundry

- I spend \$1.75 per wash cycle
 - Looks discrete, right?
 - **Converted to Iraqi Dinar: 2080.75 د.إ**
 - **Converted to Euro: € 1.5596**

Types of Data

8) Money you spend to wash 1 load of laundry

- I spend \$1.75 per wash cycle
 - Looks discrete, right?
 - **Converted to Iraqi Dinar: 2080.75 د.إ**
 - **Converted to Euro: € 1.5596**
 - **Converted to Mexican Peso: \$ 40.285735**

Types of Data

8) Money you spend to wash 1 load of laundry

- I spend \$1.75 per wash cycle
 - Looks discrete, right?
 - Converted to Iraqi Dinar: 2080.75 د.إ
 - Converted to Euro: € 1.5596
 - Converted to Mexican Peso: \$ 40.285735
- While we might think of currency in "discrete" US Dollar (or Dinar, or Euro, or Peso) units, any amount of money can be converted from one currency to another
 - Currency conversion doesn't respect the fact that you don't carry around $1/12^{th}$ cents
 - So we generally think of currency as a continuous variable

How Do We Describe Our Data?

For categorical data, we might want to:

- Know how many observations are in each category
 - **Counts**
- Know how large a category is compared to all observations
 - **Percentages**
 - $\frac{\text{\# observations in a category}}{\text{Total observations}} \times 100$

How Do We Describe Our Data?

For quantitative data, we might want to:

- Know what the smallest value of a variable is
 - **Minimum**
- Know what the largest value of a variable is
 - **Maximum**
- Know the distance between the maximum and the minimum
 - **Range**
- Know what the center of our data is...

Journey to the Center of the Data

- **Mean**

- The "typical value" of the data

- $\bar{x} = \frac{\text{sum of data points}}{\text{number of data points}} = \frac{\sum_{i=1}^n X_i}{n}$

- **Median**

- The "exact middle" of the data
 - 50% of data points are below the median, 50% of data points are above the median

- **Mode**

- The most common value in the data

Journey to the Center of the Data

- **Mean**

- The "typical value" of the data

- $\bar{x} = \frac{\text{sum of data points}}{\text{number of data points}} = \frac{\sum_{i=1}^n X_i}{n}$

- **Median**

- The "exact middle" of the data
 - 50% of data points are below the median, 50% of data points are above the median

- **Mode**

- The most common value in the data

Median

2, 5, 6, 8, 8, 8, 10, 13, 14, 16, 16, 19, 20, 21, 25

Median

~~2, 5, 6, 8, 8, 8, 10, 13, 14, 16, 16, 19, 20, 21, 25~~

Median

~~2, 5, 6, 8, 8, 8, 10, 13, 14, 16, 16, 19, 20, 21, 25~~

Median

~~2, 5, 6, 8, 8, 8, 10, 13, 14, 16, 16, 19, 20, 21, 25~~

Median

~~2, 5, 6, 8, 8, 8, 10,~~ 13, ~~14, 16, 16, 19, 20, 21, 25~~

Median

~~2, 5, 6, 8, 8, 8, 10,~~ 13, ~~14, 16, 16, 19, 20, 21, 25~~

2, 5, 6, 8, 8, 8, 10, 10, 13, 14, 16, 16, 19, 20, 21, 25

Median

~~2, 5, 6, 8, 8, 8, 10,~~ 13, ~~14, 16, 16, 19, 20, 21, 25~~

~~2, 5, 6, 8, 8, 8, 10, 10, 13, 14, 16, 16, 19, 20, 21, 25~~

Median

~~2, 5, 6, 8, 8, 8, 10,~~ 13, ~~14, 16, 16, 19, 20, 21, 25~~

~~2, 5, 6, 8, 8, 8, 10,~~ 10, 13, ~~14, 16, 16, 19, 20, 21, 25~~

Median

~~2, 5, 6, 8, 8, 8, 10,~~ 13, ~~14, 16, 16, 19, 20, 21, 25~~

~~2, 5, 6, 8, 8, 8, 10,~~ 10, 13, ~~14, 16, 16, 19, 20, 21, 25~~

$$(10+13)/2 = 11.5$$

How Do We Describe Our Data?

- We might also want to know how much our observations **vary**
- Nice to compare how far away observations are from the center of the data (**deviation**)

$$x_i - \bar{x}$$

How Do We Describe Our Data?

- We might also want to know how much our observations **vary**
- Nice to compare how far away observations are from the center of the data (**deviation**)

$$x_i - \bar{x}$$

- But some observations are above the mean and some are below...
 - If we just add the deviations up some of them might cancel out
 - Square the deviations to remove the sign:

$$(x_i - \bar{x})^2$$

How Do We Describe Our Data?

- We might also want to know how much our observations **vary**
- Nice to compare how far away observations are from the center of the data (**deviation**)

$$x_i - \bar{x}$$

- But some observations are above the mean and some are below...
 - If we just add the deviations up some of them might cancel out
 - Square the deviations to remove the sign:

$$(x_i - \bar{x})^2$$

- But every observation has a different squared deviance. How do we get one metric?

How Do We Describe Our Data?

