

## Unit 1: Introduction to data

### 2. Exploratory data analysis

Sta 104 - Summer 2018, Term 1

Duke University, Department of Statistical Science

Prof. White

Slides posted at <https://www2.stat.duke.edu/courses/Summer18/sta104.001-1/>

- ▶ **Individual:** 15 minutes, Sakai quiz: Webcam turned on for proctoring. No communication with others, but you can ask me questions if a problem is unclear by sending a private message in Webex (75%)
- ▶ **Team:** 10 minutes, Sakai quiz: Discuss within team to decide the best answer (25%)

1

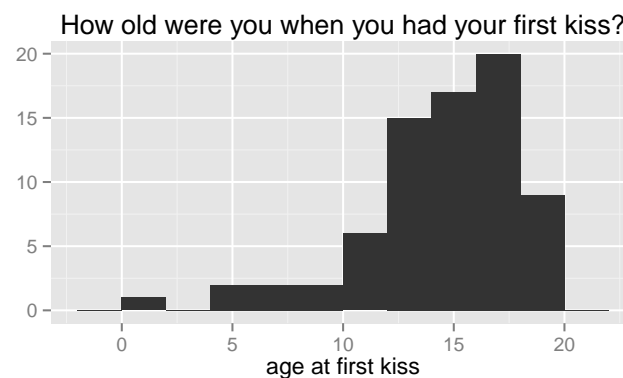
## Announcements

- ▶ PS 1 is assigned on the course website, start working on it

2

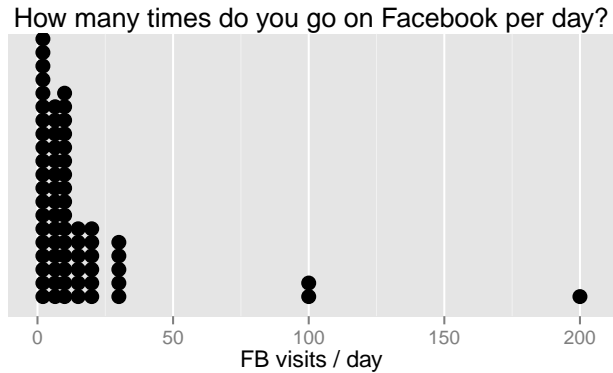
## From a past Sta 101 survey...

Do you see anything out of the ordinary?



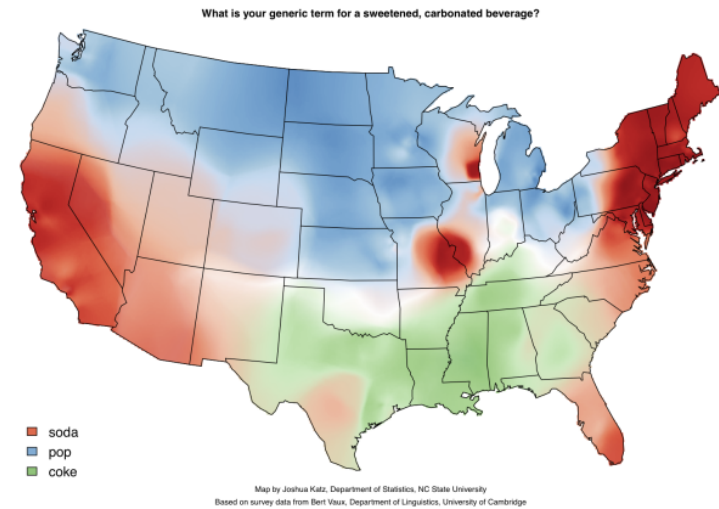
3

How are people reporting lower vs. higher values of FB visits?



4

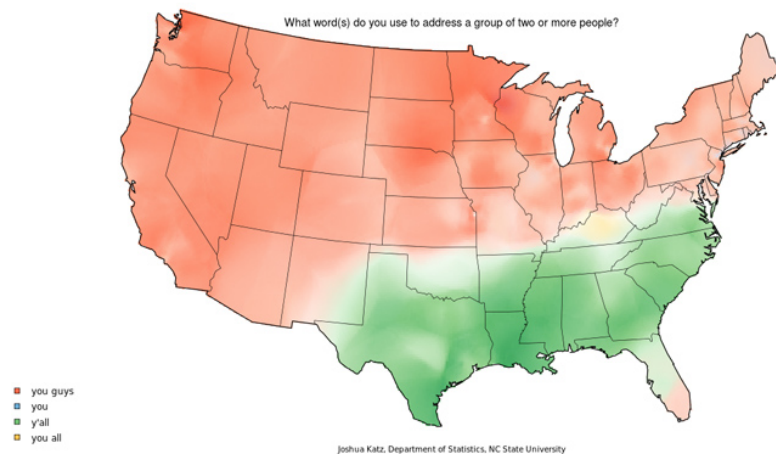
Describe the spatial distribution of preferred sweetened carbonated beverage drink.



