# Summarizing Quantitative Data: Part I

➢ Histograms

➢ Mean

➢ Standard Deviation

➢ Z-Scores

➢ Empirical Rule

**University of Pittsburgh**

---

## Motivation: Describing Quantitative Data

- **Scenario:** Take a random sample of 13 songs from Spotify and record the length (in seconds)

  138 162 178 197 204 209 216 222 231 245 262 273 297

- **Question:** What type of data has been collected?

- **Answer:** _____
  - Data is _____ and it makes sense to _____

- **Question:** What population does this data represent?

- **Answer:** _____

---

## Motivation: Describing Quantitative Data

- **Scenario:** Take a random sample of 13 songs from Spotify and record the length (in seconds)

  138 162 178 197 204 209 216 222 231 245 262 273 297

  Average in the sample is 218 seconds.

- **Question:** Can we conclude that the mean of all observations in the population is greater than 210 seconds?

- **Answer:** _____
  - While the sample may be _____, it is _____ to definitively conclude that the population mean is greater than 210
  - This is an _____ question that can only be answered by looking at not only the _____, but also the _____ of the observations

# DESCRIBING A SINGLE QUANTITATIVE VARIABLE

- **Shape:** reports observations that tend to be more or less common
  - Symmetry/Skewness
  - Modality

- **Center:** measure of a "typical" observation
  - Mean
  - Median

- **Spread:** measure of how different the observations tend to be
  - Standard deviation
  - Range
  - Interquartile range

# HISTOGRAM

- **Histogram:** graphical display that summarizes quantitative data using bars of different heights to show frequencies in predetermined ranges

- To construct a histogram:
  1. Divide the range of data into _____ (or _____)
  2. Count the _____ that fall into each bin
  3. Draw a _____ equal to the _____ for each bin using vertical axis for _____ and horizontal axis for _____
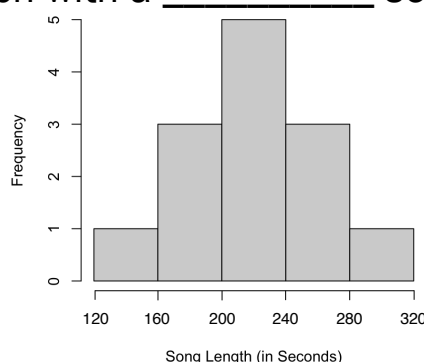
# EXAMPLE: HISTOGRAM

- **Scenario:** Take a random sample of 13 songs from Spotify and record the length (in seconds)

    138 162 178 197 204 209 216 222 231 245 262 273 297

- **Question:** What does the histogram look like?

- **Answer:** ___ bins, each with a _____ seconds

| Bins | Count |
|------|-------|
| 120 to < 160 | |
| 160 to < 200 | |
| 200 to < 240 | |
| 240 to < 280 | |
| 280 to < 320 | |



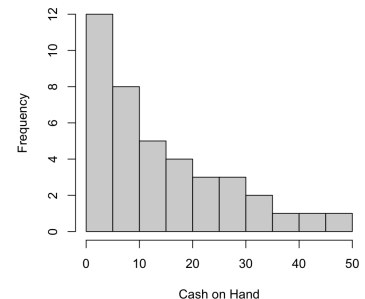*Note: Software typically does the best job of deciding how many bins to use and how wide they should be.*

# EXAMPLE: HISTOGRAM

- **Scenario:** Random sample of 40 people were asked how much cash they are currently carrying

- **Question:** What percentage of people were carrying between $10 and $25?

- **Answer:** ____ (_____ = ____)

- **Question:** What percentage of people were carrying between $12 and $24?

- **Answer:** _____

  - $12 and $24 are _____ (i.e. are not values that _____ two bars)

# DESCRIBING A HISTOGRAM

## Symmetry

- **Symmetric:** similar number of observations on either side of the center

- **Left-skewed:** longer left tail
  - Indicates a few unusually small observations

- **Right-skewed:** longer right tail
  - Indicates a few unusually large observations

