

# L02-Visualizing data and describing data with numbers

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Today's objectives

L02-Visualizing  
data and  
describing data  
with numbers

Start working with data visualization - Introducing ggplot

Visualizations for  
categorical data

Describing your data with numbers:

Introducing ggplot

1. Investigate measures of centrality
  - ▶ mean and median, and when they're the same vs. different
2. Investigate measures of spread
  - ▶ IQR, standard deviation, and variance
3. Create a visualization of the “five number summary”
  - ▶ boxplots using ggplot
4. Calculate the variance and standard deviation

Visualizing quantitative  
variables  
Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere  
Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Visualizations for categorical data

# Visualization of categorical data

L02-Visualizing  
data and  
describing data  
with numbers

- ▶ What is the best way to visualize one categorical variable at a time?
- ▶ Generally speaking, it is not a good idea to use pie charts

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere  
Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

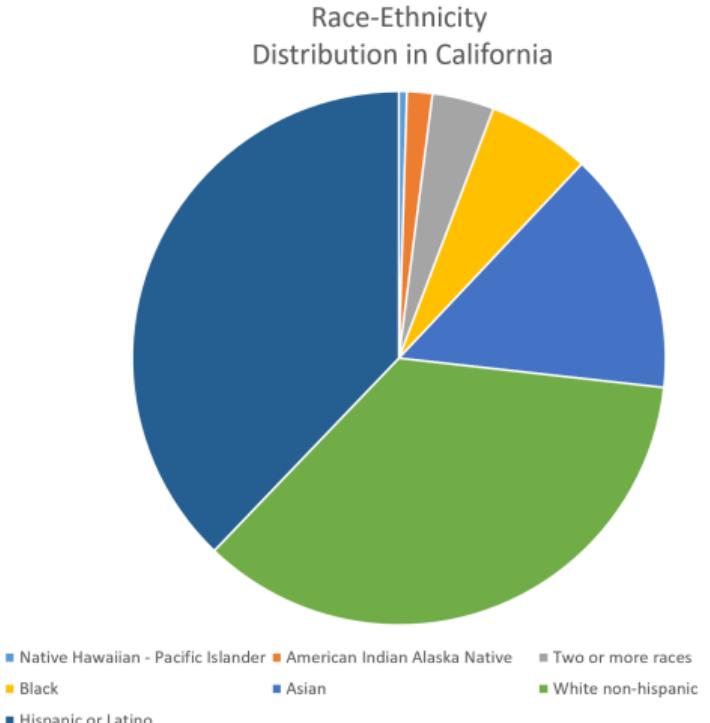
Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Visualization of categorical data

Can you judge the area of the slices?



Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

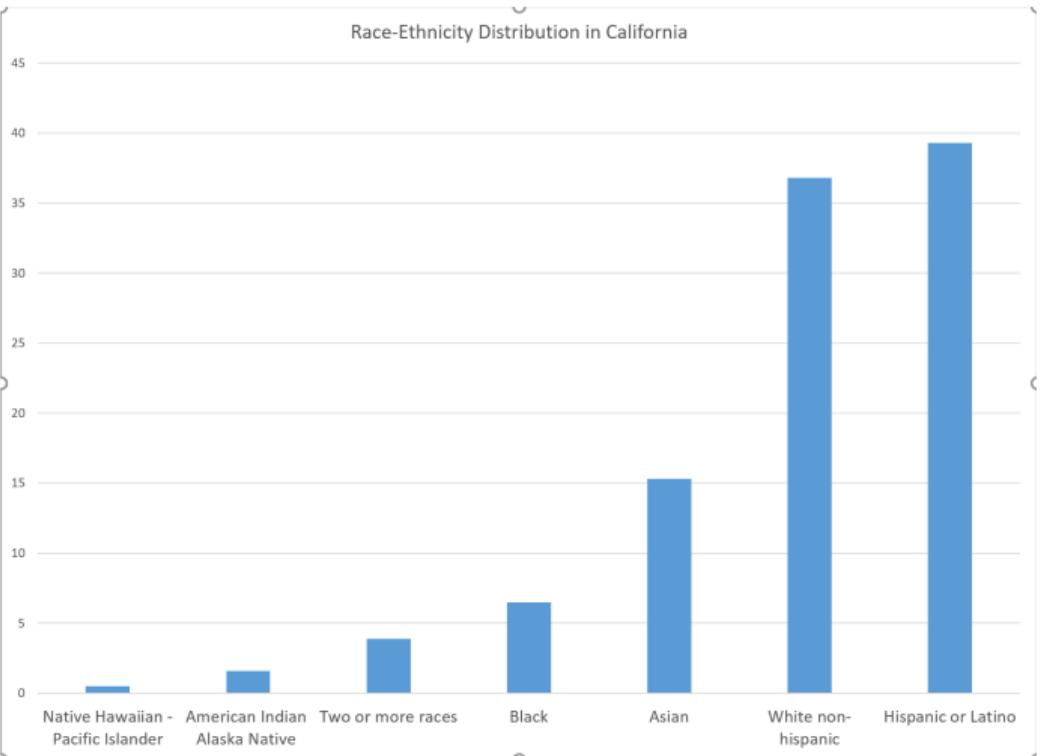
Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Visualization of categorical data

L02-Visualizing  
data and  
describing data  
with numbers



Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Visualization of categorical data

L02-Visualizing  
data and  
describing data  
with numbers

- ▶ We prefer bar graphs (also called bar charts) for the display of categorical data.
- ▶ Bar charts display the number or percent of data for each level of the categorical variable being plotted

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

## Example: infectious disease data

- ▶ Task: Make a bar chart of the percent of cases on infectious disease for each category of disease.
- ▶ First, read and view the infectious disease data from Baldi and Moore:

```
id_data <- read_csv("Ch01_ID-data.csv")
```

```
## Rows: 7 Columns: 4
## -- Column specification -----
## Delimiter: ","
## chr (2): disease, type
## dbl (2): number_cases, percent_cases
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this m
```

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Example: infectious disease data

id\_data

```
## # A tibble: 7 x 4
##   disease           type number_cases percent_cases
##   <chr>            <chr>      <dbl>        <dbl>
## 1 Chlamydia         STI       174557       66.4
## 2 Gonorrhea          STI       44974        17.1
## 3 Pertussis          Pertussis  11219        4.27
## 4 Campylobacteriosis Foodborne  7919        3.01
## 5 Early syphilis    STI       7191        2.74
## 6 Salmonellosis     Foodborne  5361        2.04
## 7 Other              Other      11559       4.40
```

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Example: infectious disease data

L02-Visualizing  
data and  
describing data  
with numbers

- ▶ Note the variables `number_cases` and `percent_cases`
- ▶ What do you want the bar chart to display? What is the x and y variables for a bar chart?

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Introducing ggplot

Visualizations for  
categorical data

## Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere  
Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

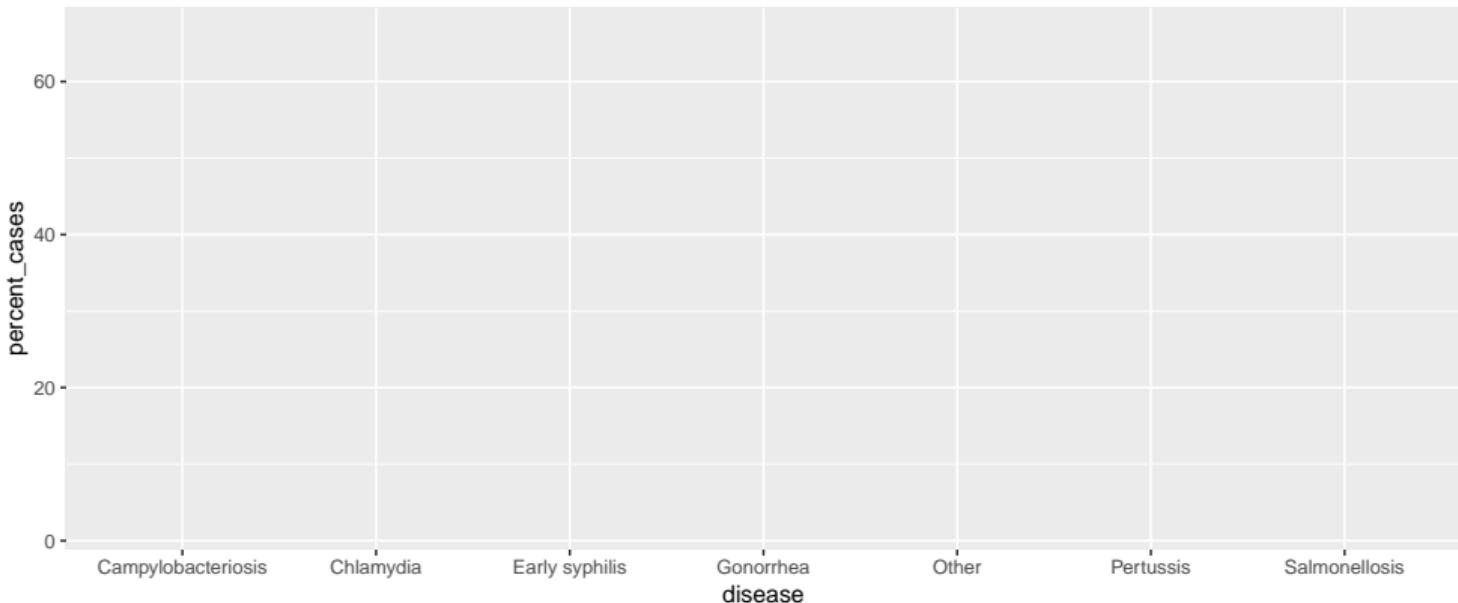
Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# First step to building a ggplot(): set up the canvas

- ▶ The first line of code below pulls in the ggplot package
- ▶ The second line of code below specifies the data set and what goes on the x and y axes



Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Next choose a function

L02-Visualizing  
data and  
describing data  
with numbers

- We will use a `geom_` function to create our chart

`ggplot()`'s `geom_bar()` makes a bar chart

Visualizations for  
categorical data

**Introducing ggplot**

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere  
Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Syntax for bar charts

L02-Visualizing  
data and  
describing data  
with numbers

```
ggplot(id_data, aes(x = disease, y = percent_cases)) +  
  geom_bar(stat = "identity")
```

stat = “identity” tells geom\_bar that we supplied a y variable that is exactly what we want to plot.

We do not need geom\_bar() to calculate the number or percent for us.

