

Chapter 2 - Summarizing Data

Ramnivas Singh

02/18/2021

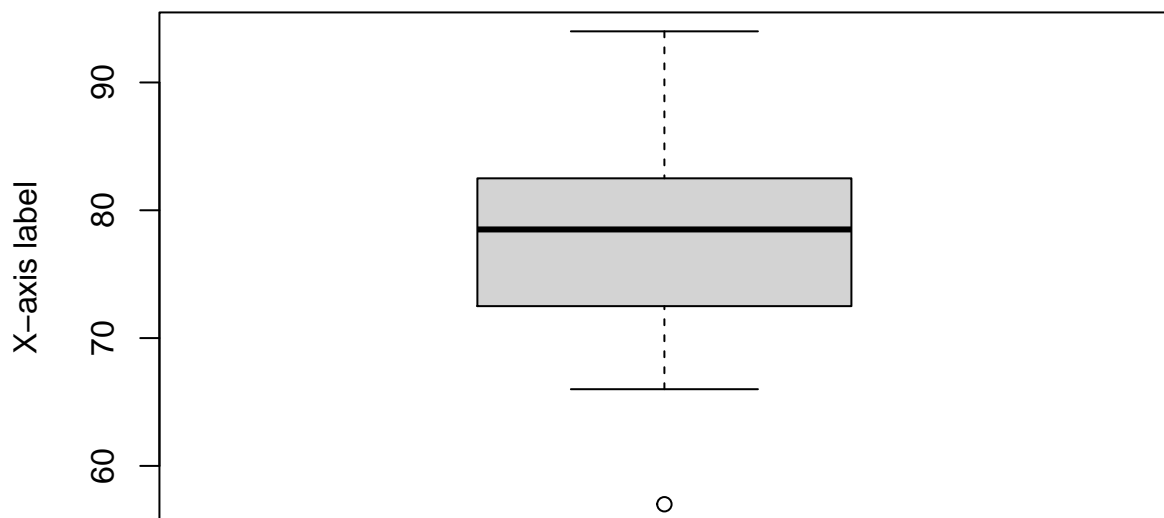
Stats scores. (2.33, p. 78) Below are the final exam scores of twenty introductory statistics students.

57, 66, 69, 71, 72, 73, 74, 77, 78, 78, 79, 79, 81, 81, 82, 83, 83, 88, 89, 94

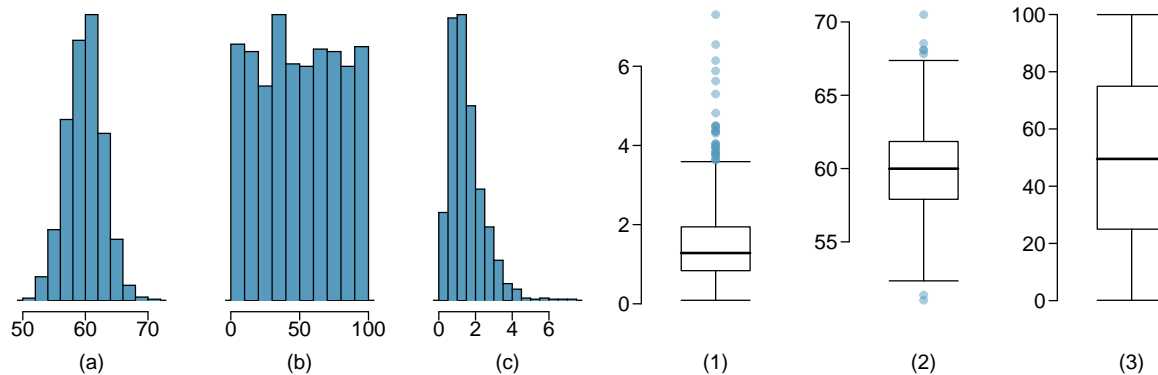
Create a box plot of the distribution of these scores. The five number summary provided below may be useful.

