Summarizing Quantitative Data: Part 2

- Median and Quartiles
- Five Number Summary
- ➢ Boxplots
- Comparing Histograms and Boxplots
- Independent Comparisons
- Dependent Comparisons



3

PERCENTILE

- Percentile: value at or below which a given percentage of a variable's values fall
 - If the population mean is known, then exact percentiles can be found
 - If we do not know the population mean and only have a sample, then percentiles are only approximations
 - To calculate a percentile in a sample by hand:
 - Sort the data from smallest to largest
 - Calculate the location of percentile: $L_p = (n + 1) \times \frac{p}{100}$
 - Count to position L_p in ordered data
 - If L_p is an integer, use that value as the percentile
 - If L_p is not an integer, weight the values around L_p according to the decimal part

EXAMPLE: PERCENTILES

Scenario: Sample of 20 final exam scores from a large class
33 54 59 62 65 67 69 71 73 74 76 77 79 82 83 85 88 91 94 98
 Question: What is the 60th percentile?
• Answer: • Location: <i>L</i> ₆₀ = =
 12th and 13th Largest Scores:
• 60 th Percentile: <i>P</i> ₆₀ = =
Question: What is the interpretation of this value?
• Answer: of students scored or equal to while scored or equal to

SPECIAL PERCENTILES		
 Median: measure of center equal to the 50th percentile Half of observations lie on either side Odd sample size: Middle number 		
 Even sample size: Average of middle two numbers 		
• 1 st quartile (Q_1): value equal to the 25 th percentile		
• 25% of data is smaller while 75% is larger		
 3rd quartile (Q₃): Value equal to the 75rd percentile 75% of data is smaller while 25% is larger 		
• Quicker Method: Rather than calculating percentiles by hand, we can use the summary function in R.		
SPECIAL PERCENTILES		
• Scenario: Sample of 20 final exam scores from a large class		
33 54 59 62 65 67 69 71 73 74 76 77 79 82 83 85 88 91 94 98		
 Question: What are the the median, 1st quartile, and 3rd quartile? R Code: summary(scores\$score, quantile.type = 6) 		
Min. 1st Qu. Median Mean 3rd Qu. Max. 33.0 65.5 75.0 74.0 84.5 98.0		
• Answer:		
• Median:		
• First Quartile:		
Third Quartile:		
Measures of Spread		
 Range: difference between maximum and minimum values Drawback: One outlier can drastically skew the perception of how spread out the data is 		
 Interquartile range: measures the spread of the middle 50% of the data 		
 IQR = Q₃ - Q₁ Difference between third and first quartiles Benefit: Ignores outliers 		



 Five number summary: quick, straightforward way of summarizing a set of quantitative data that consists of the Minimum, 1st quartile, median, 3rd quartile, and maximum
 Boxplot: graphical display of quantitative data that displays the five number summary and outliers Bottom whisker: Ranges from smallest non-outlier to first quartile Bottom of box: Drawn at first quartile Line through box: Drawn at median Top of box: Drawn at third quartile Top whisker: Ranges from third quartile to largest non-outlier Outliers: represented by asterisks (*)
EXAMPLE: BOXPLOT
 Scenario: Sample of 20 final exam scores from a large class
33 54 59 62 65 67 69 71 73 74 76 77 79 82 83 85 88 91 94 98
The five number summary is 33, 65.5, 75, 84.5, 98.
 Task: Construct the boxplot. Answer:
EXAMPLE: CHANGING AN OBSERVATION
 Scenario: Suppose the student in the sample who scored 33 was replaced by a student who scored 63. Question: What would happen to the mean, median, and standard deviation? Answer: Mean:
54 59 62 <mark>63</mark> 65 67 69 71 73 74 76 77 79 82 83 85 88 91 94 98

Five Number Summary and Boxplot

COMPARING GRAPHS AND MEASURES OF CENTER				
Histogram	Boxplot	Mean vs. Median		
Approximately symmetric	Whiskers about same length	Mean \approx Median		
Left-skewed	Lower whisker longer	Mean < Median		
Right-skewed	Upper whisker longer	Mean > Median		
 Note: Use the mean to describe the center unless the data is severely skewed. 				
 Question: Why does the skewness impact the relationship between the mean and median? 				
Answer:				
If the data is	-skewed, small values pull	the mean		
• If the data is	skewed, large values pull 1	ine mean		
EXAMPLE: COMPAR	ing Graphs and Mi	EASURES OF CENTER ¹⁴		
 Scenario: Boxplot be in the same city. 	elow displays assessed v	alues of 188 houses		
 Question: What would 	ld the shape of the histo	ogram be?		
• Answer:		800000 - •		
Many asses	sed values	•		
• values ter	nd to be spread out	600000 -		
• Ouestion: What is th	a relationshin between			
the mean and media	an?	400000 -		
Answer:				
• Mean: \$;	Median: \$297,618	200000 -		
• observat	tions in tail pull mea	an		
COMPARING ONE	QUANTITATIVE AND	ONE CATEGORICAL		
 Three situations may arise when comparing the relationship between one quantitative and one categorical variable: Two-sample comparison Several sample comparison Paired sample comparison 				
 To compare the dist levels of a categoric Graphically: Side-by Numerically: Five number summ Means and standard 	ributions of the quantita al variable: -side boxplots aries rd deviations	tive variable across		

EXAMPLE: TWO-SAMPLE COMPARISON			
• Scenario: Uber driver tracks the lengths of his passengers' rides for one month and whether the passenger took the ride for business or personal reasons.			
 Question: Which variable is of each type? 			
Answer:			
Quantitative:			
• Ouestion: Do the observations in the categories overlap?			
• Answer:			
Each ride is a and can only be taken for of			
the two reasons			
• Takeaway: These two samples are said to be			
EXAMPLE: TWO-SAMPLE COMPARISON			
• Scenario: Uber driver tracks the lengths of his passengers' rides for one month and whether the passenger took the ride for business or personal reasons.			
Group Mean Std. Dev. Q1 Median Q3			
Business 7.59 3.29 5.08 7.30 9.52 Image: Second seco			
• Question: What do the graphs and statistics reveal about the centers of the groups?			
• Answer: Typical Uber ride for reasons is than a typical Uber ride for reasons			
• are both larger			
Longest personal rides are ride			
EXAMPLE: TWO-SAMPLE COMPARISON			
• Scenario: Uber driver tracks the lengths of his passengers' rides for one month and whether the passenger took the ride for business or personal reasons.			
Group Mean Std. Dev. Q1 Median Q3			
Business 7.59 3.29 5.08 7.30 9.52 Personal 3.33 1.36 2.33 3.35 4.28			
• Question: What do the graphs and statistics reveal about the spreads of the groups?			
• Answer: Lengths of business rides are than			
the lengths of personal rides			
and(4.44 vs. 1.95) are Business ranges from			
• Dusiness ranges from compared to for personal			