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

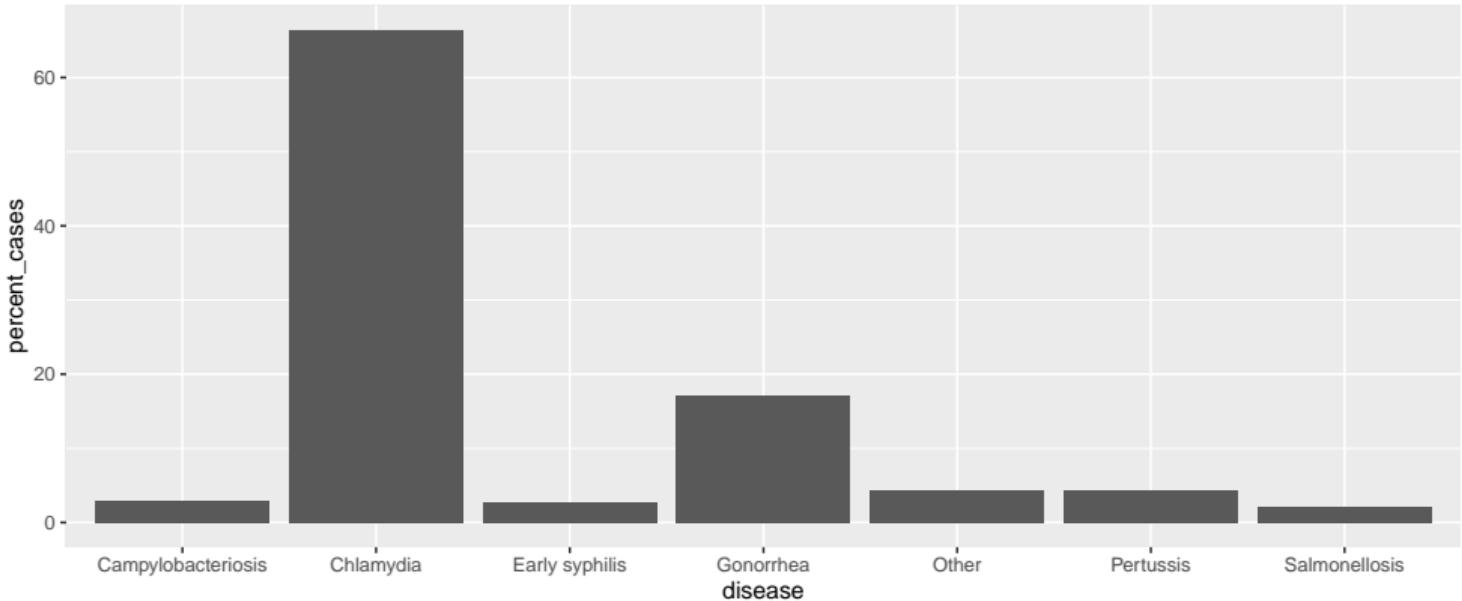
Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# ggplot()'s geom\_bar() makes a bar chart



Visualizations for  
categorical data

## Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

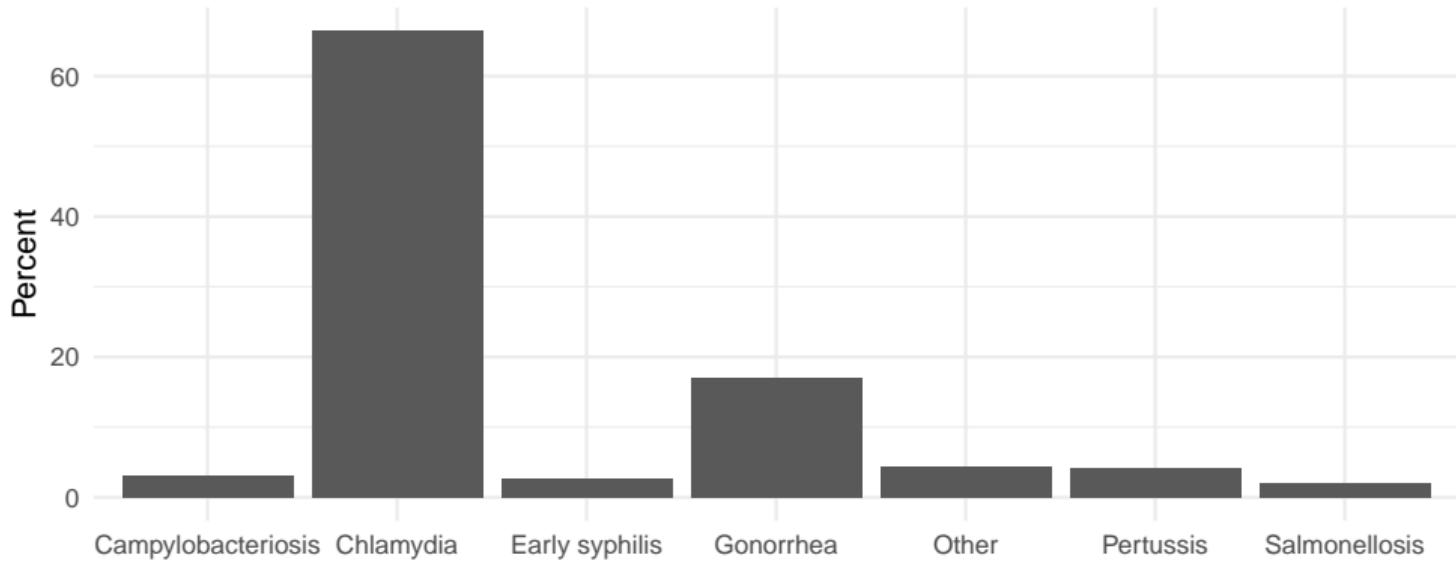
Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

## some additions to ggplot for style



base\_size controls the font size on these plots

theme\_minimal affects the “look” of the plot it removes the grey background and adds grey gridlines

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

fct\_reorder reorders disease according to value of percent\_cases

```
library(forcats)  ##notice that we have to pull in the package forcats here!
id_data <- id_data %>%
  mutate(disease_ordered = fct_reorder(disease, percent_cases, .desc = T))
```

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere  
[solution](#)

Mean vs Median: Outliers  
and sample size, skew,  
shape

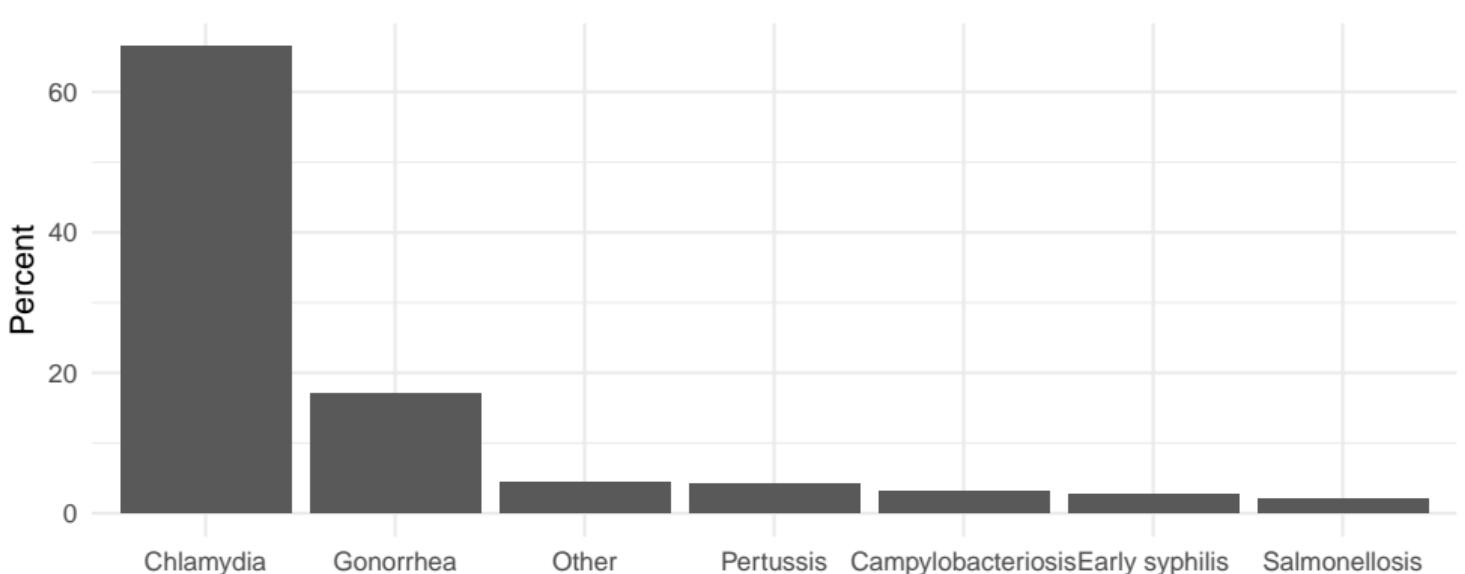
Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Re-ordered plot



Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

Use `aes(fill = type)` to link the bar's fill to the disease type

```
geom_bar(stat = "identity", aes(fill = type)) +  
theme(legend.position = "top")
```

Visualizations for  
categorical data

**Introducing ggplot**

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere  
Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

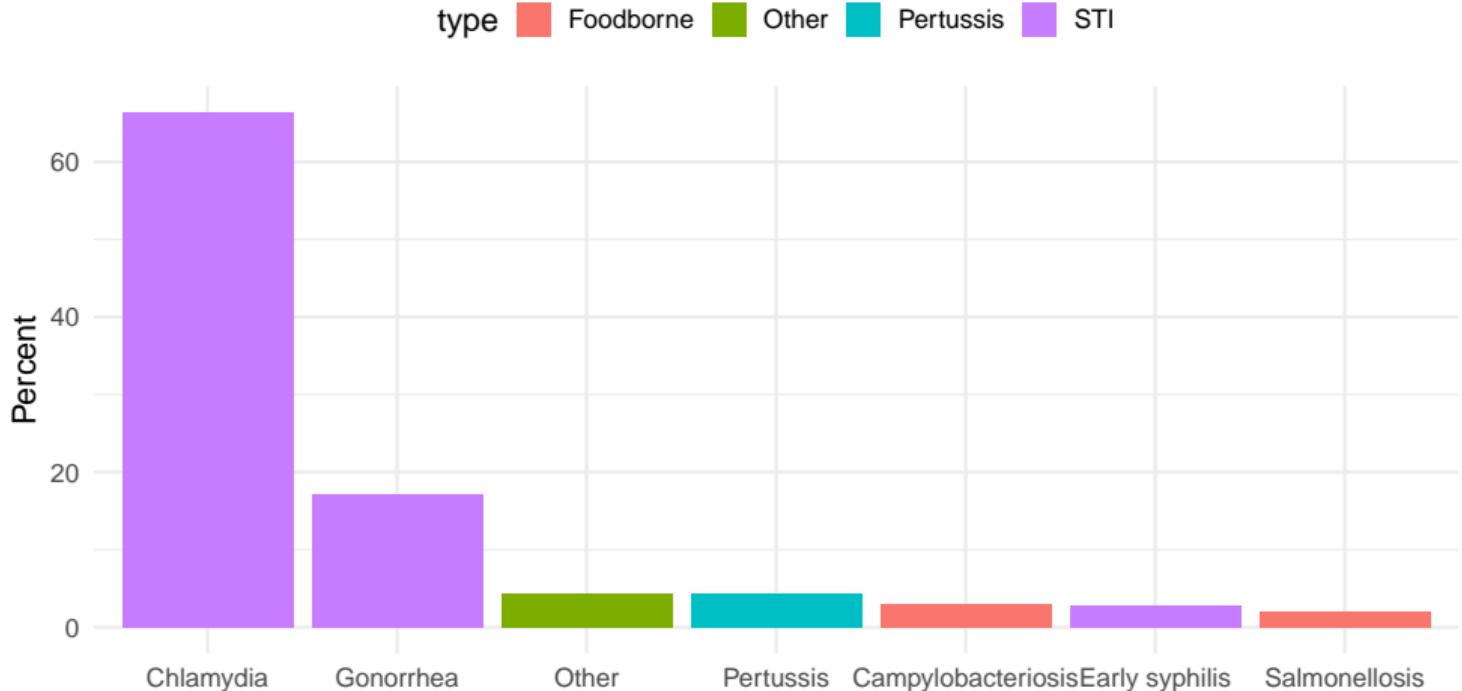
Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Use aes(fill = type) to link the bar's fill to the disease type