- We might also want to know how much our observations **vary**
- Nice to compare how far away observations are from the center of the data (**deviation**)

$$x_i - \bar{x}$$

- But some observations are above the mean and some are below...
 - If we just add the deviations up some of them might cancel out
 - Square the deviations to remove the sign:

$$(x_i - \bar{x})^2$$

- But every observation has a different squared deviance. How do we get one metric?
 - Take an average!
 - **Variance** (s^2) = $\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$
 - **Standard deviation** (s) = $\sqrt{s^2} = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}}$

Variance & Standard Deviation

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

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

- When variance/standard deviation is **large**, the observations vary from the mean a lot
- When variance/standard deviation is **small**, the observations generally stay close to the mean

Data Visualization

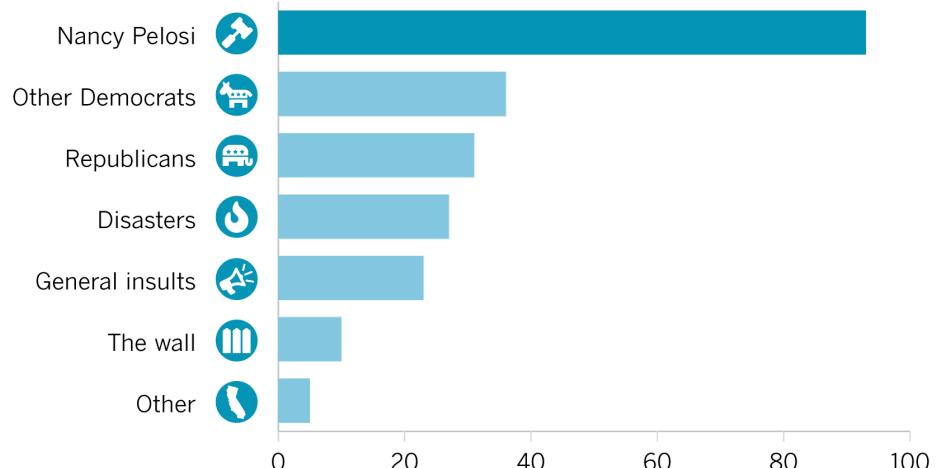
- Summarizing data with metrics are nice, but people are much better at identifying differences/patterns **graphically**

Data Visualization

- Summarizing data with metrics are nice, but people are much better at identifying differences/patterns **graphically**

- **Bar plot**

- Categories on one axis
- Counts/percentages on the other axis
- AKA: bar graph, bar chart, barplot

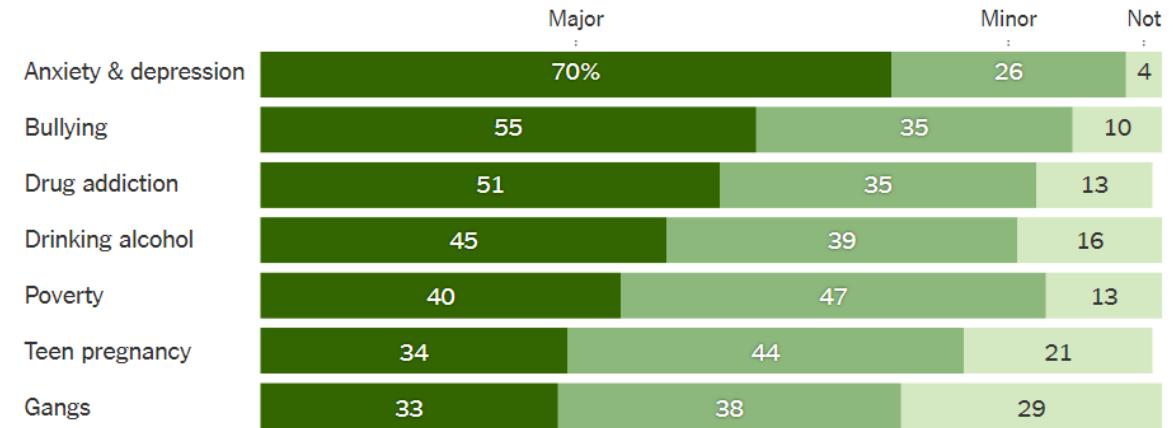


The California-related topics the president tweets most frequently about (Sept. 17, 2019), via Priya Krishnakumar
at The LA Times

Data Visualization

- Stacked bar plot

- Categories on one axis
- Percentage on other axis
- Bar divided into percentages of how common 2nd categorical variable is within your axis category
- Full bar should add up to 100%

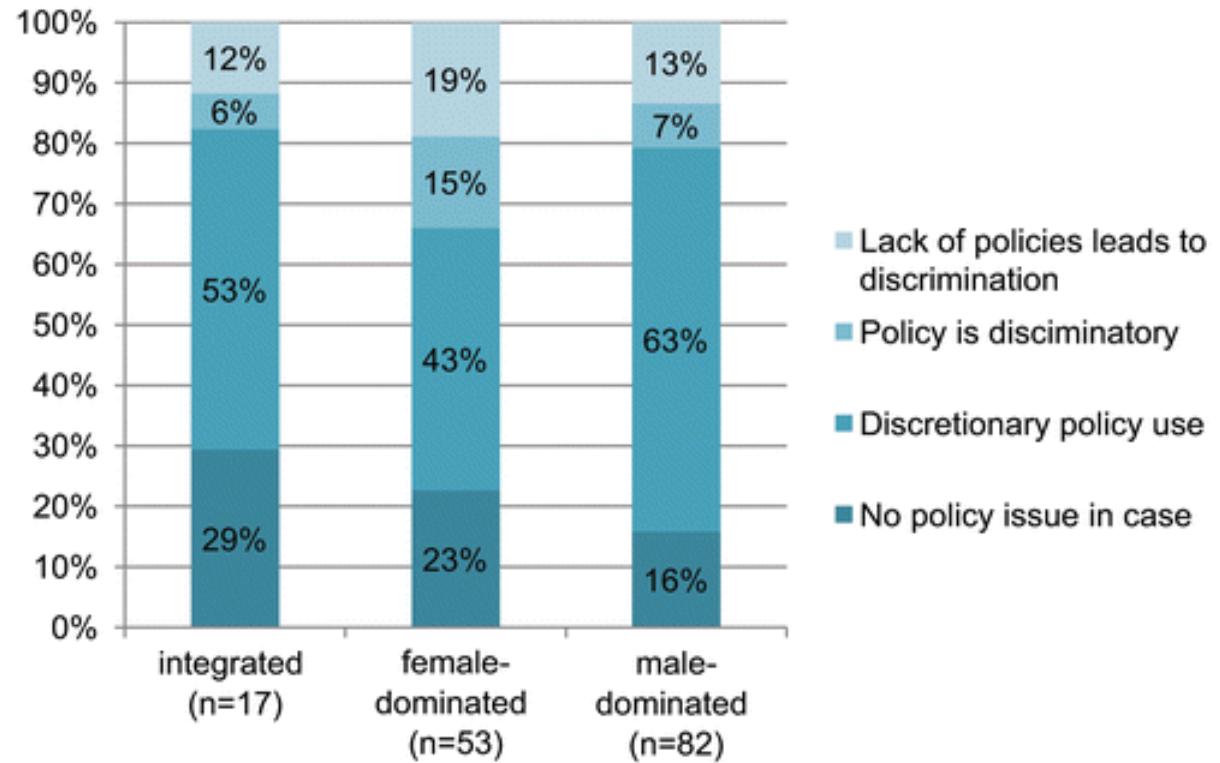


Teenagers Say Depression and Anxiety Are Major Issues Among Their Peers (Feb. 29, 2019),

via Karen Zraick at [The New York Times](#)