## Modality

- **Unimodal:** one distinct peak

- **Bimodal:** two distinct peaks

- **Multimodel:** more than two distinct peaks

- **Uniform:** bars are all of similar height

### Important Special Case
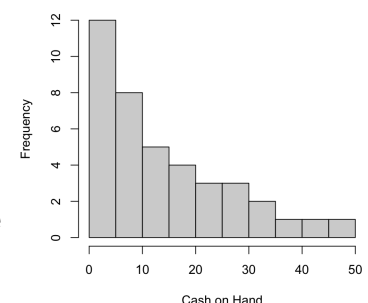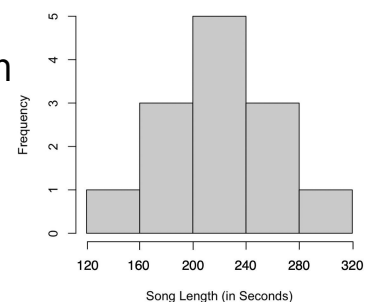**Normal:** symmetric and unimodal (bell-shaped)

# EXAMPLE: HISTOGRAM

- **Scenario:** Histogram of song lengths on top; histogram of amount of cash on hand on bottom

- **Question:** What are the histogram shapes?

- **Answer:**

  - **Song Lengths:** _____
    - Typical song is about _____ minutes long, but _____ _____ to be _____
    - Implies data is _____

  - **Cash on Hand:** _____
    - Most people carry _____
    - As the amount of cash _____, the number of people carrying that amount _____

# MEAN

- **Mean:** the average of all values in a set of quantitative data
  - Measure of the center of the data, or where a typical observation is likely to fall

  - **Population mean:** Denoted by $\mu$ (Greek letter "mu")
    - Parameter → Generally unknown

  - **Sample mean:** Denoted by $\bar{x}$ (Pronounced "x-bar")
    - Statistic → Can always be calculated for quantitative data
    - Good approximation of $\mu$

$$\bar{x} = \frac{x_1 + x_2 + \cdots + x_n}{n}$$

# EXAMPLE: MEAN

- **Scenario:** High temperature in Pittsburgh recorded each weekday during a week in June returned the following observations

  75  76  78  80  81
- **Question:** What is the sample mean high temperature in June?
- **Answer:**

  ___ = _____ = _____ =

- **Question:** What is the population mean high temperature in June?
- **Answer:** _____
  - Only have a _____ of high temperatures→ _____
  - Know probably _____

# STANDARD DEVIATION

- **Standard deviation:** measure of how "spread out" the data is
  - Can be used to measure how far away an observation is from the mean

  - **Population standard deviation:** Denoted by $\sigma$ (Greek letter "sigma")
    - Parameter → Generally unknown

  - **Sample standard deviation:** Denoted by $s$
    - Statistic → Can always be calculated for quantitative data
    - Good approximation of $\sigma$

$$s = \sqrt{\frac{(x_1 - \bar{x})^2 + \cdots + (x_n - \bar{x})^2}{n - 1}}$$

# EXAMPLE: CALCULATING STANDARD DEVIATION

- **Scenario:** High temperature in Pittsburgh recorded each weekday during a week in June returned the following observations

  75  76  78  80  81

- **Question:** What is the sample standard deviation?

- **Answer:**

$$s = \sqrt{\frac{(\underline{\quad})^2 + (\underline{\quad})^2 + (\underline{\quad})^2 + (\underline{\quad})^2 + (\underline{\quad})^2}{\underline{\quad}}}$$

$$= \sqrt{\frac{\underline{\ } + \underline{\ } + \underline{\ } + \underline{\ } + \underline{\ }}{\underline{\ }}}$$

$$= \underline{\quad} = \underline{\quad}$$
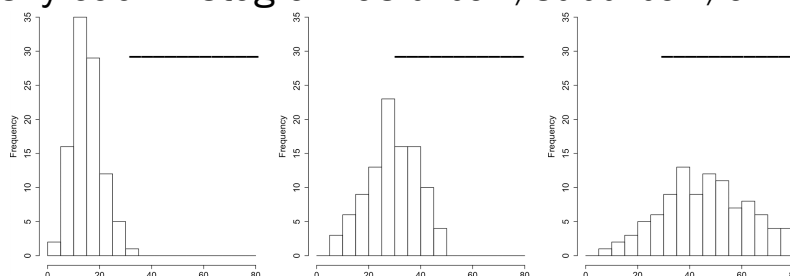
# EXAMPLE: UNDERSTANDING STANDARD DEVIATION

- **Scenario:** High temperature in Pittsburgh recorded each weekday during a week in June returned the following observations

- **Question:** What does the standard deviation tell us about the high temperatures in Pittsburgh in June?

- **Answer:** With a _____ high temperature of ___, it would not be _____ to see the high temperature on a _____ as low as ___ or as high as ___.
  - Both are within _____ of the mean

- **Question:** What is the population standard deviation?

