# Lecture 3: Data Description

STA3100: Programming with Data

# **Variability**

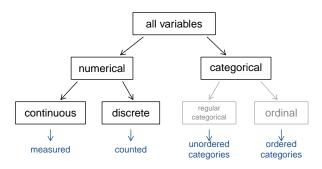
Variability is a very important concept in statistics. The higher the variability, the harder it is to draw conclusions.

Common sources of variability:

- Variability across measurements measurement error can cause the recorded length to vary from measurement to measurement.
- Variability across subjects average pulse rate varies from moment to moment and from person to person
- Variability across groups in a comparison between men and women, need to take into account variability within men and women.

Chapter 3 Data Description

# Types of variables



STA3100: Programming with Data

Lecture 3

Chapter 3 Data Description

Single Variable: Numerical descriptive measures

The most common numerical descriptive measures are:

- Central tendency measures: mean, median, mode
- Variability measures: variance, standard deviation, range, IQR

population 
$$\longrightarrow$$
 parameters sample  $\longrightarrow$  statistics

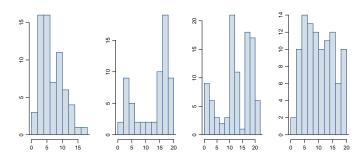
	mean	variance	SD
sample	$\bar{x}$	$s^2$	S
population	$\mu$	$\sigma^2$	$\sigma$

STA3100: Programming with Data STA3100: Programming with Data Lecture 3 Lecture 3

Chapter 3 Data Description

# Shape of a distribution: modality

The *mode* is defined as the most frequent observation in the data set. Does the histogram have a single prominent peak (unimodal), several prominent peaks (bimodal/multimodal), or no apparent peaks (uniform)?



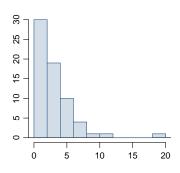
In order to determine modality, it's best to step back and imagine a smooth curve over the histogram.

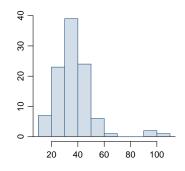
STA3100: Programming with Data

Lecture 3

# Shape of a distribution: unusual observations

Are there any unusual observations or potential *outliers*?

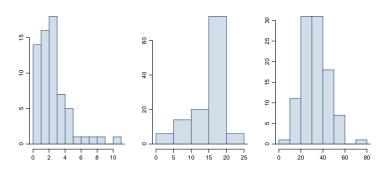




Chapter 3 Data Description

# Shape of a distribution: skewness

Is the histogram right skewed, left skewed, or symmetric?



Histograms are said to be skewed to the side of the long tail.

STA3100: Programming with Data

## How would you expect each of the following to be distributed?

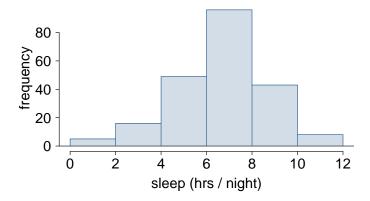
- Weights of adult males
- Salaries of people from Florida
- Exam scores on an easy test
- Birthdays of classmates (day of the month)

Lecture 3 Lecture 3 STA3100: Programming with Data STA3100: Programming with Data

Chapter 3 Variability

## Variability in data

How would you describe the amount of variability in the number of hours of sleep students get per night?



STA3100: Programming with Data

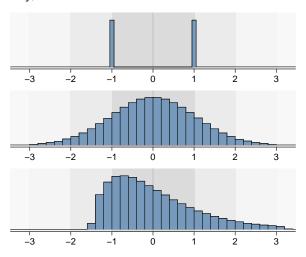
ecture 3

9 / 34

#### Chapter 3 Varial

# Describing distributions

When describing distributions make sure to talk about the shape, center, spread, and if any, unusual observations.



Chapter 3 Variabilit

# Sample Variance and Standard Deviation

STA3100: Programming with Data

[1] 2.02

Lecture 3

10 / 24

# Q1, Q3, and IQR

- The 25<sup>th</sup> percentile is also called the first quartile, *Q1*.
- The 50<sup>th</sup> percentile is also called the median.
- The 75<sup>th</sup> percentile is also called the third quartile, *Q3*.

```
summary(d$study_hours)
  Min. 1st Qu. Median Mean 3rd Qu. Max. NA's
  2.00 10.00 15.00 17.45 20.00 69.00 20.00
quantile(d$study_hours, prob=0.25, na.rm=TRUE)
  25\%
  10
```

• Between Q1 and Q3 is the middle 50% of the data. The range these data span is called the *interquartile range*, or the *IQR*.

$$IQR = 20 - 10 = 10$$

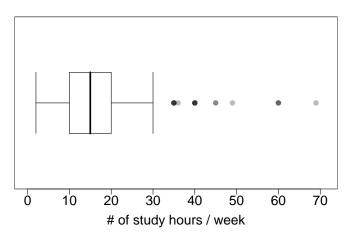
STA3100: Programming with Data Lecture 3 11/34 STA3100: Programming with Data

Lecture 3

12/34

## Box plot

The box in a box plot represents the middle 50% of the data, and the thick line in the box is the median.