Visualizations for  
categorical data

## Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

Visualizations for  
categorical data

Introducing ggplot

**Visualizing quantitative  
variables**

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Visualizing quantitative variables

# Visualize quantitative variables using histograms

L02-Visualizing  
data and  
describing data  
with numbers

- ▶ Histograms look a lot like bar charts, except that the bars touch because the underlying scale is continuous and the order of the bars matters
- ▶ In order to make a histogram, the underlying data needs to be **binned** into categories and the number or percent of data in each category becomes the height of each bar.
- ▶ the **bins** devide the entire range of data into a series of intervals and counts the number of observations in each interval
- ▶ the intervals must be consecutive and non-overlapping and are almost always chosen to be of equal size

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

## Example: opioid state prescription rates

L02-Visualizing  
data and  
describing data  
with numbers

- ▶ The textbook gives an example using data from 2012.
- ▶ In the data folder, there is updated data from 2018. It came from the paper: “Opioid Prescribing Rates by Congressional Districts, United States, 2016”, by Rolheiser et al. link

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere  
Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

## Example: opioid state prescription rates

Problem: To determine the extent to which opioid prescribing rates vary across US congressional districts.

Plan: In an observational cross-sectional framework using secondary data, they constructed 2016 congressional district-level opioid prescribing rate estimates using a population-weighted methodology.

Data: In the data structure we have State as the unit of analysis, and measured prescription rates as the variable of interest

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere  
Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

## Example: opioid state prescription rates

```
opi_data <- read.csv("Ch01_opioid-data.csv")  
head(opi_data)
```

##	Rank	State	Mean	Median	SD	Min	Max	Num_Districts
## 1	1	AL	121.31	113.09	21.87	105.58	166.69	7
## 2	2	AR	115.22	115.13	8.59	104.80	125.79	4
## 3	3	TN	108.12	108.26	19.16	73.60	133.00	9
## 4	4	MS	105.64	106.25	17.36	83.90	126.14	4
## 5	5	LA	98.38	98.88	10.34	83.22	112.65	6
## 6	6	KY	98.13	85.76	26.72	77.62	147.00	6

- ▶ Mean provides the mean prescribing rate per 100 individuals. Thus, a mean of 121.31 implies that in Alabama, there were 121.31 opioid prescriptions per 100 persons, an average across the 7 congressional districts.

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Histogram of opioid prescription rates

L02-Visualizing  
data and  
describing data  
with numbers

- ▶ Task: Make a histogram of the average prescribing rates across US states
- ▶ What is the x variable? What is the y variable?
- ▶ What geom should be used?

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Histogram of opioid prescription rates - default is 30 bins

L02-Visualizing  
data and  
describing data  
with numbers

```
ggplot(data = opi_data, aes(x = Mean)) +  
  geom_histogram(col = "white") +  
  labs(x = "Mean opioid prescription rate (per 100 individuals)",  
       y = "Number of states") +  
  theme_minimal(base_size = 15)
```

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

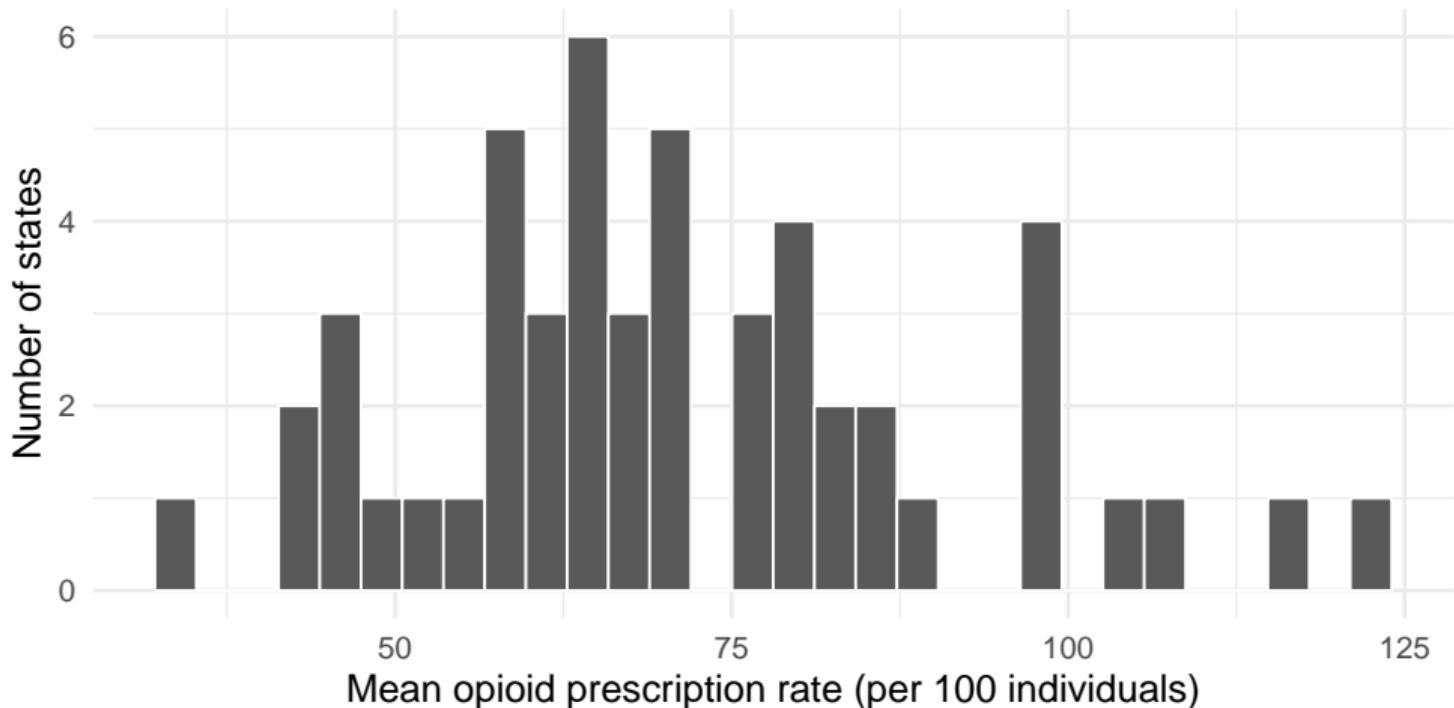
Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Histogram of opioid prescription rates

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

same graph, change the bins `geom_histogram(binwidth = 5)`

```
ggplot(data = opi_data, aes(x = Mean)) +  
  geom_histogram(col = "white", binwidth = 5) +  
  labs(x = "Mean opioid prescription rate (per 100 individuals)",  
       y = "Number of states") +  
  theme_minimal(base_size = 15)
```

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

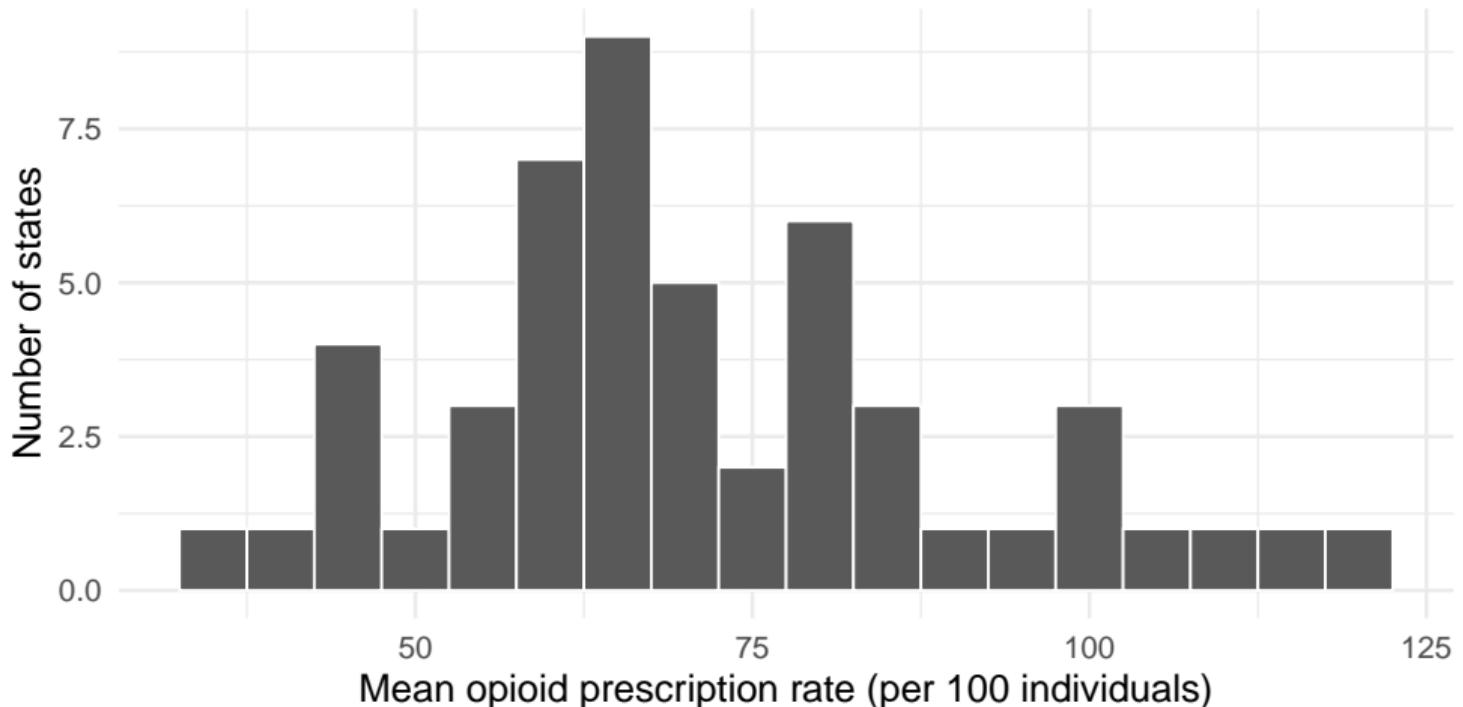
Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

same graph, change the bins `geom_histogram(binwidth = 5)`



Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

change the bins again `geom_histogram(binwidth = 10)`

```
ggplot(data = opi_data, aes(x = Mean)) +  
  geom_histogram(col = "white", binwidth = 10) +  
  labs(x = "Mean opioid prescription rate (per 100 individuals)",  
       y = "Number of states") +  
  theme_minimal(base_size = 15)
```