<http://spark.rstudio.com/jkatz/SurveyMaps>

5

What is missing in this visualization?



<http://spark.rstudio.com/jkatz/SurveyMaps>

6

Describing distributions of numerical variables

- ▶ **Shape:** skewness, modality
- ▶ **Center:** an estimate of a *typical* observation in the distribution (mean, median, mode, etc.)
  - Notation:  $\mu$ : population mean,  $\bar{x}$ : sample mean
- ▶ **Spread:** measure of variability in the distribution (standard deviation, IQR, range, etc.)
- ▶ **Unusual observations:** observations that stand out from the rest of the data that may be suspected outliers

7

## Clicker question

Which of these is most likely to have a roughly symmetric distribution?

- (a) salaries of a random sample of people from North Carolina
- (b) weights of adult females
- (c) scores on an well-designed exam
- (d) last digits of phone numbers

8

## Clicker question

How do the mean and median of the following two datasets compare?

Dataset 1: 30, 50, 70, 90

Dataset 2: 30, 50, 70, 1000

- (a)  $\bar{X}_1 = \bar{X}_2$ ,  $median_1 = median_2$
- (b)  $\bar{X}_1 < \bar{X}_2$ ,  $median_1 = median_2$
- (c)  $\bar{X}_1 < \bar{X}_2$ ,  $median_1 < median_2$
- (d)  $\bar{X}_1 > \bar{X}_2$ ,  $median_1 < median_2$
- (e)  $\bar{X}_1 > \bar{X}_2$ ,  $median_1 = median_2$

9

## Standard deviation and variance

- ▶ Most commonly used measure of variability is the *standard deviation*, which roughly measures the average deviation from the mean
  - Notation:  $\sigma$ : population standard deviation,  $s$ : sample standard deviation
- ▶ Calculating the standard deviation, for a population (rarely, if ever) and for a sample:

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

- ▶ Square of the standard deviation is called the *variance*.

10

## More on SD

Why divide by  $n - 1$  instead of  $n$  when calculating the sample standard deviation?

Lose a “degree of freedom” for using an estimate (the sample mean,  $\bar{x}$ ), in estimating the sample variance/standard deviation.

Why do we use the squared deviation in the calculation of variance?

- ▶ To get rid of negatives so that observations equally distant from the mean are weighed equally.
- ▶ To weigh larger deviations more heavily.

11

## Clicker question

True / False: The range is always at least as large as the IQR for a given dataset.

- ☒ (a) Yes
- ☐ (b) No

Is the range or the IQR more robust to outliers?

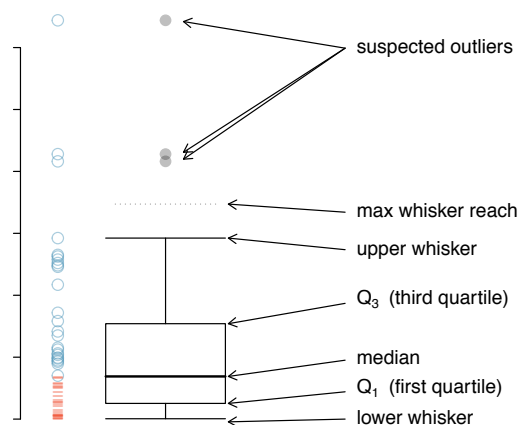
12

- ▶ Mean and standard deviation are easily affected by extreme observations since the value of each data point contributes to their calculation.
- ▶ Median and IQR are more robust.
- ▶ Therefore we choose median&IQR (over mean&SD) when describing skewed distributions.

13

## Box plot

A *box plot* visualizes the median, the quartiles, and suspected outliers. An *outlier* is defined as an observation more than  $1.5 \times \text{IQR}$  away from the quartiles.



14

## Application exercise: 1.1 Distributions of numerical variables

See the course website for instructions.

15

1. Always start your exploration with a visualization
2. When describing numerical distributions discuss shape, center, spread, and unusual observations
3. Robust statistics are not easily affected by outliers and extreme skew
4. Use box plots to display quartiles, median, and outliers