Data Visualization

- Stacked bar chart
 - What's going on here?



Gender Discrimination at Work: Connecting Gender Stereotypes, Institutional Policies, and Gender

Composition of Workplace (Dec. 5, 2011), via Donna Bobbitt-Zehner at [Gender & Society](#)

Data Visualization

- Pie chart
 - Can also visualize percentages
 - No real axes
 - Total circle = 100%
 - Each wedge corresponds to percentage of total observations correspond to one category
 - People are much worse at judging volume by circle wedge than volume by box, though

Remove
to improve
the pie chart edition

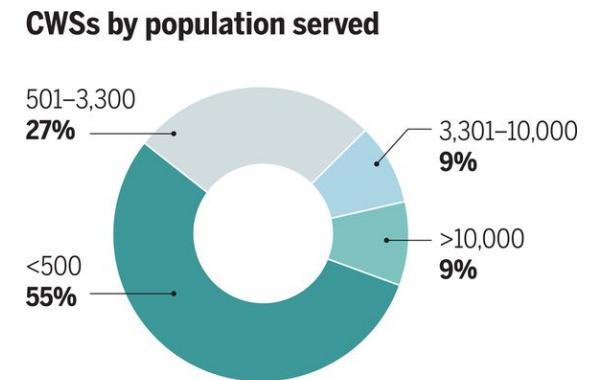
Created by Darkhorse Analytics

www.darkhorseanalytics.com

Salvaging the Pie (Sept. 26, 2014), via Joey Cherdarchuk at [Darkhorse Analytics](#)

Data Visualization

- Pie chart & stacked bar plot
 - What's going on here?
 - CWS = community water system

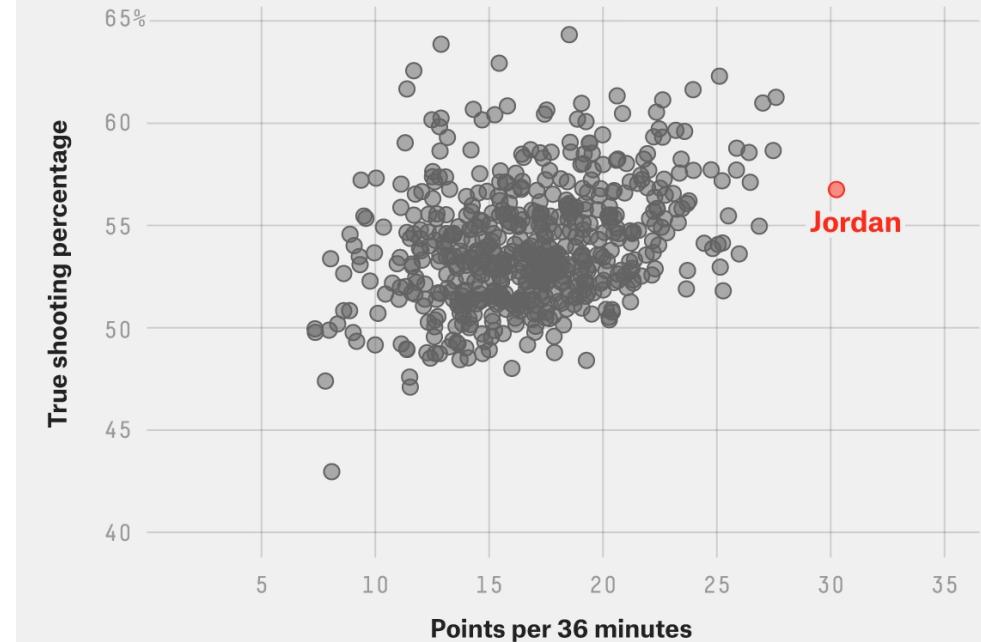


Data Visualization

- Scatter plot
 - Quantitative variable on one axis
 - Another quantitative variable on other axis
 - Observations shown as points

MJ pushed the boundaries of efficient scoring

Career true shooting percentage vs. points per 36 minutes (adjusted for pace) for NBA players with at least 15,000 minutes, 1976-2020



FiveThirtyEight

SOURCE: BASKETBALL-REFERENCE.COM

Why Michael Jordan Was The Best (April 17, 2020), via Neil Paine at FiveThirtyEight

Data Visualization

- Scatter plot
 - What's going on here?



Data Visualization

- **Line plot**

- Time on one axis (usually x-axis)
- Quantitative variable on other axis
- Line connecting each observation (usually the peaks in the line)
- Good for observing trends over time

Fans who skip shortened seasons eventually return

Average attendance at MLB regular-season games per year since 1970



Do Baseball's Labor Fights Drive Fans Away? (June 12, 2020), via Travis Sawchik at FiveThirtyEight

Data Visualization

- Line plots
 - What's going on here?