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

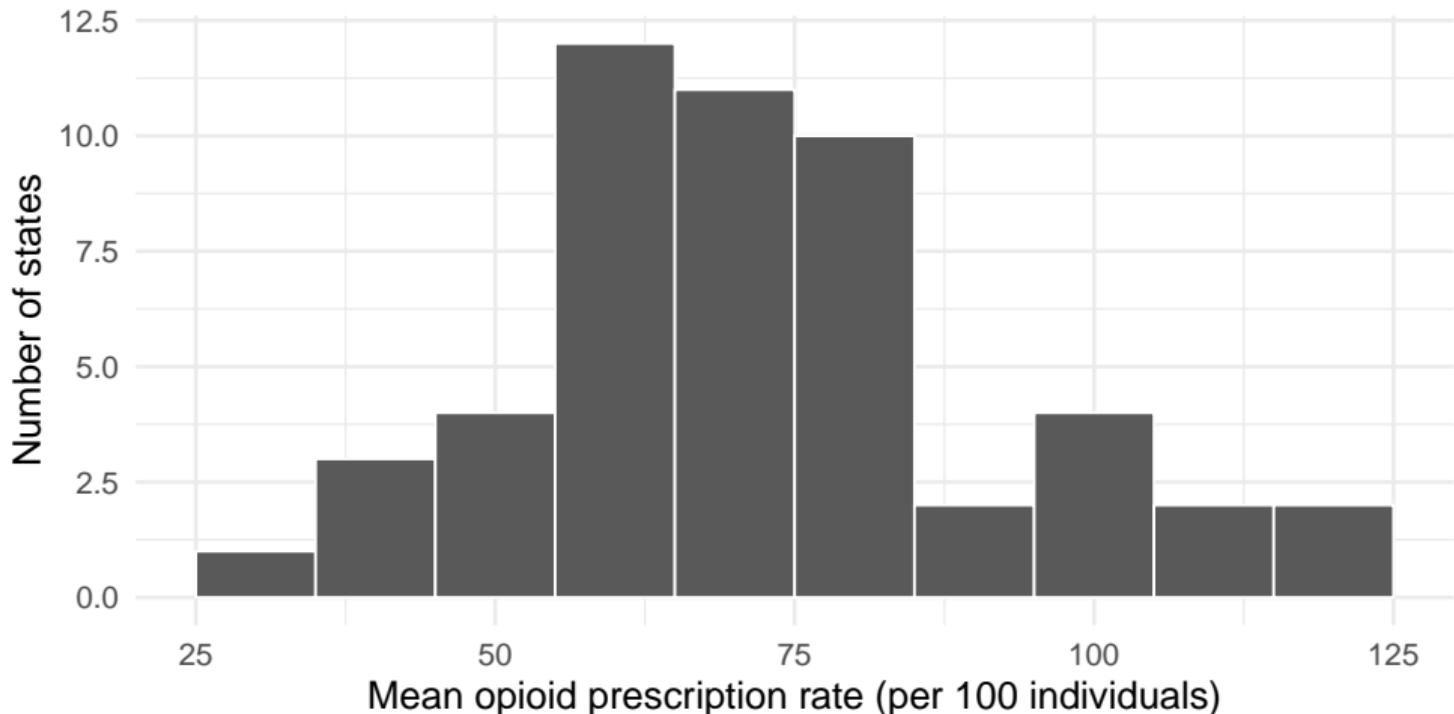
Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

change the bins again `geom_histogram(binwidth = 10)`



Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

## Describing your distribution - what are we looking for?

# Shape, Center, Spread

L02-Visualizing  
data and  
describing data  
with numbers

- ▶ When we examine histograms, we can make comments on a distribution's:
  - ▶ Shape: Is the distribution **symmetric** or **skewed** to the left or right?
  - ▶ Center: Does the histogram have one peak (**unimodal**), or two (**bimodal**) or more?
  - ▶ Spread: How spread out are the values? What is the range of the data?
  - ▶ Outliers: Do any of the measurements fall outside of the range of most of the data points?

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

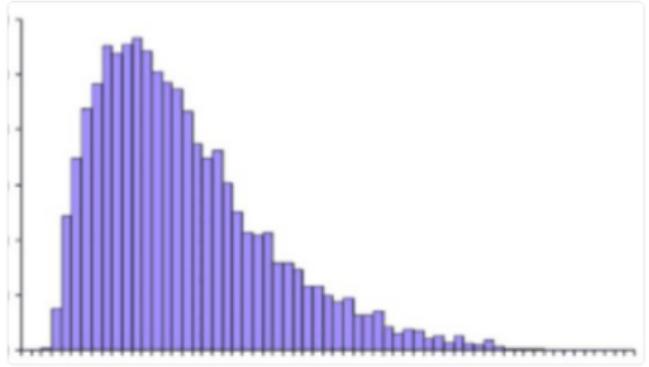
Box plots

# Is this skewed left or skewed right?

L02-Visualizing  
data and  
describing data  
with numbers



Jesse Singal @jessesingal · 13h  
THIS IS NOT NORMAL



72 323 1.5K

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

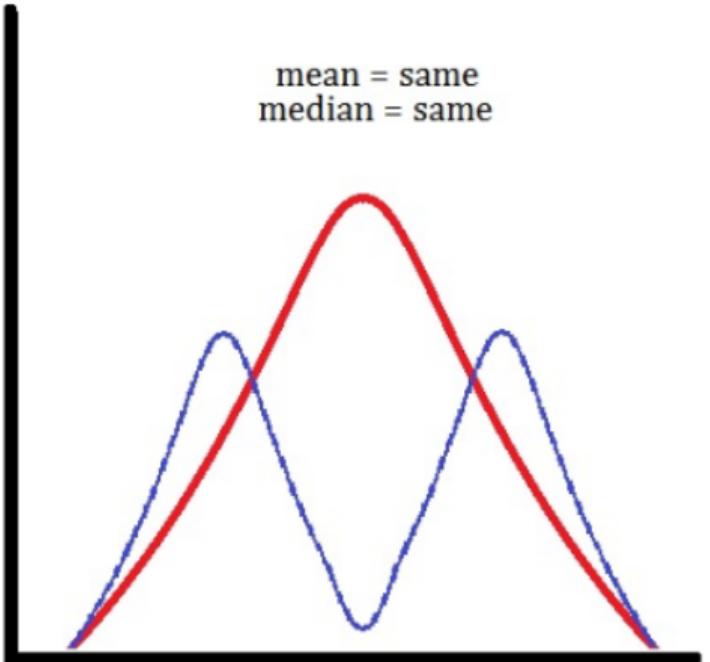
Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Center - one hump or two?

L02-Visualizing  
data and  
describing data  
with numbers



Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

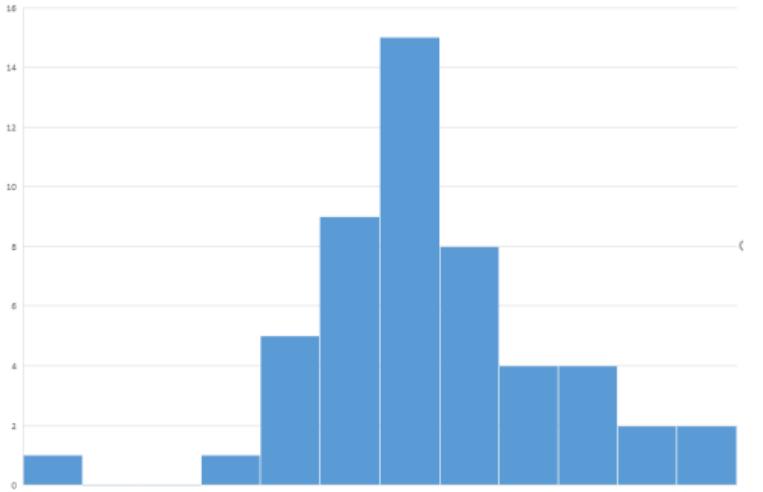
Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Outlier

L02-Visualizing  
data and  
describing data  
with numbers



Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

**Measures of central  
tendency**

Statistics is Everywhere  
Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

## Measures of central tendency

# Measures of central tendency

L02-Visualizing  
data and  
describing data  
with numbers

- Most common: [mean](#) and [median](#)

[Visualizations for  
categorical data](#)

[Introducing ggplot](#)

[Visualizing quantitative  
variables](#)

[Describing your distribution  
- what are we looking for?](#)

[Measures of central  
tendency](#)

[Statistics is Everywhere](#)

[Discussion](#)

[Mean vs Median: Outliers  
and sample size, skew,  
shape](#)

[Measures of spread](#)

[Example: Hospital cesarean  
delivery rates](#)

[Sample variance and  
standard deviation](#)

[Box plots](#)

# The arithmetic mean

L02-Visualizing  
data and  
describing data  
with numbers

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

$$\bar{x} = \sum_{i=1}^n \frac{x_i}{n}$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# The median

L02-Visualizing  
data and  
describing data  
with numbers

- ▶ Half of the measurements are larger and half are smaller.
  - ▶ What is the median if there is an odd number of observations?
  - ▶ An even number?

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere  
Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Statistics is Everywhere

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

**Statistics is Everywhere**

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

San Francisco

## Apartments for rent in San Francisco: What will \$3,400 get you?



From Hoodline.com

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Bay Area rent

L02-Visualizing  
data and  
describing data  
with numbers

Business > Real Estate > News

## Here's how much rent has gone up in the Bay Area since 2010

In San Francisco, tenants are paying 70 percent more than in 2010

[Facebook](#) [Twitter](#) [Email](#) [Print](#)



Now sitting at \$3,680, average rent in San Francisco has soared 70 percent since 2010 while home prices climbed an eye-popping 95 percent and median income crept up a comparatively modest 61 percent. Across the bay in Oakland, rent climbed even more — 108 percent. [Mercury News article](#)

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Discussion

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

**Discussion**

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# When are these measures approximately equal?

