

Unit 1: Introduction to data

2. Exploratory data analysis

Sta 101 - Summer 2017

Duke University, Department of Statistical Science

Prof. Burris Slides posted at <http://www2.stat.duke.edu/courses/Summer17/sta101.001-2/>

Team Pascal: Lilly Gu, Kortney Dry, Trevon Lee, and Kaylee Brillhart

Team Gauss: Paul Hill, Catalin Mateas, Domonique Panton, and Edom Tilahun

Team Bayes: Mark Birmingham, Rachel Rubin, Samuel Johnson, and Casey Martinez

Team Fisher: Kyra Lambert, Aaron Therien, Anne Slusser, and Maliik Marcin

Team Pearson: Xiaofu Zhang, Michael Bivona, Sam June, and Jessica Findlay

Team Tukey: Allison Florian, Elizabeth Bonnell, Nicholas Solfanelli, and Sayvon Sampson

1

Readiness assessment

- *Individual:* 10 minutes, using bubble sheets

RA 1

Name: _____

1. ☐ A ☐ B ☐ C ☐ D
2. ☐ A ☐ B ☐ C ☐ D
3. ☐ A ☐ B ☐ C ☐ D
4. ☐ A ☐ B ☐ C ☐ D
5. ☐ A ☐ B ☐ C ☐ D
6. ☐ A ☐ B ☐ C ☐ D
7. ☐ A ☐ B ☐ C ☐ D
8. ☐ A ☐ B ☐ C ☐ D
9. ☐ A ☐ B ☐ C ☐ D
10. ☐ A ☐ B ☐ C ☐ D

- *Team:* 7 minutes, using bubble sheets (1 per team)

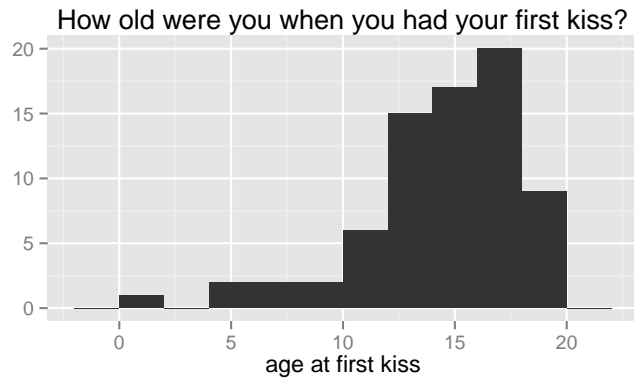
2

Announcements

- PS 1 is assigned on the course website. I recommend you start working on it

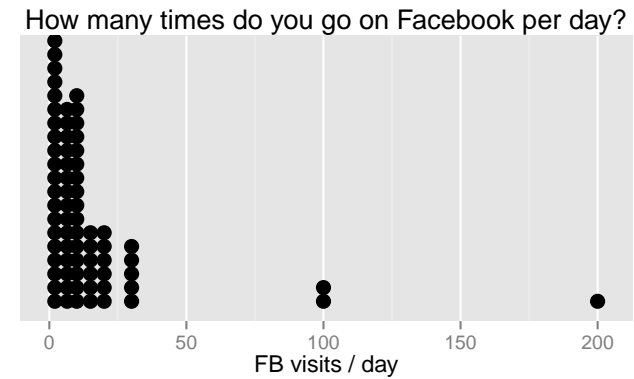
3

Do you see anything out of the ordinary?



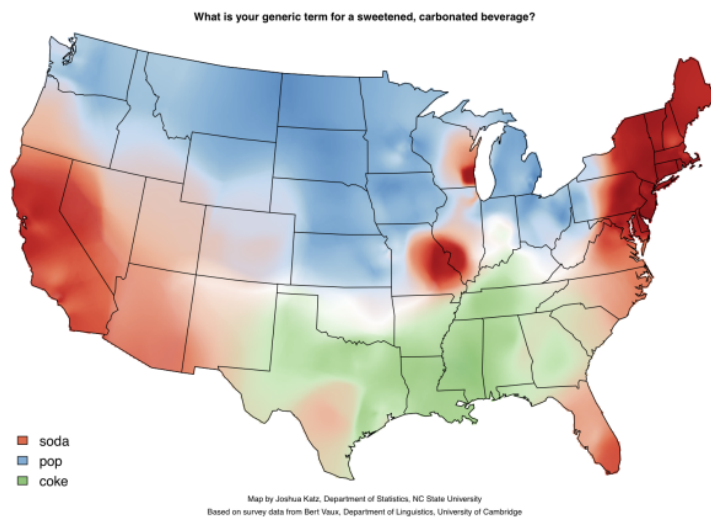
4

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



5

Describe the spatial distribution of preferred sweetened carbonated beverage drink.



6

What is missing in this visualization?



7

- ▶ *Shape*: skewness, modality
- ▶ *Center*: an estimate of a *typical* observation in the distribution (mean, median, mode, etc.)
 - Notation: μ : 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

8

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

9

Mean vs. median

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) $\bar{x}_1 > \bar{x}_2$, $median_1 = median_2$

10

Standard deviation and variance

- ▶ Most commonly used measure of variability is the *standard deviation*, which roughly measures the average deviation from the mean
 - Notation: σ : 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*.

11

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.

12

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?

13

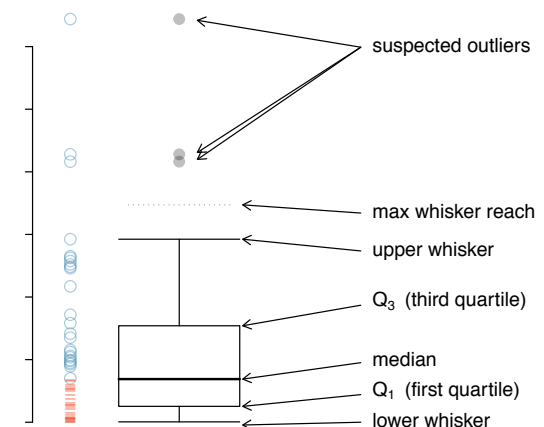
Robust statistics

- ▶ 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.

14

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.



15

Application exercise: 1.2 Distributions of numerical variables

See the course website for instructions.

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