

12-1 Descriptive Statistics

Objective:

I can describe a distribution by its shape, outliers, center, and spread.

Vocabulary:

Population: Set of all

Sample: A subset of the population

Parameter: Measures of a population

-Use $\mu = \text{population mean}$

$\sigma = \text{population standard deviation}$

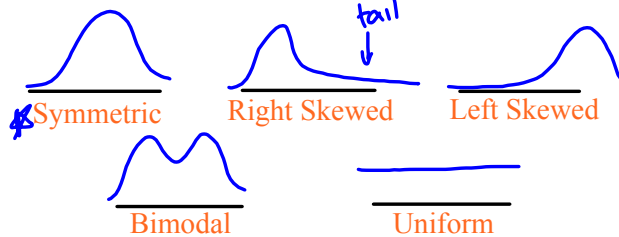
★ Statistics: Measures of a sample

-Use $\begin{cases} \bar{x} = \text{sample mean} \\ s = \text{sample standard deviation} \end{cases}$



"Remember your S.O.C.S"

1. Shape :



2. Outlier : Data far away from the rest of the data. Formula to come ...

3. Center : Measures of central tendency:

1. Mean - arithmetic average of the data
2. Median - Middle value when placed in order, or average of the two middle values
3. Mode - Most frequently occurring value(s)

4. Spread : Measure of the variability in the data

Mean - Median - Mode ?

The average on the test was an 84 - Mean

The average test score puts you in the
middle of the class - Median

The average American student starts
college at 18- Mode

Find the mean, median, and mode for the following set of data:

12, 14, 10, 1, 9, 13, 17, 14, 16

Mean: $11.\bar{7}$ or 11.8
(\bar{x})

Mode: 14

Med: 13

Is there an outlier for the following set? If so, find the mean, median, and mode without the outlier and describe how it affects the data.

Test scores from a class: 70, 70, 75, 75, 90, 70, 80, 85, 65, 95, 70, 85, 90, 70, 20 ←

Mean: 74 ← 77.85

Median: 75 ← 75

Mode: 70

The salaries of the LA Lakers (who makes more than a million a year) for the 2013-2014 season

| | |
|-----------------------------|--------------------------|
| Kobe Bryant: \$30,453,805 | Pau Gasol: \$19,285,850 |
| Steve Nash: \$9,300,500 | Jordan Hill: \$3,563,600 |
| Chris Kaman: \$3,183,000 | Jodie Meeks: \$1,550,000 |
| MarShon Brooks: \$1,210,080 | Nick Young: \$1,106,942 |
| Jordan Farmar: \$1,106,942 | Chris Duhon: \$1,500,000 |

Mean:

Median:

Mode:

Range:

Why do we have all of these measures?

Example: On a cul-de-sac, you have 5 houses built for:

\$200,000, \$200,000, \$200,000, \$200,000,
\$1,200,000

Find the median and the mean? Which one is a better measure?

\$200,000

\$400,000

Median better b/c outlier

Find the standard deviation: Weights in grams of 30 loon chicks

79.5 87.5 88.5 89.2 91.6 84.5 82.1 82.3 85.1 89.8
84.0 84.8 88.2 88.2 82.9 89.8 89.2 94.1 88.0 91.1
91.8 87.0 87.7 88.0 85.4 94.4 91.3 86.3 85.7 86.0

Spread: When we use the median to measure center, we use 5-Number Summary

Range = maximum - minimum

Quartiles split the data into **fourths**

First Quartile (Q_1) = the median of the lower half of the data

Second Quartile = the median

Third Quartile (Q_3) = the median of the upper half of the data

Interquartile Range (IQR) measures the spread between Q_1 and Q_3

$$\text{IQR} = Q_3 - Q_1$$

Five number summary = {minimum, Q_1 , median, Q_3 , maximum}

~~★~~ Box Plot

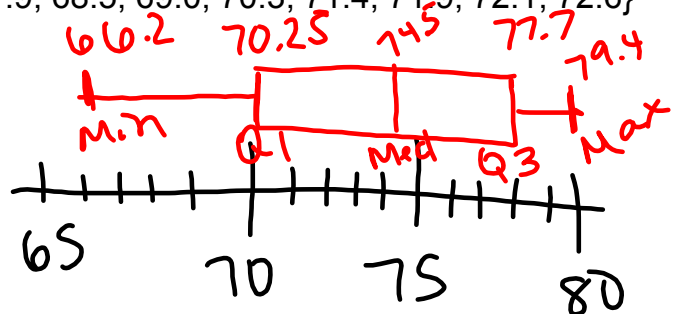
Find the five number summary for the male and female life expectancies in South American nations and compare. Then draw its boxplot.

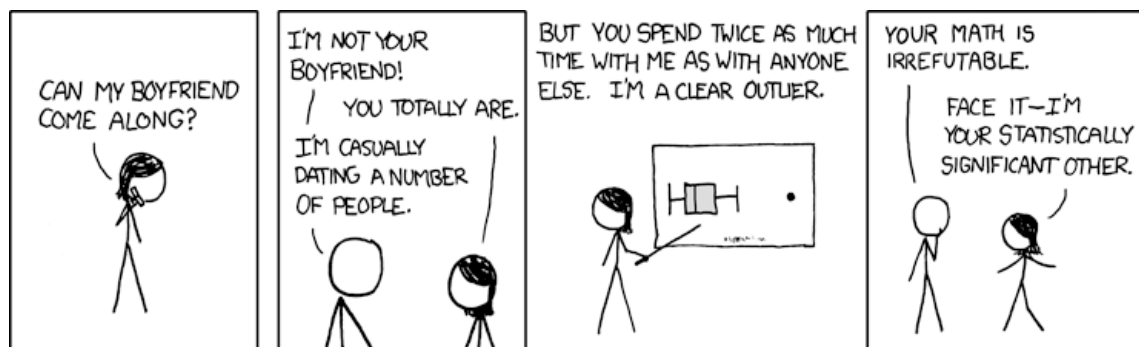
A **box plot** (sometimes called box and whisker plot) is a graph that depicts the five number summary of a data set.

~~★~~ females: {66.2, 66.7, 67.7, 72.8, 74.3, 74.4, 74.6, 76.5, 76.6, 78.8, 79.0, 79.4}

males: {59.0, 60.5, 61.5, 66.7, 67.9, 68.5, 69.0, 70.3, 71.4, 71.9, 72.1, 72.6}

Min: 66.2
[Q_1 : 70.25
Med: 74.5
 Q_3 : 77.7
Max: 79.4





Box and Whisker plots allow us to get a good visual of outliers: a number that makes one of the whiskers noticeably longer than the box:

RULE OF THUMB: a number is considered an outlier if it is more than $1.5 \times \text{IQR}$ below Q_1 or above Q_3

Is 61 an outlier in Roger Maris's home run data?

Five number summary = {5, 11, 19.5, 30.5, 61}

Min Q_1 Med Q_3 Max

$$\text{IQR} = Q_3 - Q_1$$

$$30.5 - 11 = 19.5$$

$$\text{IQR} \times 1.5 = 19.5 \times 1.5 = 29.25$$

$$Q_3 + 29.25 =$$

$$30.5 + 29.25 = 59.75 < 61$$

↑
outlier

$$Q_1 - 16.5 =$$

$$11 - 16.5 = -5.5$$