- ▶ Answer: When the data has one peak and is roughly **symmetric**
  - ▶ In this case, the mean  $\approx$  median, so provide either one in a summary
- ▶ **Skewed data**
  - ▶  $\text{mean} \neq \text{median}$
  - ▶ Right-skewed data will commonly have a \_\_\_\_\_ mean than median
  - ▶ Left-skewed data will commonly have a \_\_\_\_\_ mean than median
  - ▶ Which statistic should we report?

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

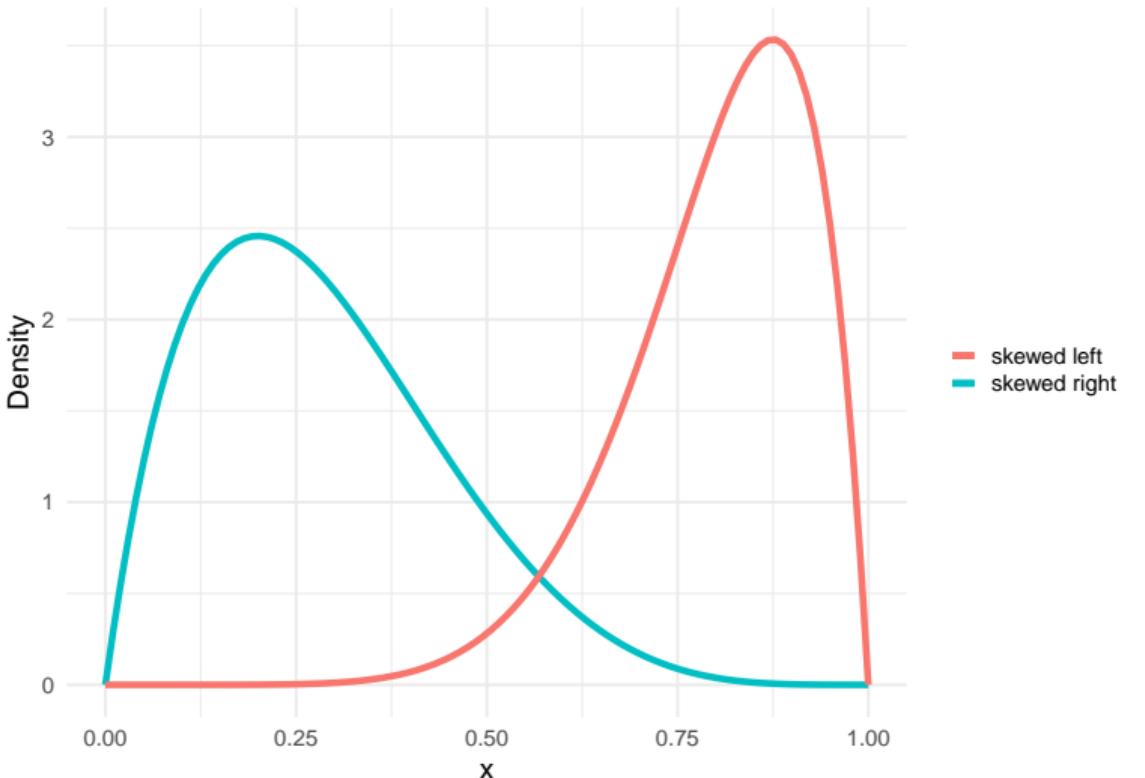
Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Skewed data



Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Apartment rent in SF

L02-Visualizing  
data and  
describing data  
with numbers

Problem: We want to understand how much it costs for a new resident to rent a 1 bedroom apartment in San Francisco

Plan: Take a sample of 1000 apartment units listed for rent (currently available) and ask the rental price (excluding utilities)

Data: Here I will present data that I simulated in r using a mean value published on [rentjungle.com](http://rentjungle.com) - you will not be expected to do this or be tested on it.

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

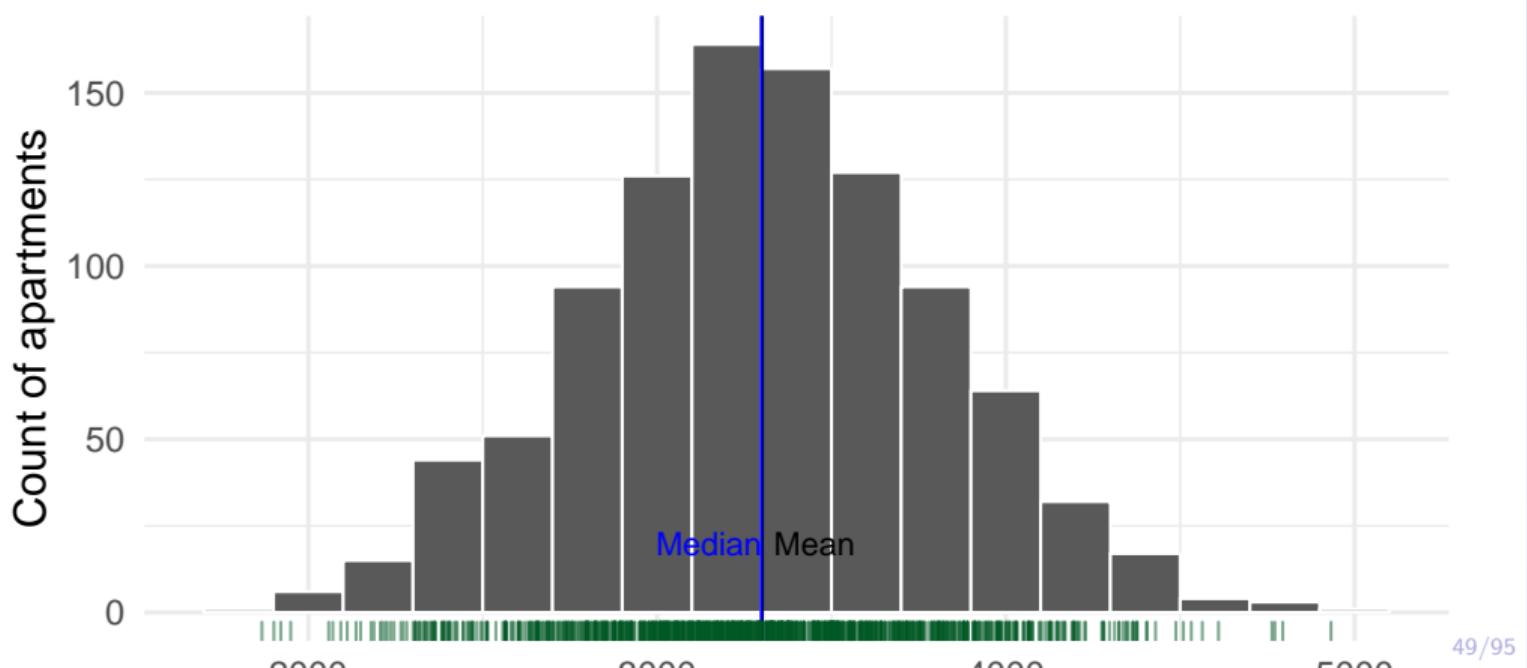
Sample variance and  
standard deviation

Box plots

## Example: Apartment rent in SF

Suppose that the distribution of rent prices looked like this. The green ticks underneath the histograms shows you the exact rent values that contribute data to each bin.

Symmetric distribution in rental prices (\$)



Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Example: Apartment rent in SF

L02-Visualizing  
data and  
describing data  
with numbers

From last lecture: We describe this distribution in terms of center, shape and spread:

- ▶ Center: Where is the center of the distribution?
- ▶ Shape: Is this distribution unimodal or bimodal?
- ▶ Spread: How much variability is there between the lowest and highest rent values?

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Example: Apartment rent in SF

Summarizing numerically: Center:

```
# in base R
```

```
mean(rent_data[, "sym"])
```

```
## [1] 3301.662
```

```
median(rent_data[, "sym"])
```

```
## [1] 3298.832
```

*# using the summarize function and a pipe operator*

```
rent_data %>% summarize(  
  mean=mean(sym),  
  median = median(sym))
```

```
##      mean    median
```

```
## 1 3301.662 3298.832
```

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

**Mean vs Median: Outliers  
and sample size, skew,  
shape**

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

## Mean vs Median: Outliers and sample size, skew, shape

# When are the mean and median approximately equal?

L02-Visualizing  
data and  
describing data  
with numbers

- ▶ If your data has one peak (unimodal), is roughly symmetric, and does not have outliers
  - ▶  $\text{mean} \approx \text{median}$ , so provide either one in a summary

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

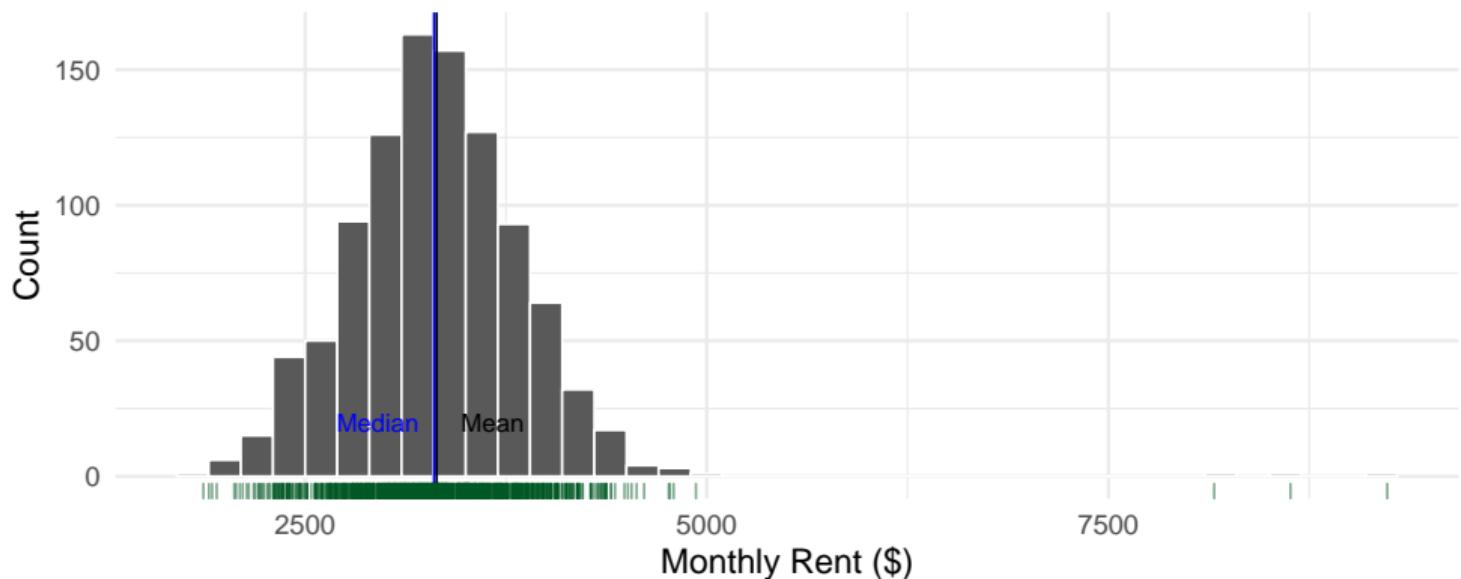
Sample variance and  
standard deviation

Box plots

## Example: Apartment rent in SF

Now suppose that there were three rents within the data set with much larger values than the rest of the distribution. Here is the plot for this updated data.