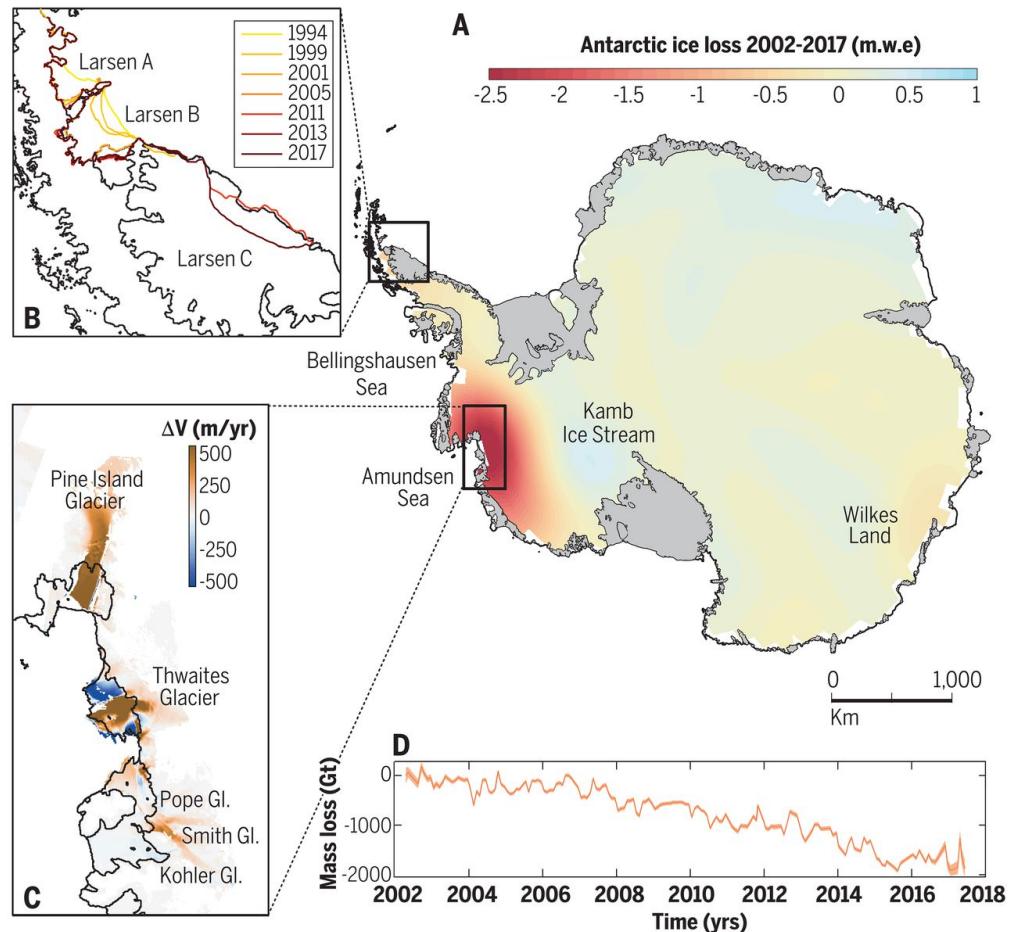
Thanksgiving dinner cost survey

Slight increases in many traditional Thanksgiving dinner ingredients were offset by cheaper turkey prices this year. The American Farm Bureau asks volunteers to record the prices of these items nationwide. Quantities of food purchased enough to feed a family of 10.

Item	Price	2018	2019	Change	Percent change
Turkey (16 pounds)	\$21.71	\$20.80		-\$0.91	-4%
Pumpkin pie mix (30 oz.)	\$3.33	\$3.32		-\$0.01	-0.3%
Whole milk (1 gallon)	\$2.92	\$3.10		+\$0.18	+6.16%
Veggie tray (1 pound)	\$0.75	\$0.79		+\$0.04	+5.33%
			\$2.01		

Data Visualization

- Line plots
 - What's going on here (part D)?



History, mass loss, structure, and dynamic behavior of the Antarctic Ice Sheet (March 20, 2020),

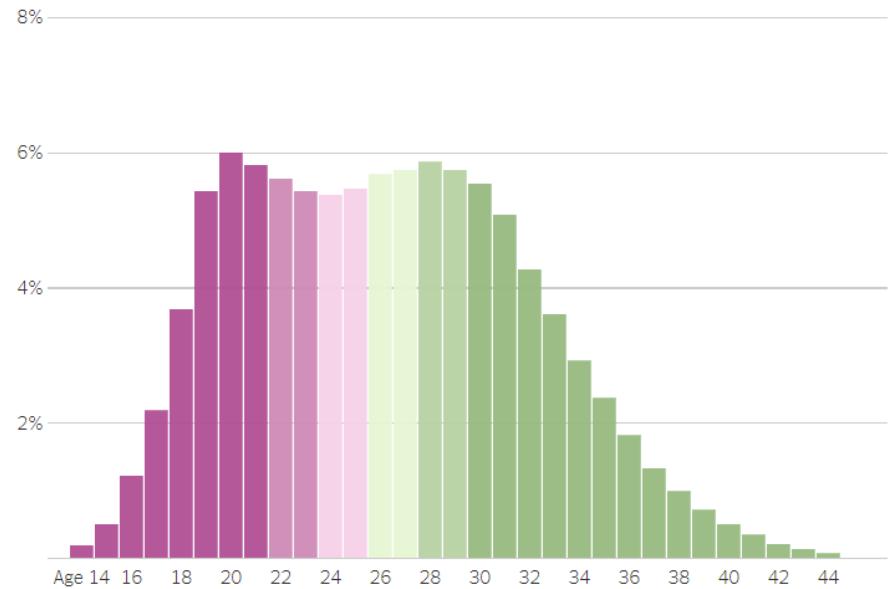
via Robin E. Bell & Helene Seroussi at [Science](#)

Data Visualization

- **Histogram**

- Like a bar plot, but instead of counts of observations in categories we have counts of observations in intervals (bins) of a quantitative variable
- You can change the interval width (and the # of bars) for histograms; you can't combine or divide categories for bar plots
- Good for figuring out the **data's shape**