STA3100: Programming with Data

Lecture 3

## Whiskers and outliers

STA3100: Programming with Data

• Whiskers of a box plot can extend up to 1.5 \* IQR away from the quartiles.

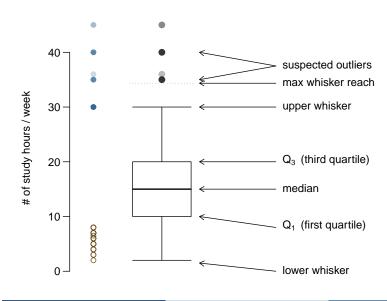
max upper whisker reach : 
$$Q3 + 1.5 * IQR = 20 + 1.5 * 10 = 35$$
 max lower whisker reach :  $Q1 - 1.5 * IQR = 10 - 1.5 * 10 = -5$ 

• The whiskers are calculated as the min/max value between the data points within the max whisker reach and the max whisker reach limit:

upper whisker : 
$$min(max(x), Q3 + 1.5 * IQR) = 30$$
  
lower whisker :  $max(min(x), Q1 - 1.5 * IQR) = 2$ 

• An outlier is defined as an observation beyond the maximum reach of the whiskers. It is an observation that appears extreme relative to the rest of the data.

# Anatomy of a box plot



STA3100: Programming with Data

Outliers (cont.)

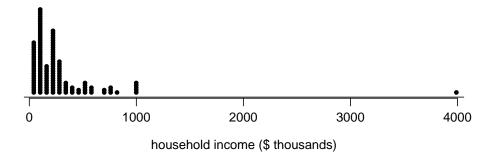
Why is it important to look for outliers?

Lecture 3 Lecture 3 STA3100: Programming with Data

Chapter 3 Robust statistics

# Extreme observations

How would sample statistics such as mean, median, SD, and IQR of household income be affected if the largest value was replaced with \$10 million? What if the smallest value was replaced with \$10 million?



STA3100: Programming with Data

Lecture 3

17 / 34

Chapter 3

Robust statistics

## Robust statistics

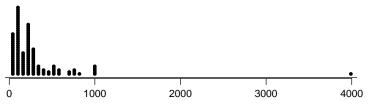
Median and IQR are more robust to skewness and outliers than mean and SD. Therefore.

- for skewed distributions it is more appropriate to use median and IQR to describe the center and spread
- for symmetric distributions it is more appropriate to use the mean and SD to describe the center and spread

If you would like to estimate the typical household income for a UF student, would you be more interested in the mean or median income?

### Chapter 3 Robust statistics

## Robust statistics



household income (\$ thousands)

	robust		not r	not robust	
scenario	median	IQR	$\bar{x}$	S	
original data	165K	150K	211K	180K	
move largest to \$10 million	165K	150K	398K	1,422K	
move smallest to \$10 million	190K	163K	4,186K	1,424K	

STA3100: Programming with Data

Lecture 3

18 / 34

Chapter 3 Correlation

# More than one variable: Calculating the correlation

Pearson correlation coefficient:

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

- *Correlation* describes the strength and direction of the *linear* relationship between two variables.
- It takes values between -1 (perfect negative relationship) and +1 (perfect positive relationship).
- A value of 0 indicates no linear relationship.
- Using R:

cor(x, y)

STA3100: Programming with Data Lecture 3 19 / 34 STA3100: Programming with Data Lecture 3

Describing the relationship: Shape, Direction, and Strength

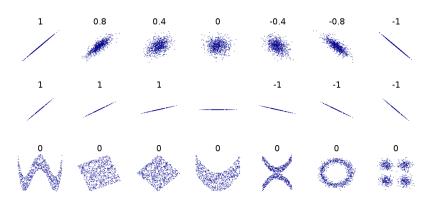
The scatterplot below shows the relationship between HS graduate rate in

all 50 states and the District of Columbia vs the % of residents who live

below the poverty line (< \$22,350 for a family of 4).

# Quantifying the relationship

## Correlation means linear association!



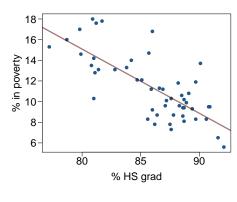
STA3100: Programming with Data

STA3100: Programming with Data

% HS grad

# Guessing the correlation

Which of the following is the best guess for the correlation between % in poverty and % HS grad?



-0.75

-0.1

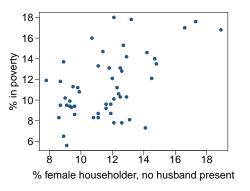
0.02

-1.5

How would you describe the relationship between % in poverty and % HS grad?