Symmetric, but with outliers on the right, n=1000



- With 1000 sampled points the outliers do not have a large effect on the mean

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variablesDescribing your distribution  
- what are we looking for?Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

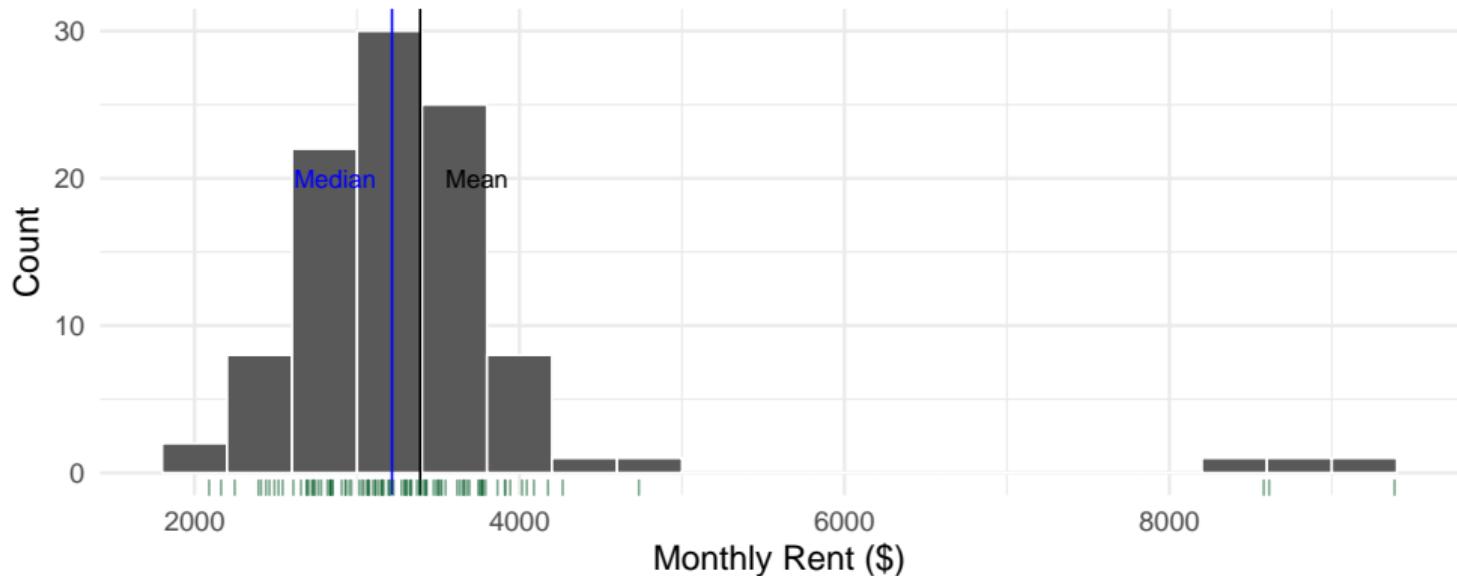
Example: Hospital cesarean  
delivery ratesSample variance and  
standard deviation

Box plots

## Example: Apartment rent in SF

Imagine instead, there were only 100 sampled points. Here, the outliers have a larger effect on the mean. **The mean is not resistant to outliers.**

Symmetric, but with outliers on the right, n=100



Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

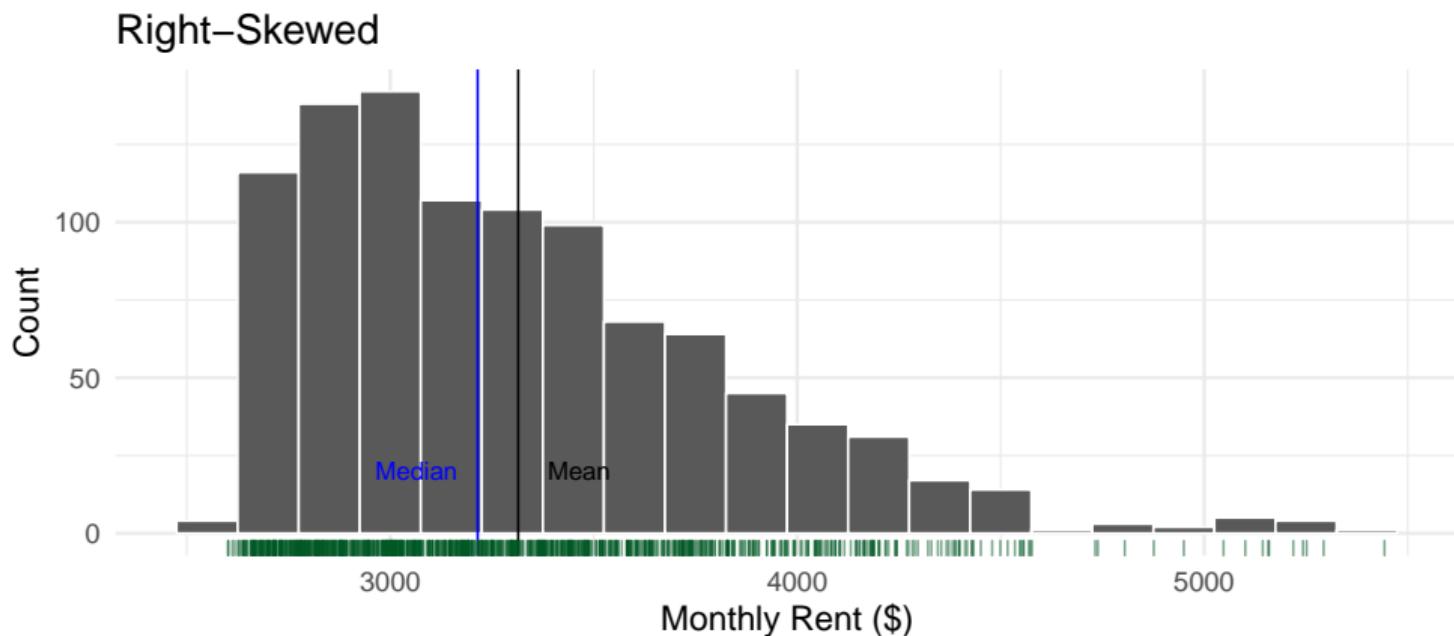
Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

## Example: Apartment rent in SF

Consider instead what happens if there are many high-end apartments in the area. Here is the histogram of data for this example:



Why is the mean larger than the median in this case?

Skewed data -  $\text{mean} \neq \text{median}$ . Data with a long right tail will commonly

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

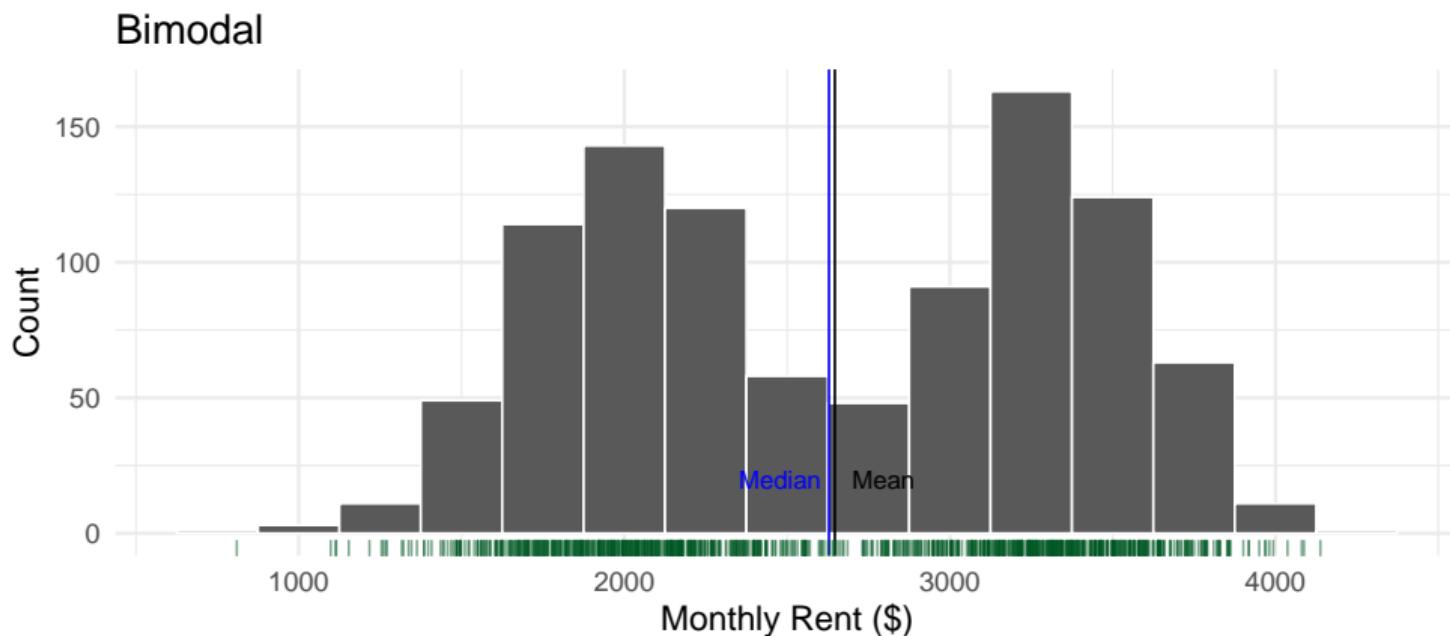
Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

## Example: Apartment rent in SF

Now, suppose that the sample of estimates did not look like the distribution in the previous example. Instead, it looked like this:



Describe the distribution. How does it differ from the first plot? Would you want to provide the mean or median for these data?

