

## Unit 3: Foundations for inference

### 4. MT 1 Review

Sta 101 - Spring 2019

Duke University, Department of Statistical Science

Dr. Abrahamsen

Slides posted at <https://stat.duke.edu/courses/Spring19/sta101.002>

### Exam Details:

- ▶ When: Thursday, Feb 14 in Class
- ▶ What to bring:
  - Scientific calculator (graphing calculator ok, No Phones!)
  - Cheat sheet (can be typed - Definitions and equations only!)
- ▶ Provided: Z table

### Exam Format:

- ▶ 4 written questions - 60 pts
- ▶ 5 T/F - 2pts each
- ▶ 10 Multiple choice - 3pts each

## Key Terms

- ▶ Population
- ▶ Parameter
- ▶ Statistic
- ▶ Simple Random Sample
- ▶ Stratified Sample
- ▶ Cluster Sample
- ▶ Multistage Sample
- ▶ Experiment
- ▶ Observational Study
- ▶ Control
- ▶ Placebo
- ▶ Confounding Variable

## Data Collection, Observational Studies & Experiments

<i>ideal experiment</i>	Random assignment	No random assignment	<i>most observational studies</i>
Random sampling	Causal conclusion, generalized to the whole population.	No causal conclusion, correlation statement generalized to the whole population.	Generalizability
No random sampling	Causal conclusion, only for the sample.	No causal conclusion, correlation statement only for the sample.	No generalizability
<i>most experiments</i>	Causation	Correlation	<i>bad observational studies</i>

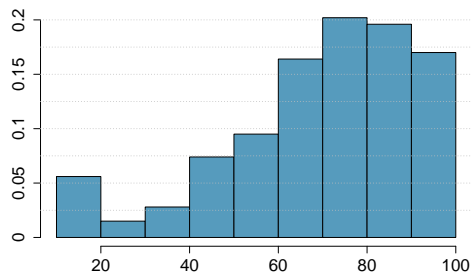
### Describing Distributions of Numerical Variables:

- ▶ *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

4

#### Clicker question

Which of the following is false?



- (a) The box plot would have outliers only on the lower end.
- (b) The median is between 70 and 80.
- (c) More than 25% of the data is above 90.
- (d) More than 50% of the data have positive Z scores.
- (e) The mean is likely to be smaller than the median.

6

### Robust statistics:

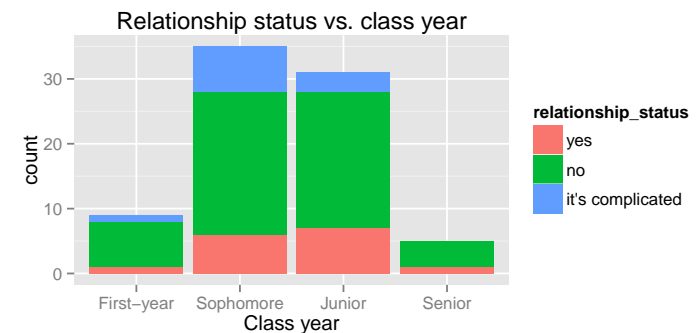
- ▶ Mean and standard deviation are easily affected by extreme observations since the value of each data point contributes to their calculation.
- ▶ Median and IQR are more robust.
- ▶ Therefore we choose median & IQR (over mean & SD) when describing skewed distributions.

5

### More Exploratory Data Analysis

#### Use segmented bar plots for visualizing relationships between 2 categorical variables

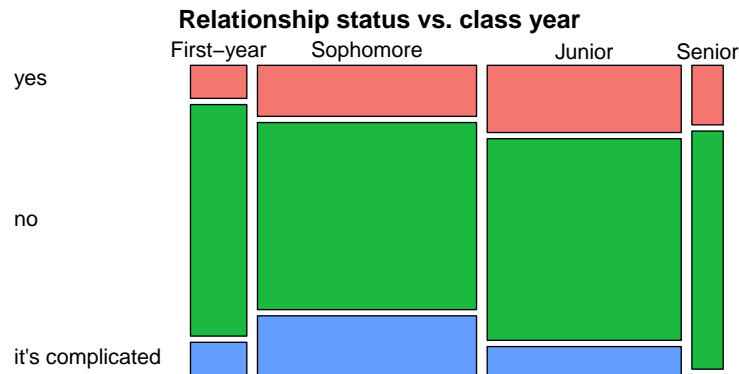
What do the heights of the segments represent? Is there a relationship between class year and relationship status? What descriptive statistics can we use to summarize these data? Do the widths of the bars represent anything?