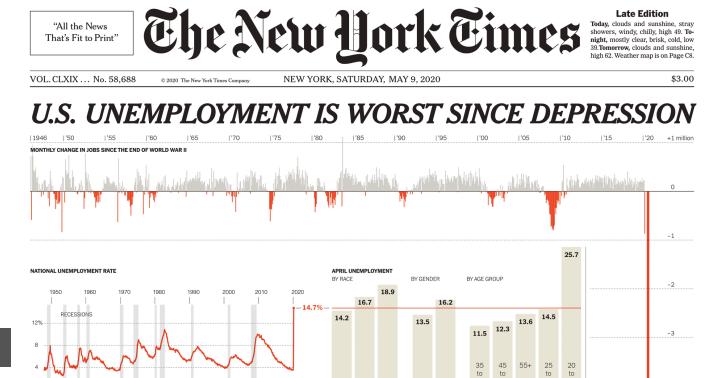
Ages of first-time mothers in 2016



The age of first-time mothers in 2016 (Aug. 4, 2018), via Quoctrung Bui at The New York Times

Data Visualization

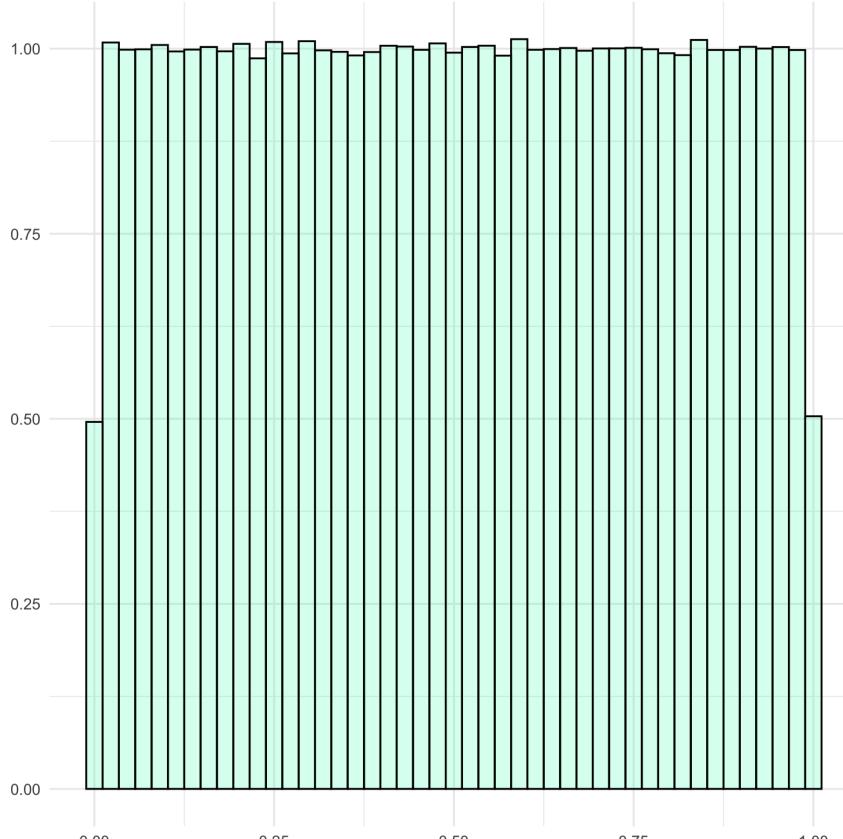
- Histogram
 - What's going on here?