Min	Q1	Q2 (Median)	Q3	Max
57	72.5	78.5	82.5	94

```
##      scores
## Min.    :57.00
## 1st Qu.:72.75
## Median :78.50
## Mean   :77.70
## 3rd Qu.:82.25
## Max.   :94.00
```



Mix-and-match. (2.10, p. 57) Describe the distribution in the histograms below and match them to the box plots.



- (a) Normal Distribution matches boxplot (2)
- (b) Uniform Distribution matches boxplot (3)
- (c) Right skewed Distribution matches boxplot (1)

Distributions and appropriate statistics, Part II. (2.16, p. 59) For each of the following, state whether you expect the distribution to be symmetric, right skewed, or left skewed. Also specify whether the mean or median would best represent a typical observation in the data, and whether the variability of observations would be best represented using the standard deviation or IQR. Explain your reasoning.

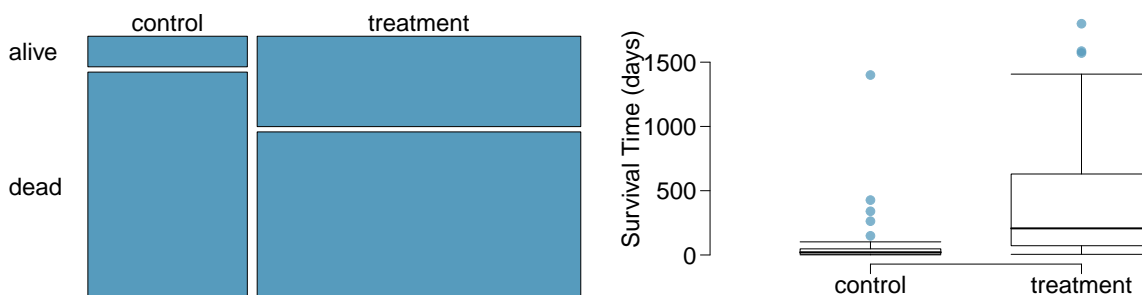
- (a) Housing prices in a country where 25% of the houses cost below \$350,000, 50% of the houses cost below \$450,000, 75% of the houses cost below \$1,000,000 and there are a meaningful number of houses that cost more than \$6,000,000.
- (b) Housing prices in a country where 25% of the houses cost below \$300,000, 50% of the houses cost below \$600,000, 75% of the houses cost below \$900,000 and very few houses that cost more than \$1,200,000.
- (c) Number of alcoholic drinks consumed by college students in a given week. Assume that most of these students don't drink since they are under 21 years old, and only a few drink excessively.
- (d) Annual salaries of the employees at a Fortune 500 company where only a few high level executives earn much higher salaries than the all other employees.

Answers: (a) Right skewed distribution, the meaningful number of houses will skew the numbers to the right. The median would be best to represent the observations because of the right skew affecting the average. The IQR is best used for the variability of observations

- (b) Symmetrical distribution. The mean would be best represent the observations. The standard deviation is best used for the variability of observations
- (c) Left skewed as most of the students under 21 don't drink. The median would be best to represent the observations because of the left skew affecting the average. The IQR is best used for the variability of observations
- (d)

Annual salaries of the employees at a Fortune 500 company where only a few high level executives earn much higher salaries than the all other employees

Heart transplants. (2.26, p. 76) The Stanford University Heart Transplant Study was conducted to determine whether an experimental heart transplant program increased lifespan. Each patient entering the program was designated an official heart transplant candidate, meaning that he was gravely ill and would most likely benefit from a new heart. Some patients got a transplant and some did not. The variable *transplant* indicates which group the patients were in; patients in the treatment group got a transplant and those in the control group did not. Of the 34 patients in the control group, 30 died. Of the 69 people in the treatment group, 45 died. Another variable called *survived* was used to indicate whether or not the patient was alive at the end of the study.