7

**...or use a mosaic plot**

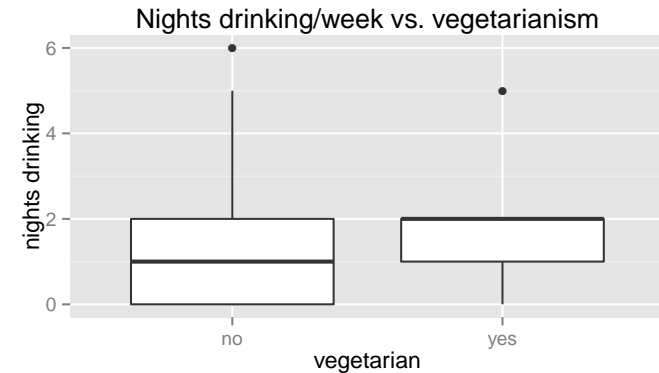
What do the widths of the bars represent? What about the heights of the boxes? Is there a relationship between class year and relationship status? What other tools could we use to summarize these data?



8

**Use side-by-side box plots to visualize relationships between a numerical and categorical variable**

How do drinking habits of vegetarian vs. non-vegetarian students compare?



9

**Key Ideas:**

- ▶ Observed differences may be due to random chance
- ▶ Test whether difference is significant using simulations

10

CPR is a procedure commonly used on individuals suffering a heart attack when other emergency resources are not available. The chest compressions involved with this procedure can also cause internal injuries. Blood thinners that are often given to help release a clot that is causing the heart attack may also negatively affect such internal injuries. An experiment was designed to evaluate if blood thinners have an impact on survival after a heart attack. Patients were randomly divided into a treatment group (received a blood thinner) or the control group (no blood thinner). The outcome variable of interest was whether the patients survived for at least 24 hours.

11

Form hypotheses for this study in plain and statistical language. Let  $p_c$  represent the true survival proportion in the control group and  $p_t$  represent the survival proportion for the treatment group.

$H_0$ : Blood thinners do not have an overall survival effect, i.e. the survival proportions are the same in each group.  $p_t - p_c = 0$ .

$H_A$ : Blood thinners do have an impact on survival.  $p_t - p_c \neq 0$ .

#### Clicker question

Given these hypotheses, what is the sample statistic?

$$H_0 : p_t - p_c = 0 \quad H_A : p_t - p_c \neq 0$$

	Survived	Died	Total
Control	11	39	50
Treatment	14	26	40
Total	25	65	90

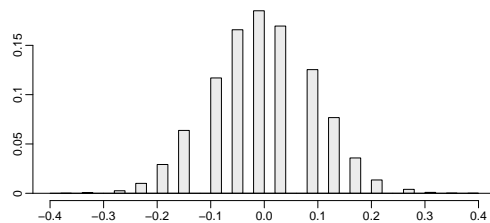
- (a)  $(11 / 25) - (39 / 65) = -0.16$
- (b)  $(14 / 40) - (11 / 50) = 0.13$
- (c)  $(14 / 90) - (11 / 90) = 0.033$
- (d)  $(40 / 90) - (50 / 90) = -0.111$

12

13

#### Clicker question

A randomization test was conducted to evaluate these hypotheses. Based on the randomization distribution below, what is the conclusion?



These data

- (a) provide convincing evidence that blood thinners
- (b) provide convincing evidence that blood thinners do not
- (c) do not provide convincing evidence that blood thinners
- (d) do not provide convincing evidence that blood thinners do not have an impact on survival.