Data Shapes

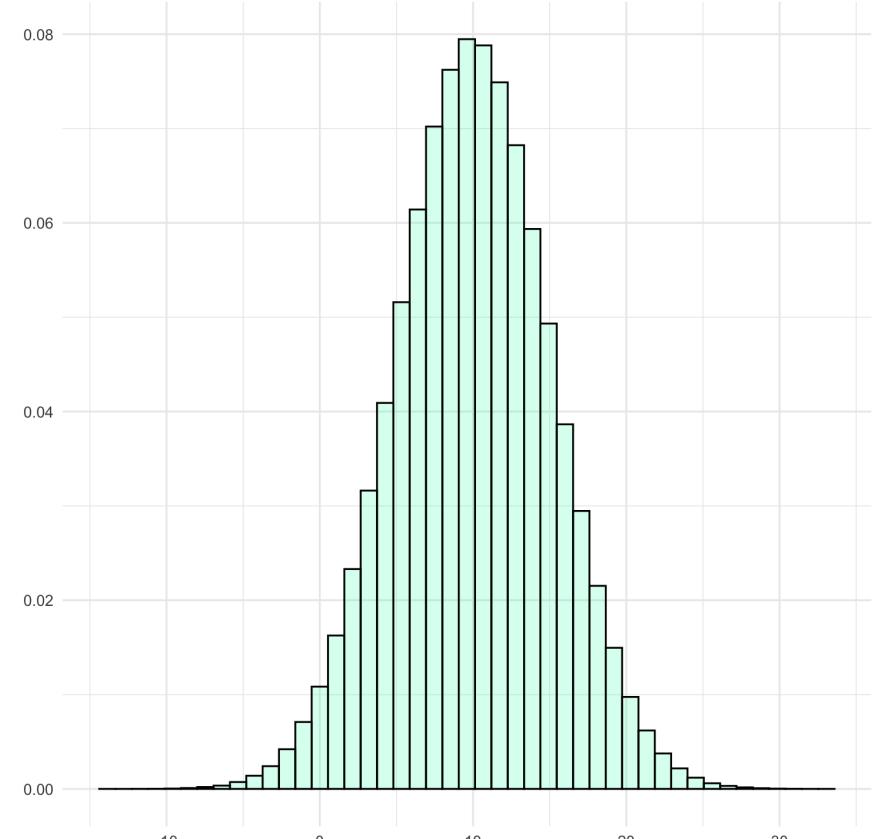
Uniform

- No defined peaks



Unimodal

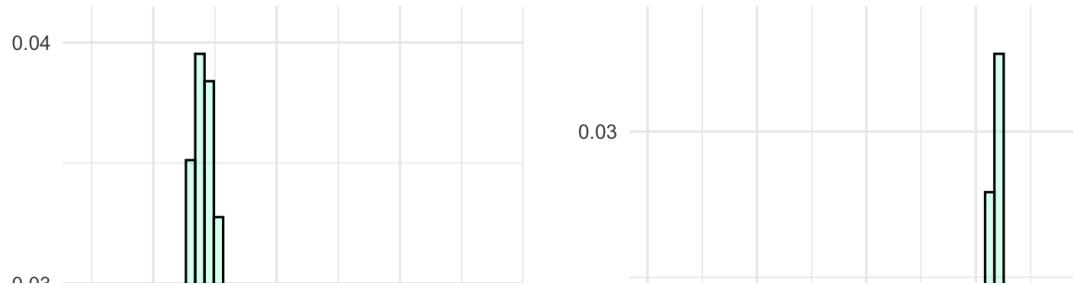
- One defined peak



Data Shapes

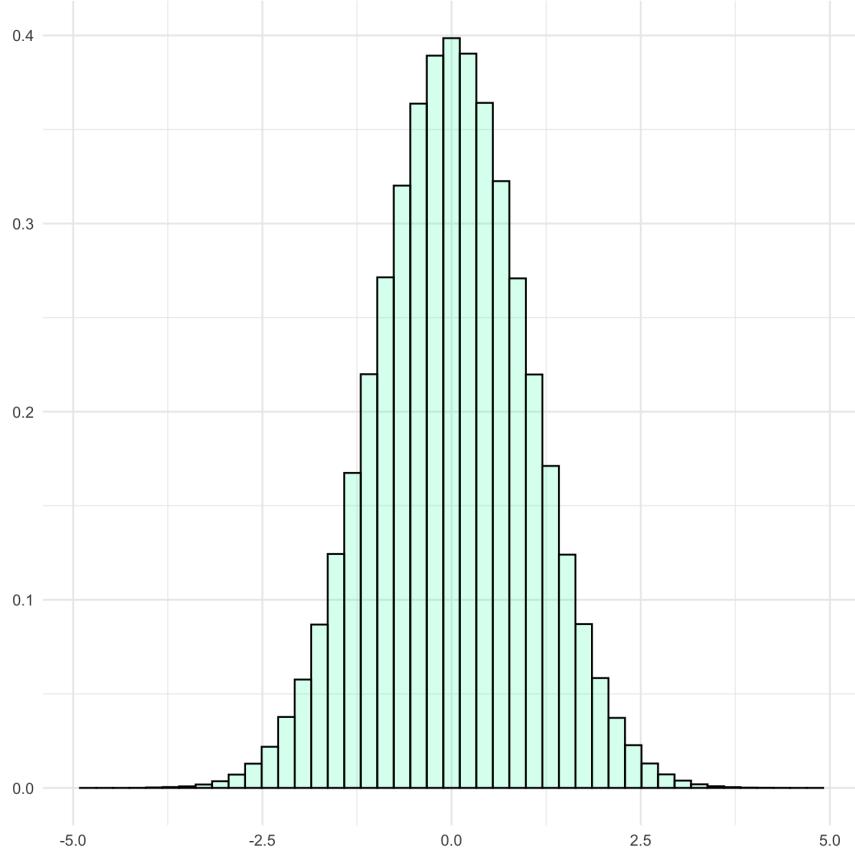
Multimodal

- Multiple defined peaks

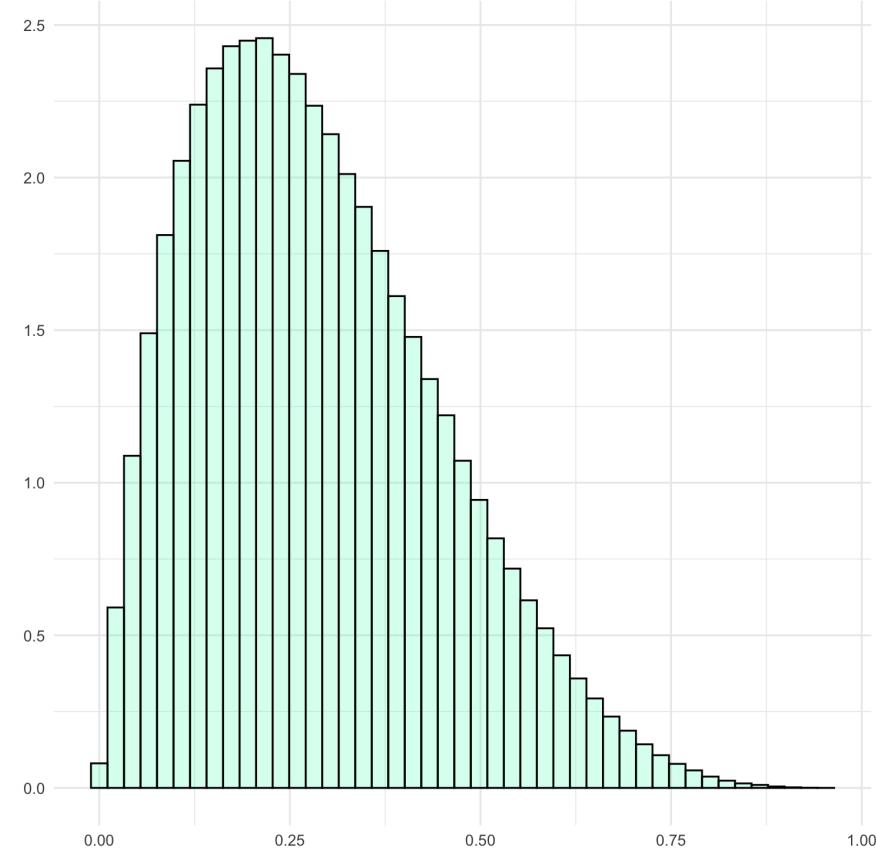


Data Shapes

Symmetrical

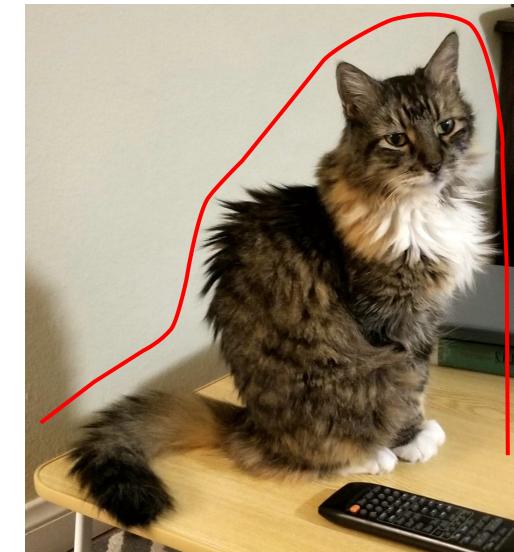


Skewed



Skew

- Skew is what happens when the mean and median aren't exactly the same
- Right (positive) skew
 - Mean to the right (**greater than**) median
- Left (negative) skew
 - Mean to the left (**less than**) median
- When in doubt, look for the "tail"