- Based on the mosaic plot, is survival independent of whether or not the patient got a transplant? Explain your reasoning.
- What do the box plots below suggest about the efficacy (effectiveness) of the heart transplant treatment.
- What proportion of patients in the treatment group and what proportion of patients in the control group died?
- One approach for investigating whether or not the treatment is effective is to use a randomization technique.

Answers: (a) Survival is dependent of the patient receiving the transplant because the mosaic plot shows the patient who got transplant are more survival days than the patients who did not get transplant

- Box plots show that patients who got the transplant have increased lifespan compared to who didn't received the treatment
- Proportion of patients in the treatment group and proportion of patients in the control group died

```
transplantProp<-subset(heartTr, heartTr$transplant=="treatment", select=c("survived"))
summary(transplantProp)
```

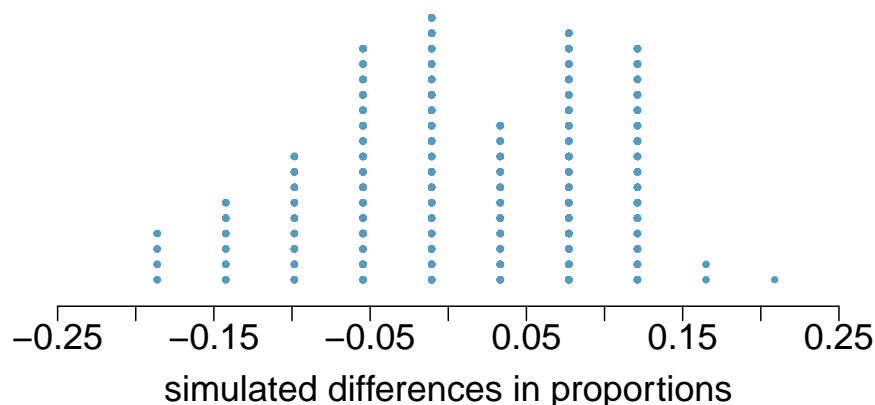
```
## survived
## alive:24
## dead :45
```

- What are the claims being tested?
- The paragraph below describes the set up for such approach, if we were to do it without using statistical software. Fill in the blanks with a number or phrase, whichever is appropriate.

We write *alive* on _____ cards representing patients who were alive at the end of the study, and *dead* on _____ cards representing patients who were not. Then, we shuffle these cards and split them into two groups: one group of size _____ representing treatment, and another group of size _____ representing control. We calculate the difference between the proportion of *dead* cards in the treatment and control groups (treatment - control) and record this value. We repeat this 100 times to build a distribution centered at

_____. Lastly, we calculate the fraction of simulations where the simulated differences in proportions are _____. If this fraction is low, we conclude that it is unlikely to have observed such an outcome by chance and that the null hypothesis should be rejected in favor of the alternative.

iii. What do the simulation results shown below suggest about the effectiveness of the transplant program?



Answers: (i) The claims can be tested on 2 ways i.e null hypothesis and actual hypothesis. Null hypothesis: Survivility is not dependent on if the patient had a transplant or not. Actual hypothesis: Survivility is dependent on the transplant.

(ii) We write *alive* on _____28_____ cards representing patients who were alive at the end of the study, and *dead* on **75**_____ cards representing patients who were not. Then, we shuffle these cards and split them into two groups: one group of size **69**_____ representing treatment, and another group of size _____34_____ representing control. We calculate the difference between the proportion of *dead* cards in the treatment and control groups (treatment - control) and record this value. We repeat this 100 times to build a distribution centered at _____0_____. Lastly, we calculate the fraction of simulations where the simulated differences in proportions are **0.23**_____. If this fraction is low, we conclude that it is unlikely to have observed such an outcome by chance and that the null hypothesis should be rejected in favor of the alternative.

(iii) As per the graph instances are low with the fraction is 0.23. Which is very rare event so we can reject null hypothesis and conclude, having heart transplant affect on the survivility.