- **Answer:** _____
  - Don't have high temperature for _____
  - Probably _____

# EXAMPLE: UNDERSTANDING STANDARD DEVIATION

- **Scenario:** Commute time to work based on location of residence.

- **Task:** Classify each histogram as urban, suburban, or rural.



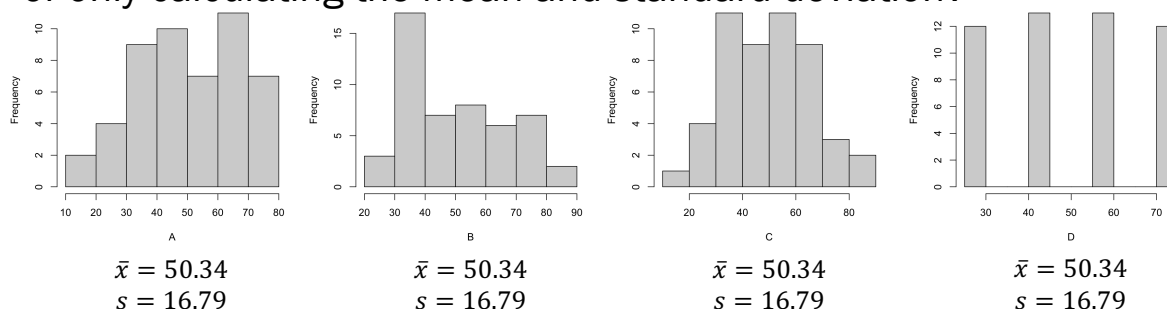- **Task:** Approximate the mean and standard deviation of each group.

- **Answer:**

| Variable | Urban | Suburban | Rural |
|---|---|---|---|
| Mean | | | |
| Standard Deviation | | | |

# IMPORTANCE OF GRAPHICAL AND NUMERIC DISPLAYS

- **Question:** Why is it important to look at a graph of the data instead of only calculating the mean and standard deviation?



| A | B | C | D |
|---|---|---|---|
| $\bar{x} = 50.34$ | $\bar{x} = 50.34$ | $\bar{x} = 50.34$ | $\bar{x} = 50.34$ |
| $s = 16.79$ | $s = 16.79$ | $s = 16.79$ | $s = 16.79$ |

- **Answer:** Statistics do not _____ the _____ of a dataset
  - All four graphs have the _____, but drastically _____
    - A: _____       • C: _____
    - B: _____       • D: _____

---

# STANDARDIZING AND Z-SCORES

- **Standardization:** the process of calculating how many standard deviations an observation is from its mean
  - Results in a Z-score

- **Z-score:** in a sample, a standardized value calculated as:

$$Z = \frac{x - \bar{x}}{s}$$

  - Positive Z-score → Observation _____ than mean
  - Negative Z-score → Observation _____ than mean
  - _____ indicates how unusual the observation is relative to the rest of the data

---

# EXAMPLE: USING Z-SCORES

- **Scenario:** SAT scores have a mean of 1000 and standard deviation 217 while ACT scores have a mean of 21 and standard deviation 5. Student A scored 1080 on the SAT, while Student B scored 28 on the ACT.

- **Question:** Which student performed better relative to the mean of the test they took?

- **Answer:** Student ____

  - **A:** $Z = $ _____ = _____ standard deviations _____ the mean

  - **B:** $Z = $ _____ = _____ standard deviations _____ the mean

- **Takeaway:** Cannot look at _____ between an observation and the mean when comparing – must consider the _____.

# Z-Score Interpretations and Outliers

- Observations that are:
    - Within 1 standard deviation of the mean
        - Quite common and not unusual in any way
    - Between 1 and 2 standard deviations from the mean
        - Relatively common, but slightly unusual
    - Between 2 and 3 standard deviations from the mean
        - Unusual, but not unheard of
    - More than 3 standard deviations from the mean
        - Exist, but are highly unusual and unlikely

- Outlier: an observation that is unusually far away from the rest of the data
    - *Rule of Thumb #1: More three standard deviations away from the mean*

# Example: Z-Score Interpretations

- Scenario: SAT scores have a mean of 1000 and a standard deviation of 217 while ACT scores have a mean of 21 and a standard deviation of 5. Student A scored 1080 on the SAT while Student B scored 28 on the ACT.
- Question: How should we describe each student's performance?
- Answer:
    - A: _____ but _____
    - B: Slightly _____ but still _____