Left skew (and good bean)

Five Number Summary

- Getting a 10,000-foot view of the data's shape
- The five number summary contains:
 - Minimum
 - **1st quartile** (25th percentile)
 - 25% of data are below the 1st quartile
 - Median (50th percentile)
 - **3rd quartile** (75th percentile)
 - 75% of data are below the 3rd quartile
 - Maximum

Data Visualization

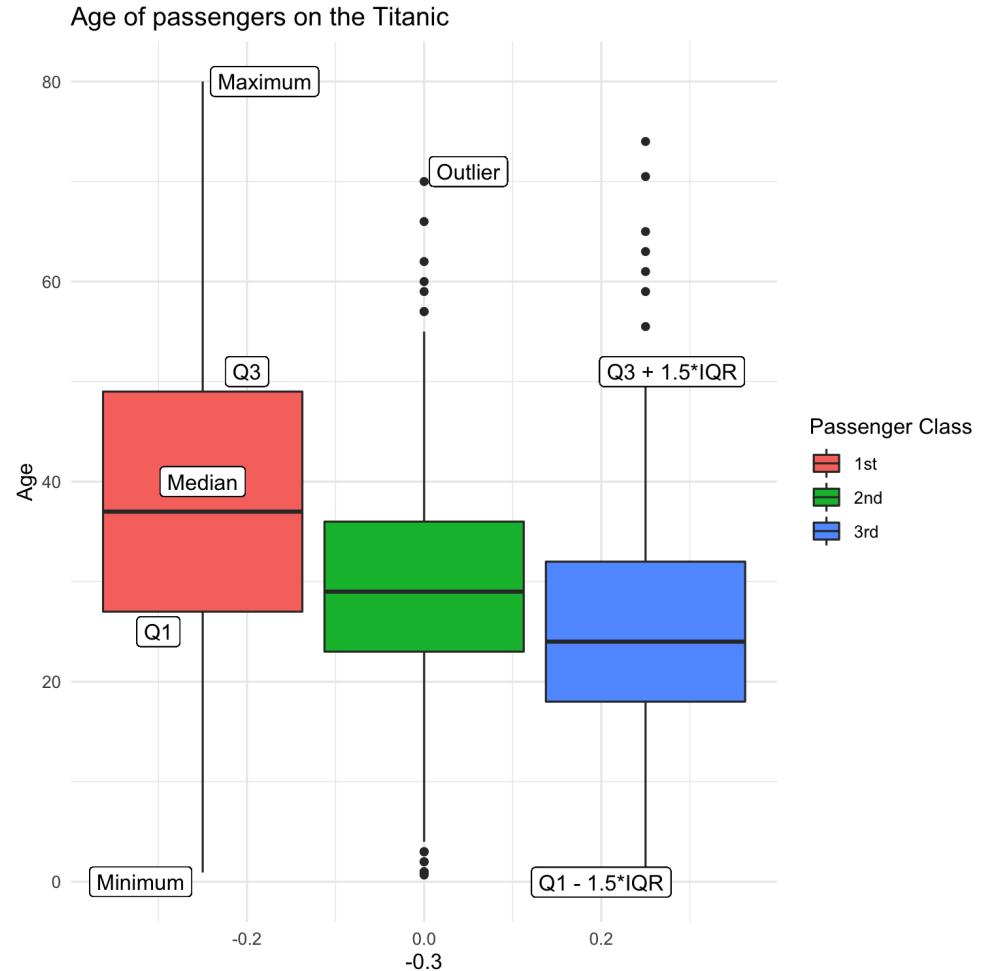
- **Boxplot**

- Bottom line of box at 1st quartile
- Top line of box at 3rd quartile
- Line in middle of box at median
- Box "whiskers" extend from minimum to maximum, or to "outlier fences"