14

## Probability and Conditional Probability

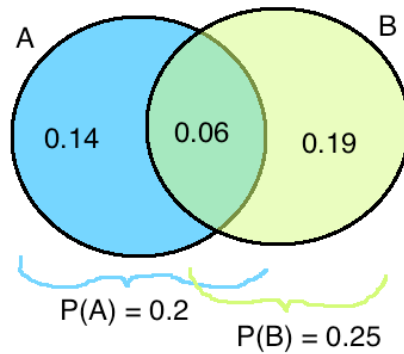
- ▶ *Disjoint (mutually exclusive) events* cannot happen at the same time
  - For disjoint A and B:  $P(A \text{ and } B) = 0$
- ▶ If A and B are *independent events*, having information on A does not tell us anything about B (and vice versa)
  - If A and B are independent:
    - $P(A | B) = P(A)$
    - $P(A \text{ and } B) = P(A) \times P(B)$
- ▶ *General addition rule*:  $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$
- ▶ *Bayes' theorem*:  $P(A | B) = \frac{P(A \text{ and } B)}{P(B)}$

15

Clicker question

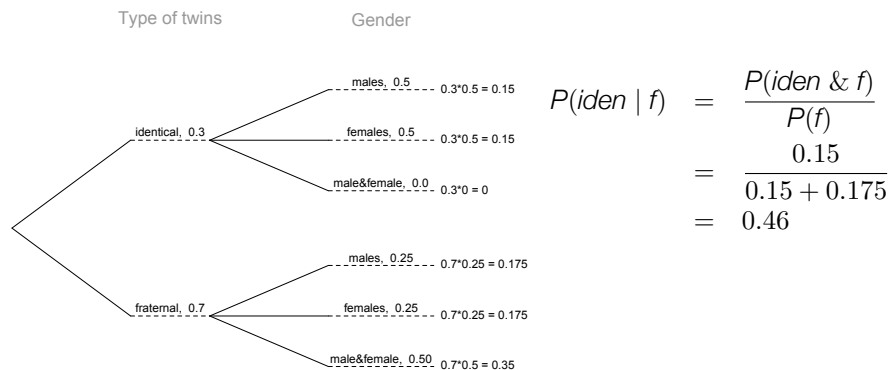
Which of the following is true?

- (a) A and B are independent.
- (b)  $P(A \text{ but not } B) = 0.2$
- (c)  $P(A | B) = 0.06 / 0.14$
- (d)  $P(A \text{ or } B) = 0.14 + 0.19 + 0.06$
- (e)  $P(\text{neither } A \text{ nor } B) = 1 - 0.06$



16

About 30% of human twins are identical and the rest are fraternal. Identical twins are necessarily the same sex – half are males and the other half are females. One-quarter of fraternal twins are both male, one-quarter both female, and one-half are mixes: one male, one female. You have just become a parent of twins and are told they are both girls. Given this information, what is the posterior probability that they are identical?



$$\begin{aligned}
 P(\text{iden} | f) &= \frac{P(\text{iden} \& f)}{P(f)} \\
 &= \frac{0.15}{0.15 + 0.175} \\
 &= 0.46
 \end{aligned}$$

18

- ▶ Probability trees are useful for organizing information in conditional probability calculations
- ▶ They're especially useful in cases where you know  $P(A | B)$ , along with some other information, and you're asked for  $P(B | A)$
- ▶ Using Bayes' theorem

$$\begin{aligned}
 P(\text{hypothesis} | \text{data}) &= \frac{P(\text{hypothesis and data})}{P(\text{data})} \\
 &= \frac{P(\text{data} | \text{hypothesis}) \times P(\text{hypothesis})}{P(\text{data})}
 \end{aligned}$$

17

Normal and Binomial Distributions

- ▶ Two types of probability distributions: discrete and continuous
- ▶ Normal distribution is unimodal, symmetric and follows the 68-95-99.7 rule
- ▶ Z scores serve as a ruler for any distribution

$$Z = \frac{\text{obs} - \text{mean}}{SD}$$

- ▶ Z score: number of standard deviations the observation falls above or below the mean

19

- ▶ The *Binomial distribution* describes the probability of having exactly  $k$  successes in  $n$  independent trials with probability of success  $p$ .

$$P(k \text{ successes in } n \text{ trials}) = \binom{n}{k} p^k (1 - p)^{(n-k)}$$

- ▶ *Expected Value:*  $np$
- ▶ *Standard Deviation:*  $\sqrt{np(1 - p)}$
- ▶ Shape of the binomial distribution approaches normal when the S-F rule is met

20

Clicker question

Which of the following probabilities should be calculated using the Binomial distribution?