# Guessing the correlation

Which of the following is the best guess for the correlation between % in poverty and % HS grad?



0.1

-0.4

0.9

0.5

STA3100: Programming with Data

Lecture 3

# Graphical description

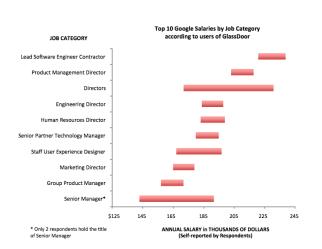
Depending on the type of measurement, common plots are pie charts, bar charts, histograms, box plots, and density plots.

- Pie charts: Used to describe any variable type.
  - Continuous numeric variables must be collapsed into bins or buckets.
  - The size of the sectors of the pie represent the relative frequency of each category.
- Bar charts: Nominal or ordinal data.
  - The variable levels are arrayed on the bottom (or left side) of the plot.
  - Bars above (or beside) the levels represent the range of values taken by the response variable or the frequency or relative frequency of the number of observations belonging to the various categories.

STA3100: Programming with Data

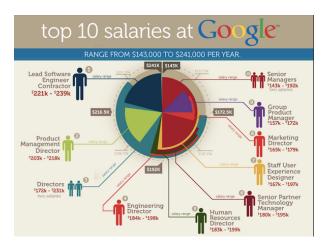
Lecture 3

## Problem with Pie Charts



http://junkcharts.typepad.com/junk\_charts/2011/10/the-massive-burden-of-pie-charts.html

## Problem with Pie Charts



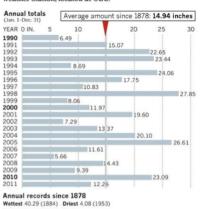
http://junkcharts.typepad.com/junk\_charts/2011/10/the-massive-burden-of-pie-charts.html

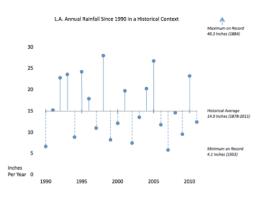
STA3100: Programming with Data

# Axis Placement - What story to tell

### L.A. annual rainfall, 1990-2011

A total of 12.26 inches of rain was recorded in 2011 at the National Weather Service's downtown Los Angeles weather station, located at USC.





STA3100: Programming with Data Lecture 3 STA3100: Programming with Data Lecture 3 Chapter 3 G

Graphical Methods

## **Univariate Plots**

# Graphical description (cont.): Numeric Variables

- *Histograms:* The heights of the bars above the bins represent the frequency or relative frequency of the various bins.
- Box plots: They identify particular percentiles of a distribution and are useful in detecting outlying observations and spread in the distribution.
- Density plots: They represent smoothed histograms describing the proportion of measurements within some distance of each point on the continuum.
- *Violin plots:* Similar to box-plot but allows to see better the shape of the distribution.

Categorical

- "Good" Bar plot
- "Bad" Pie chart
- Numerical
  - "Good" boxplot (outliers), histogram/density (shape), violin plot (shape and outliers)

Lecture 3

Chapter 3

Chapter 3 Graphical Methods

• "Bad" - Stem and leaf plot

STA3100: Programming with Data

Lecture 3

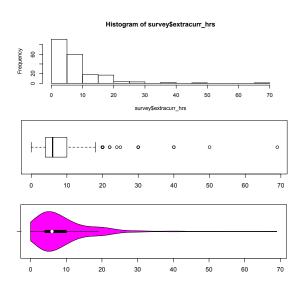
29 / 34

30 / 34

Chapter 3

Graphical Methods

## Univariate Plots - Numerical



Graphical Method

## Bivariate Plots

STA3100: Programming with Data

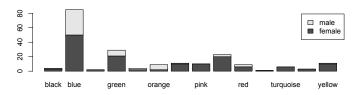
- Categorical vs Categorical
  - Stacked or grouped bar plot, mosaic plots
- Numerical vs Categorical
  - $\bullet\,$  Side by side boxplot or violin plot
- Numerical vs Numerical
  - Scatter plot, line plot

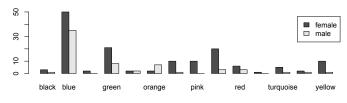
STA3100: Programming with Data Lecture 3 31/34 STA3100: Programming with Data Lecture 3

Chapter 3 Graphical Methods

# Bivariate Plots - Bar plots

## Favorite color by gender:



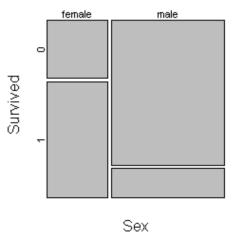


STA3100: Programming with Data Lecture 3 33 / 34

Chapter 3 Graphical Methods

# Bivariate Plots - Mosaic plot

## Titanic survivors by gender:



STA3100: Programming with Data

Lecture 3