Upper fence: $Q3 + 1.5 \times (Q3 - Q1)$

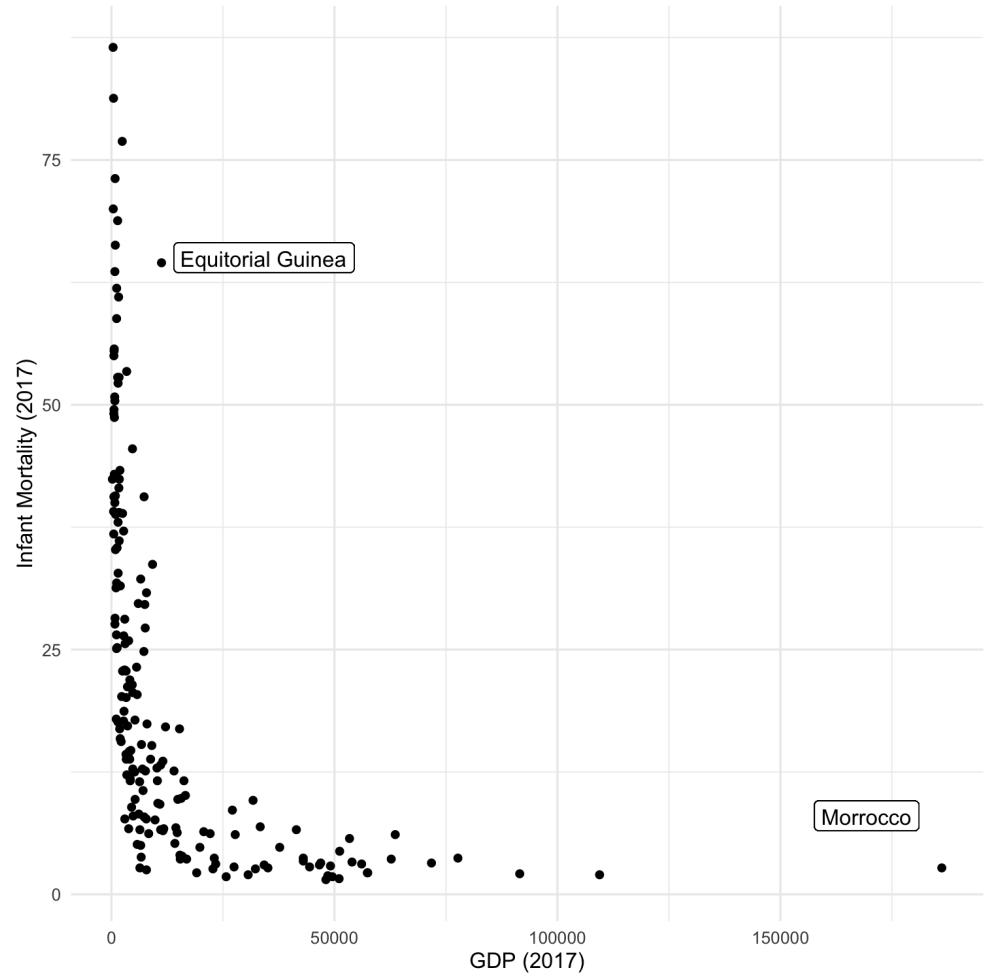
Lower fence: $Q1 - 1.5 \times (Q3 - Q1)$

- Dots are outliers
 - Observations above $Q3 + 1.5 \times IQR$ or below $Q1 - 1.5 \times IQR$



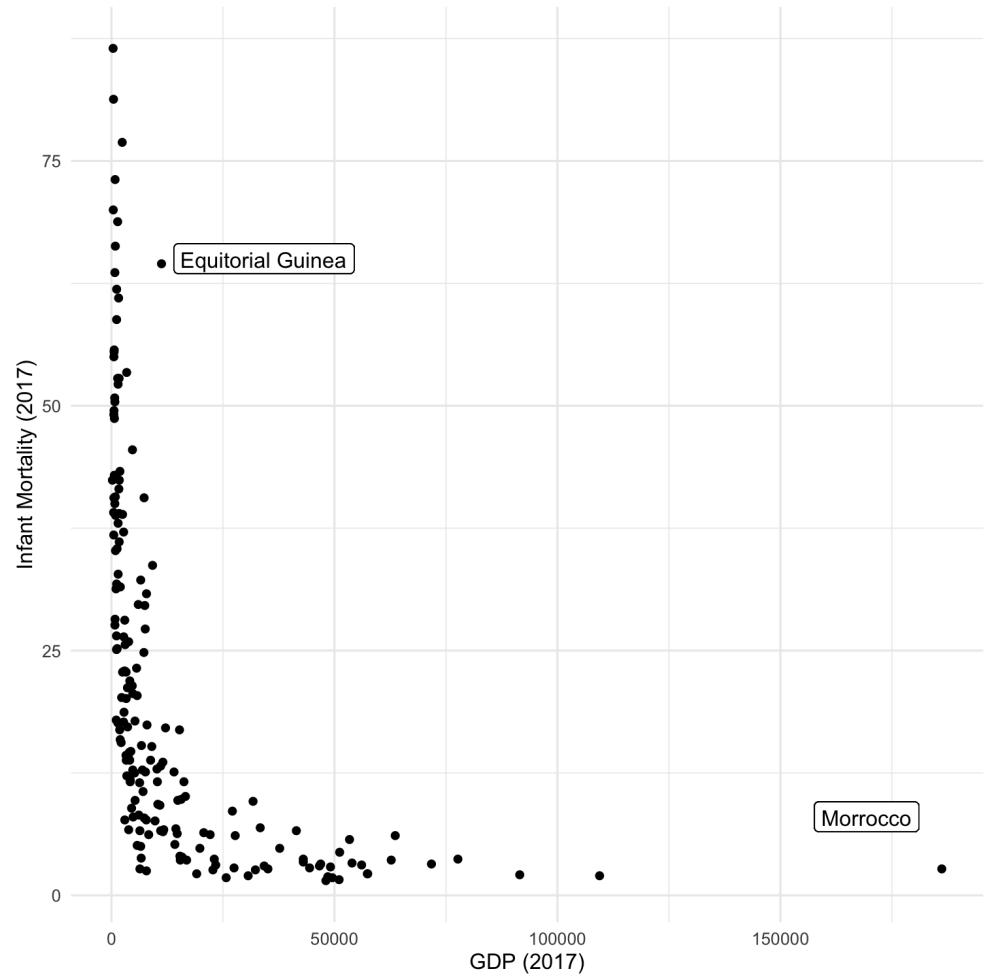
Outliers

- **Outliers** are data points that are strangely far away from the rest of the data
- Could signal that the values of those observations are wrong
 - Data entry error, recorded in inches instead of centimeters, etc.
- Could just be strange observations
- Outliers will create skew, pull your mean away from the median



Outliers

- **Outliers** are data points that are strangely far away from the rest of the data
- Could signal that the values of those observations are wrong
 - Data entry error, recorded in inches instead of centimeters, etc.
- Could just be strange observations
- Outliers will create skew, pull your mean away from the median
- What do we do with them?



Data Visualization

- Histogram vs. boxplot
 - Histograms give a more detailed view of data shape than boxplots
 - Because it's less detailed, boxplots can be better for comparing the data shapes of multiple groups