Probability that

- (a) a basketball player misses 3 times in 5 shots
- (b) train arrives on the time on the third day for the first time
- (c) height of a randomly chosen 5 year old is greater than 4 feet
- (d) a randomly chosen individual likes chocolate ice cream best

21

Why Binomial?

Suppose the probability of a miss for this basketball player is 0.40. What is the probability that she misses 3 times in 5 shots?

- ▶ One possible scenario is that she misses the first three shots, and makes the last two. The probability of this scenario is:

$$0.4^3 \times 0.6^2 \approx 0.023$$

- ▶ But this isn't the only possible scenario:

- |          |           |          |          |           |
|----------|-----------|----------|----------|-----------|
| 1. MMMHH | 3. MHMMH  | 5. HMMHM | 7. HHMMM | 9. MHHMM  |
| 2. MMHMH | 4. HMMMHH | 6. HMHMM | 8. MHHMM | 10. MMHMM |

- ▶ Each one of these scenarios has 3 Ms and 2 Hs, therefore the probability of each scenario is 0.023.
- ▶ Then, the total probability is  $10 \times 0.023 = 0.23$ .

22

... concisely

Suppose the probability of a miss for this basketball player is 0.40. What is the probability that she misses 3 times in 5 shots?

$$\begin{aligned} \binom{5}{3} \times 0.4^3 \times 0.6^2 &= \frac{5!}{3! \times 2!} \times 0.4^3 \times 0.6^2 \\ &= 10 \times 0.023 \\ &= 0.23 \end{aligned}$$

23

### Clicker question

Which of the following highlights the correct outcomes for “at most 3 misses in 5 shots”?

- (a) {0, 1, 2, 3, 4, 5}
- (b) {0, 1, 2, 3, 4, 5}
- (c) {0, 1, 2, 3, 4, 5}
- (d) {0, 1, 2, 3, 4, 5}
- (e) {0, 1, 2, 3, 4, 5}

24

## Variability in Estimates and CLT

- ▶ Sample Statistics vary from sample to sample
- ▶ CLT describes the shape, center and spread of sampling distributions

$$\bar{x} \sim N \left( \text{mean} = \mu, SE = \frac{\sigma}{\sqrt{n}} \right)$$

- ▶ CLT only applies when independence and sample size/skew conditions are met

26

### Clicker question

Which of the following is the correct calculation for “P(at most 3 misses in 5 shots)”?

Note: P(k) means P(k misses in 5 shots), calculated using the binomial formula.

- (a) P(0) + P(1) + P(2)
- (b) P(3) + P(4) + P(5)
- (c) 1 - P(0)
- (d) 1 - [ P(0) + P(1) + P(2) ]
- (e) 1 - [ P(4) + P(5) ]

25

A housing survey was conducted to determine the price of a typical home in Topanga, CA. The mean price of a house was roughly \$1.3 million with a standard deviation of \$300,000. There were no houses listed below \$600,000 but a few houses above \$3 million.

### Clicker question

Can we estimate the probability that the mean of 60 randomly chosen houses in Topanga is more than \$1.4 million?

- (a) yes
- (b) no

27

A housing survey was conducted to determine the price of a typical home in Topanga, CA. The mean price of a house was roughly \$1.3 million with a standard deviation of \$300,000. There were no houses listed below \$600,000 but a few houses above \$3 million.

What is the probability that the mean of 60 randomly chosen houses in Topanga is more than \$1.4 million?

In order to calculate  $P(\bar{X} > 1.4 \text{ mil})$ , we need to first determine the distribution of  $\bar{X}$ . According to the CLT,

$$\bar{X} \sim N\left(\text{mean} = 1.3, SE = \frac{0.3}{\sqrt{60}} = 0.0387\right)$$

$$\begin{aligned} P(\bar{X} > 1.4) &= P\left(Z > \frac{1.4 - 1.3}{0.0387}\right) \\ &= P(Z > 2.58) \\ &= 1 - 0.9951 = 0.0049 \end{aligned}$$

28

- ▶ Statistical inference methods based on the CLT require the same conditions as the CLT
- ▶ *CI: point estimate  $\pm$  margin of error*
- ▶ Calculate the sample size a priori to achieve desired margin or error

Solve for  $n$ :

$$ME = z^* \frac{s}{\sqrt{n}}$$

29