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variablesDescribing your distribution  
- what are we looking for?Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery ratesSample variance and  
standard deviation

Box plots

# Summary of measures of central tendency

L02-Visualizing  
data and  
describing data  
with numbers

- ▶ The mean and median are similar when the distribution is symmetric
- ▶ Outliers affects the mean and pull it towards their values. But they do not have a large effect on the median.
- ▶ Skewed distributions also pull the mean out into the tail.
- ▶ Measures of central tendency are not very helpful in multi-modal distributions

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Measures of spread

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

**Measures of spread**

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# The inter-quartile range (IQR)

L02-Visualizing  
data and  
describing data  
with numbers

- ▶ Q1 is the 1st quartile/the 25th percentile.
  - ▶ 25% of individuals have measurements below Q1.
- ▶ Q2 is the 2nd quartile/the 50th percentile/the median.
  - ▶ 50% of individuals have measurements below Q2.
- ▶ Q3, the 3rd quartile/the 75th percentile.
  - ▶ 75% of individuals have measurements below Q3.
- ▶ Q1-Q3 is called the **inter-quartile range (IQR)**.
  - ▶ What percent of individuals lie in the IQR?
- ▶ Know how to find Q1, Q2, and Q3 by hand for small lists of numbers

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere  
Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Quantiles using R

```
quantile(variable, 0.25)
```

```
rent_data %>% summarize(  
  Q1 = quantile(sym, 0.25),  
  median = median(sym),  
  Q3 = quantile(sym, 0.75)  
)
```

```
##           Q1    median        Q3  
## 1 2981.445 3298.832 3629.012
```

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# R's quantile function: Note

L02-Visualizing  
data and  
describing data  
with numbers

- ▶ `quantile(variable, 0.25)` will not always give the exact same answer you calculate by hand
- ▶ The R function is optimized for its statistical properties and is slightly different than the book's method
- ▶ To get the exact same answer as by hand use `quantile(data, 0.25, type = 2)`
- ▶ You may use either one in this class. Most commonly, people do not specify `type=2`

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Another measure of spread: The (full) range

L02-Visualizing  
data and  
describing data  
with numbers

- The difference between the **minimum** and **maximum** value

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

**Measures of spread**

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Concise information about spread and center: The five number summary

L02-Visualizing  
data and  
describing data  
with numbers

- ▶ The five number summary (min, Q1, median, Q3, max) is a quick way to communicate a distribution's center and spread.
- ▶ Based on the summary you can describe the full range of a dataset, where the middle 50% of the data lie, and the middle value.

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere  
Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# dplyr's summarize() to calculate the five number summary

Using our original example of rent data:

```
rent_data %>% summarize(  
  min = min(sym),  
  Q1 = quantile(sym, 0.25),  
  median = median(sym),  
  Q3 = quantile(sym, 0.75),  
  max = max(sym))
```

```
##           min         Q1      median         Q3        max  
## 1 1866.829 2981.445 3298.832 3629.012 4932.54
```

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere  
Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

## Example: Hospital cesarean delivery rates

## Example: Hospital cesarean delivery rates

These data were provided by the first author (Kozhimannil) of a manuscript published in the journal *Health Affairs*. link

From the article: Cesarean delivery is the most commonly performed surgical procedure in the United States, and cesarean rates are increasing. In its Healthy People 2020 initiative, the Department of Health and Human Services put forth clear, authoritative public health goals recommending a 10 percent reduction in both primary and repeat cesarean rates, from 26.5 percent to 23.9 percent, and from 90.8 percent to 81.7 percent, respectively.

A targeted approach to achieving such reductions might focus on hospitals with exceptionally high cesarean rates. However, adopting such a strategy requires quantification of hospital-level variation in cesarean delivery rates.

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

## Example: Hospital cesarean delivery rates

Problem: To characterize the variation in cesarean rates between Hospitals in the United States

Plan: Collect existing data from a variety of institutions for one year and compare rates of cesarean delivery. They also looked at cesarean rates among only low risk births at each institution. Why might this be important?

Data: For this article, they worked with 2009 data from 593 US hospitals nationwide

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere  
Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Example: Hospital cesarean delivery rates

L02-Visualizing  
data and  
describing data  
with numbers

We start by importing the data:

```
library(readxl)  
# this library helps with reading xlsx and xls files into R  
CS_dat <- read_xlsx("Kozhimannil_Ex_Cesarean.xlsx", sheet = 1)
```

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Example: Hospital cesarean delivery rates

```
head(CS_dat)
```

```
## # A tibble: 6 x 6
##   Births HOSP_BEDSIZE cesarean_rate lowrisk_cesarean_rate `Cesarean rate *` 
##   <dbl>      <dbl>        <dbl>            <dbl>
## 1     767          1       0.344           0.107
## 2     183          1       0.454           0.186
## 3     668          1       0.430           0.195
## 4     154          1       0.279           0.0844
## 5     327          1       0.306           0.119
## 6    2356          1       0.301           0.0662
## # ... with 1 more variable: 'Low Risk Cearean rate*100' <dbl>
```

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Example: Hospital cesarean delivery rates

L02-Visualizing  
data and  
describing data  
with numbers

```
names(CS_dat)
```

```
## [1] "Births"                      "HOSP_BEDSIZE"  
## [3] "cesarean_rate"                "lowrisk_cesarean_rate"  
## [5] "Cesarean rate *100"          "Low Risk Cearean rate*100"
```

let's take a moment to discuss variable names containing spaces

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

## Sidenote on variable names containing spaces

- ▶ Two variables in CS\_dat contain spaces.
- ▶ We generally want to remove spaces from variable names.
- ▶ Question: Which dplyr function can we use to change the variable names?
- ▶ Answer: rename(new\_name = old\_name) can be used. When the old variable name contains spaces, you need to place back ticks around it like this:

```
CS_dat <- CS_dat %>% rename(cs_rate = `Cesarean rate *100`,  
                                low_risk_cs_rate = `Low Risk Cesarean rate*100`)
```

- ▶ See this paper for tips on storing data in Excel for later analysis.

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Percentiles

# Tidy the data for analysis

L02-Visualizing  
data and  
describing data  
with numbers

For our example, we are only interested in each hospital's cesarean delivery rate, the rate for lower risk pregnancies, and the number of births at the hospital.

```
CS_dat <- CS_dat %>%  
  select(Births, cs_rate, low_risk_cs_rate) %>%  
  rename(num_births = Births)
```

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Analysis: Histogram of cesarean delivery rates across US hospitals

L02-Visualizing  
data and  
describing data  
with numbers

```
ggplot(CS_dat, aes(x = cs_rate)) +  
  geom_histogram(col = "white", binwidth = 5) +  
  labs( x = "Cesarean delivery rate (%)", y = "Count",  
        caption = "Data from: Kozhimannil, Law, and Virnig. Health Affairs. 2018;32(3)  
  geom_rug(alpha = 0.2, col = "forest green") + #alpha controls transparency  
  theme_minimal(base_size = 15)
```

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

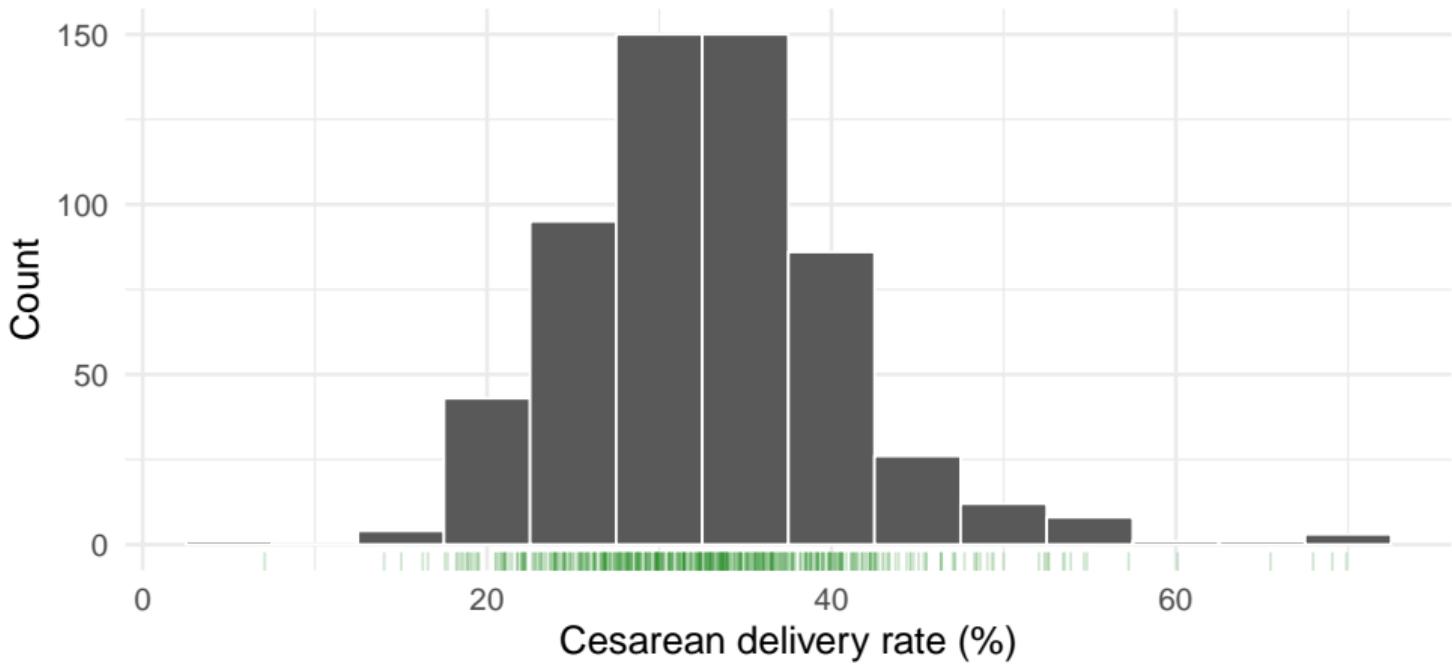
Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Histogram of cesarean delivery rates across US hospitals

L02-Visualizing  
data and  
describing data  
with numbers



Data from: Kozhimannil, Law, and Virnig. Health Affairs. 2013;32(3):527–35.

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Spread of cesarean delivery rates across US hospitals

L02-Visualizing  
data and  
describing data  
with numbers

- ▶ What can you say about this distribution? Would you expect so much variation across hospitals in their rates of cesarean delivery?
- ▶ Let's describe the **spread** of these data using the methods from Chapter 2.