- Question: Would either score be considered an outlier?
- Answer: _____: Both are within __ standard deviations of their mean
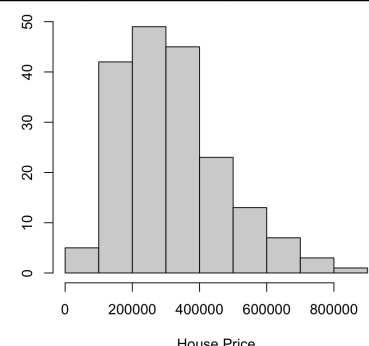
# Example: Impact of Unusual Observations

- Scenario: Histogram below displays assessed values of 188 houses in the same city.
- Question: What would happen to the mean and standard deviation if a house worth $2,000,000 were built?
- Answer: _____

| Variable | Mean | Std. Dev. |
|---|---|---|
| House Price | $317,277 | $156,071 |

- Mean: _____ of house prices is now _____ _____ but sample size only increased by __
    - New value: _____
- Standard Deviation: Data is _____
    - New value: _____

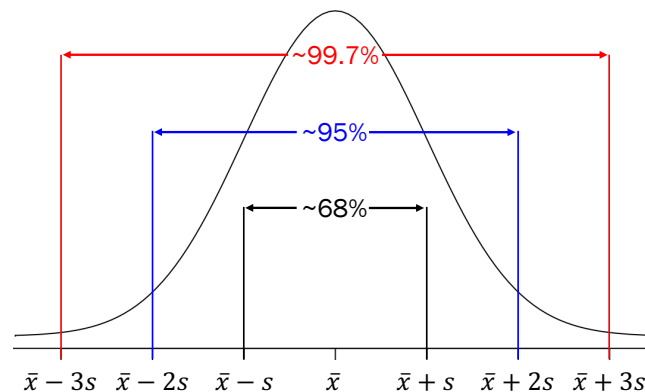- Takeaway: Mean and standard deviation are _____ by outliers.
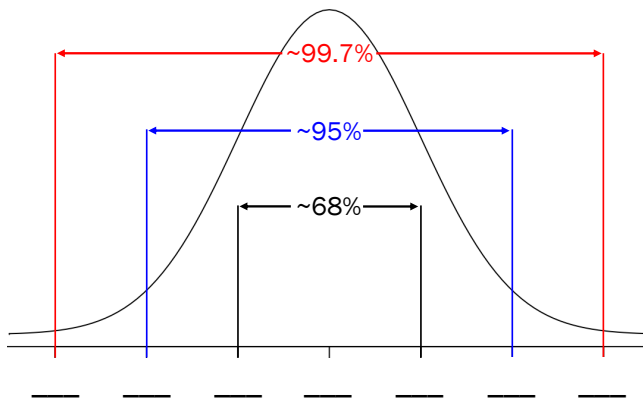
# EMPIRICAL RULE

- **Empirical Rule:** In data that is **approximately normal**, about:
  - 68% of observations are within 1 standard deviation of the mean
  - 95% of observations are within 2 standard deviations of the mean
  - 99.7% of observations are within 3 standard deviations of the mean



~99.7%
~95%
~68%

$\bar{x} - 3s \quad \bar{x} - 2s \quad \bar{x} - s \quad \bar{x} \quad \bar{x} + s \quad \bar{x} + 2s \quad \bar{x} + 3s$

---

# EXAMPLE: EMPIRICAL RULE

- **Scenario:** IQ scores are normally distributed with a mean of 100 and a standard deviation of 15
- **Question:** What are the 7 values of the Empirical Rule?
- **Answer:**



~99.7%
~95%
~68%

___  ___  ___  ___  ___  ___  ___

**Interpretations:**
- About ____% of IQ scores are between ____ and ____
- About ____% of IQ scores are between ____ and ____
- About ____% of IQ scores are between ____ and ____

---

# EXAMPLE: EMPIRICAL RULE

- **Scenario:** IQ scores are normally distributed with a mean of 100 and a standard deviation of 15
- **Question:** What scores correspond to the highest 16% of IQs?
- **Answer:** Scores _____
  - 68% of scores are _____
  - ___% are in the tails (_____ or _____)
  - _____ (____) is in the upper tail
- **Question:** What percentage of IQ scores are between 82 and 124?
- **Answer:** _____
  - Empirical Rule is _____
  - Can answer this using _____



55    70    85    100   115   130   145