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Quantiles

L02-Visualizing  
data and  
describing data  
with numbers

```
CS_dat %>% summarize(  
  Q1 = quantile(cs_rate, 0.25),  
  median = median(cs_rate),  
  Q3 = quantile(cs_rate, 0.75)  
)
```

```
## # A tibble: 1 x 3  
##       Q1   median     Q3  
##   <dbl>   <dbl>   <dbl>  
## 1  27.6    32.4   37.1
```

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere  
Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# dplyr's summarize() to calculate the five number summary

```
CS_dat %>% summarize(  
  min = min(cs_rate),  
  Q1 = quantile(cs_rate, 0.25),  
  median = median(cs_rate),  
  Q3 = quantile(cs_rate, 0.75),  
  max = max(cs_rate))  
  
## # A tibble: 1 x 5  
##       min     Q1 median     Q3     max  
##   <dbl> <dbl>  <dbl> <dbl> <dbl>  
## 1    7.09  27.6   32.4  37.1  69.9
```

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere  
Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

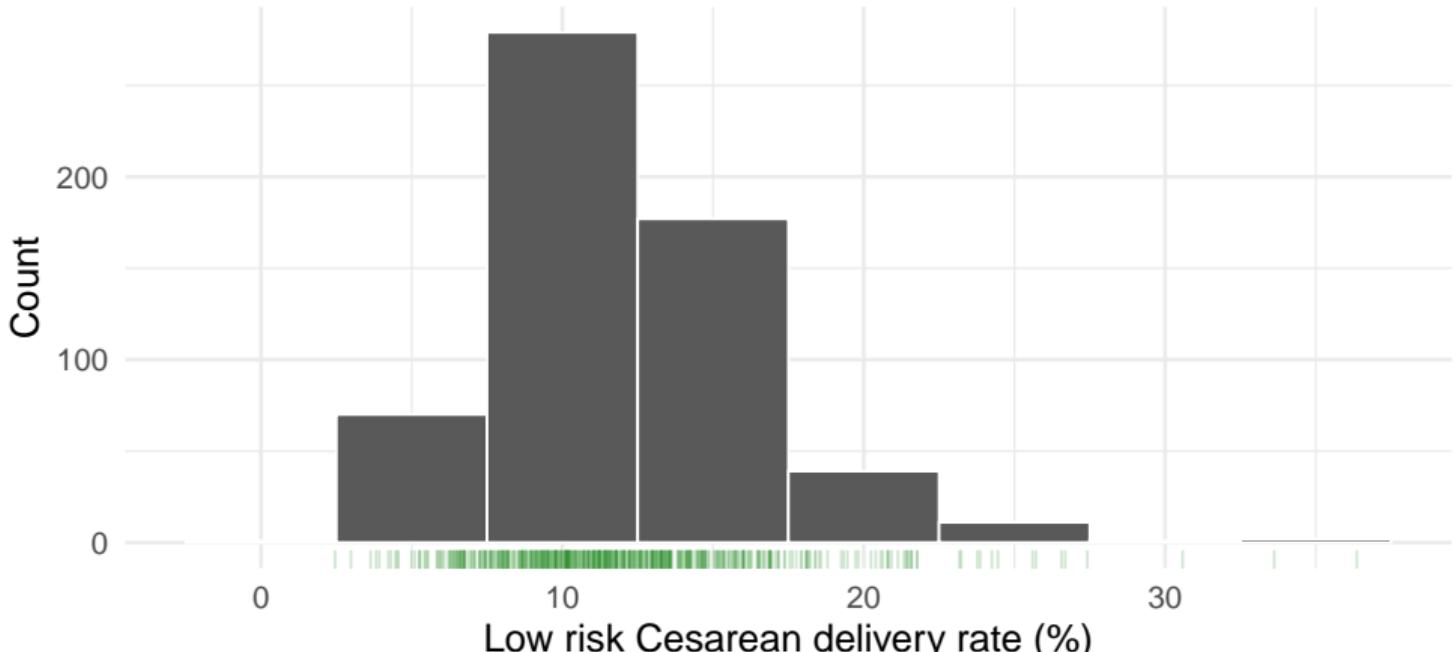
Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Histogram of low risk cesarean delivery rates across US hospitals

L02-Visualizing  
data and  
describing data  
with numbers



Data from: Kozhimannil, Law, and Virnig. Health Affairs. 2013;32(3):527–35.

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# dplyr's summarize() to calculate the five number summary

```
CS_dat %>% summarize(  
  min = min(low_risk_cs_rate),  
  Q1 = quantile(low_risk_cs_rate, 0.25),  
  median = median(low_risk_cs_rate),  
  Q3 = quantile(low_risk_cs_rate, 0.75),  
  max = max(low_risk_cs_rate))
```

```
## # A tibble: 1 x 5  
##       min     Q1 median     Q3     max  
##   <dbl> <dbl>  <dbl> <dbl> <dbl>  
## 1  2.46  9.19  11.4  14.2  36.4
```

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere  
Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Sample variance and standard deviation

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

**Sample variance and  
standard deviation**

Box plots

# Sample variance and standard deviation

Let  $s^2$  represent the variance of a sample. Then,

$$s^2 = \frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n - 1}$$

$$s^2 = \frac{1}{n - 1} ((x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2)$$

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

Let  $s$  represent the standard deviation of a sample. Then,

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

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Sample variance and standard deviation

L02-Visualizing  
data and  
describing data  
with numbers

- ▶ Some intuition on why we divide by  $n-1$ : link

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

**Sample variance and  
standard deviation**

Box plots

# dplyr's summarize() to calculate the standard deviation and the variance

```
CS_dat %>% summarize(  
  cs_sd = sd(cs_rate),  
  cs_var = var(cs_rate)  
)
```

```
## # A tibble: 1 x 2  
##   cs_sd cs_var  
##     <dbl>  <dbl>  
## 1    8.03    64.5
```

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere  
Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Example: Hospital cesarean delivery rates

L02-Visualizing  
data and  
describing data  
with numbers

What might we conclude from these data?

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

**Sample variance and  
standard deviation**

Box plots

## Example: Hospital cesarean delivery rates

From the article:

"we found that cesarean rates varied tenfold across hospitals, from 7.1 percent to 69.9 percent. Even for women with lower-risk pregnancies, in which more limited variation might be expected, cesarean rates varied fifteenfold, from 2.4 percent to 36.5 percent. Thus, vast differences in practice patterns are likely to be driving the costly overuse of cesarean delivery in many US hospitals."

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Box plots

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Box plots provide a nice visual summary of the center and spread

Also called box and whisker plots

The box:

- ▶ The centre line is the median
- ▶ The top of the box is the Q3
- ▶ The bottom of the box is the Q1

The whiskers - depends:

- ▶ The top of the top whisker is either the max value, or equal to the highest point that is below  $Q3 + 1.5 \times IQR$
- ▶ The bottom of the bottom whisker is either min value, or equal to the lowest point that is above  $Q1 - 1.5 \times IQR$
- ▶ In plots where the whiskers are **not** the min and max, the data points above and below the whiskers are the outliers

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# Box plots in R

L02-Visualizing  
data and  
describing data  
with numbers

```
ggplot(CS_dat, aes(y = cs_rate)) +  
  geom_boxplot() +  
  ylab("Cesarean delivery rate (%)") +  
  labs(title = "Box plot of the CS rates across US hospitals",  
       caption = "Data from: Kozhimannil et al. 2013.") +  
  theme_minimal(base_size = 15) +  
  scale_x_continuous(labels = NULL) # removes the labels from the x axis
```

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

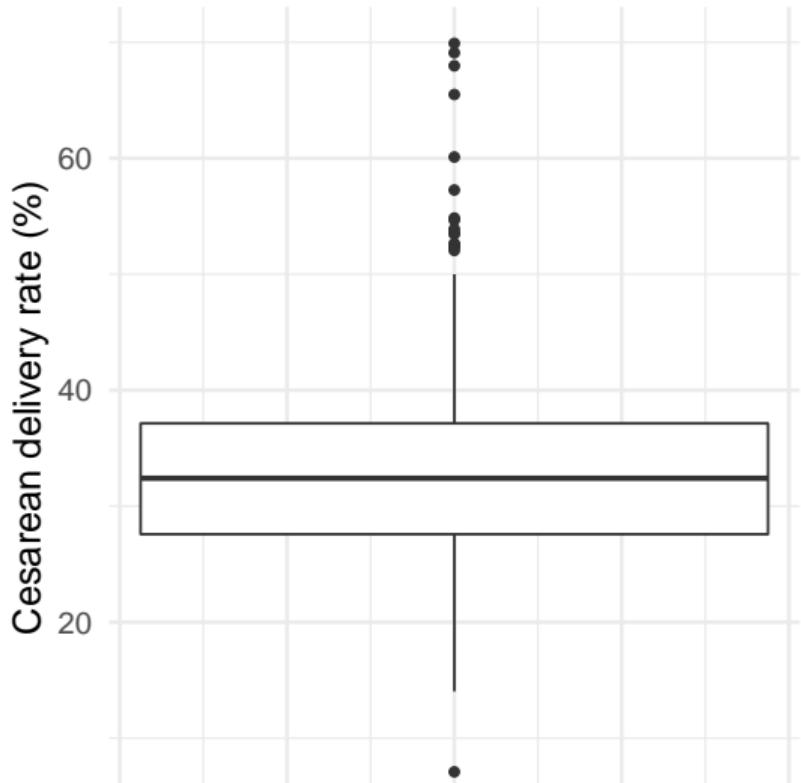
Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

Box plots provide a nice visual summary of the center and spread

### Box plot of the CS rates across US hospitals



Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# R Recap: Code we used today for graphs

L02-Visualizing  
data and  
describing data  
with numbers

1. 'ggplot' to set up a canvas for graphics
2. `geom_bar(stat = "identity")` to make a bar chart when you specify the y variable
3. `geom_histogram()` to make a histogram for which ggplot needs to calculate the count
4. `fct_reorder(var1, var2)` to reorder a categorical variable (`var1`) by a numeric variable (`var2`)
  - ▶ from the `forcats` package
5. `geom_boxplot` to make a boxplot

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere  
Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# R Recap: Code we used today to summarize data with numbers

L02-Visualizing  
data and  
describing data  
with numbers

1. `quantile(data, 0.25)`, `quantile(data, 0.75)` for Q1 and Q3, respectively
2. `min()` and `max()` for the full range of the data
3. `sd()` and `var()` for sample standard deviation and variance
4. Used the above within `summarize()` to easily output these measures

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# We only skimmed the surface!

L02-Visualizing  
data and  
describing data  
with numbers

- ▶ Here is some extra material for those of you who love data visualization.  
This material won't be tested.
  - ▶ RStudio ggplot2 cheatsheet
  - ▶ Kieran Healy's data visualization book

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

# References

L02-Visualizing  
data and  
describing data  
with numbers

1. Riddell CA, Morrison KT, Harper S, Kaufman JS. Trends in the contribution of major causes of death to the black-white life expectancy gap by US state. *Health & Place*. 2018. 52:85-100. doi: 10.1016/j.healthplace.2018.04.003.

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots

## REASONS WE'RE UP ALL NIGHT



Source: Daft Punk (research assistance by Pharrell Williams)

- ▶ from *Eric Tanoye Song Lyrics in Chart Form*

Visualizations for  
categorical data

Introducing ggplot

Visualizing quantitative  
variables

Describing your distribution  
- what are we looking for?

Measures of central  
tendency

Statistics is Everywhere

Discussion

Mean vs Median: Outliers  
and sample size, skew,  
shape

Measures of spread

Example: Hospital cesarean  
delivery rates

Sample variance and  
standard deviation

Box plots