

# Econ 325: Introduction to Empirical Economics



## Lecture 1

### Describing Data: Numerical

# Measures of Central Tendency

## Overview

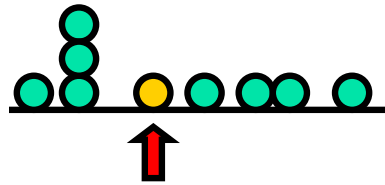
### Central Tendency

Mean

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

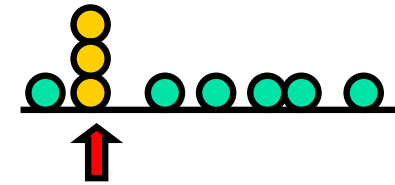
Arithmetic  
average

Median



Midpoint of  
ranked values

Mode



Most frequently  
observed value

# Arithmetic Mean

- The arithmetic mean (mean) is the most common measure of central tendency
  - For a population of N values:

$$\mu = \frac{\sum_{i=1}^N x_i}{N} = \frac{x_1 + x_2 + \cdots + x_N}{N}$$

Population values

Population size

- For a sample of size n:

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{x_1 + x_2 + \cdots + x_n}{n}$$

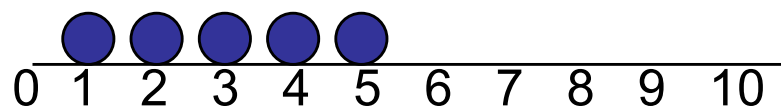
Observed values

Sample size

# Arithmetic Mean

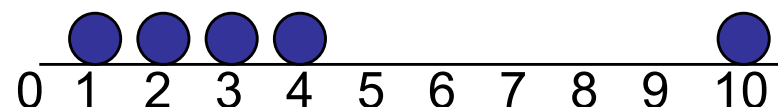
*(continued)*

- The most common measure of central tendency
- Mean = sum of values divided by the number of values
- Affected by extreme values (outliers)



Mean = 3

$$\frac{1+2+3+4+5}{5} = \frac{15}{5} = 3$$

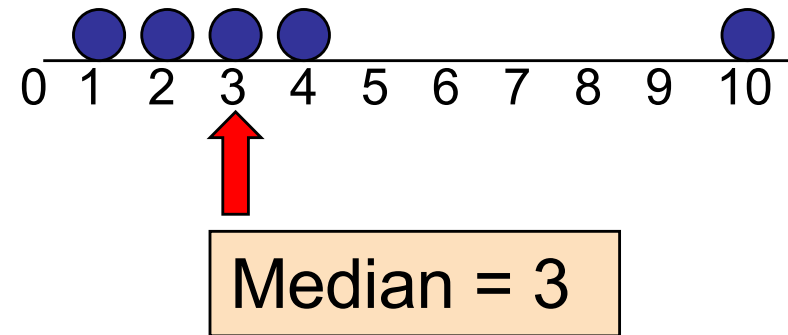
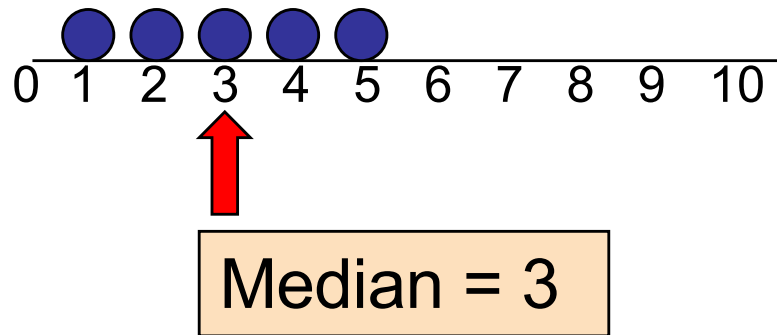


Mean = 4

$$\frac{1+2+3+4+10}{5} = \frac{20}{5} = 4$$

# Median

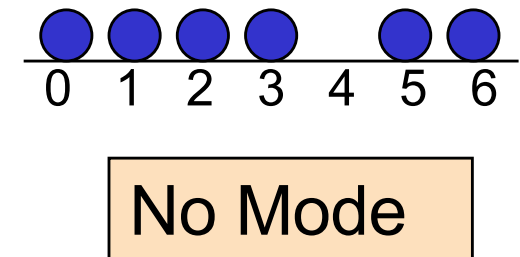
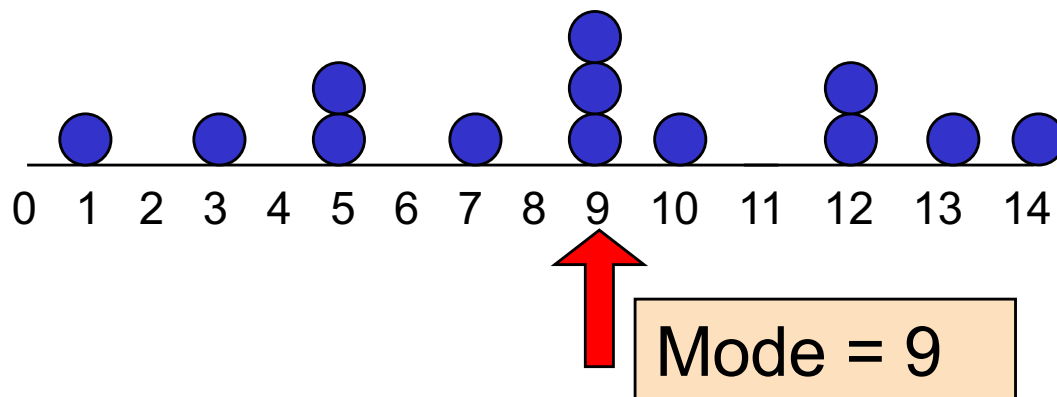
- In an ordered list, the median is the “middle” number (50% above, 50% below)



- Not affected by extreme values

# Mode

- Value that occurs most often
- Not affected by extreme values
- Used for either numerical or categorical data
- There may may be no mode
- There may be several modes





# Which measure of location is the “best”?

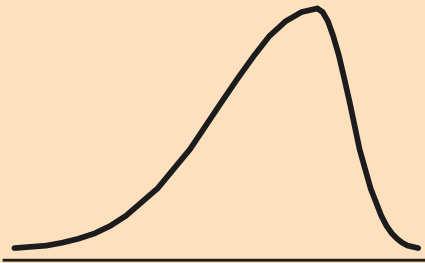
---

- **Mean** is generally used, unless extreme values (outliers) exist . . .
- Then **median** is often used, since the median is not sensitive to extreme values.
  - **Example:** Median home prices may be reported for a region – less sensitive to outliers

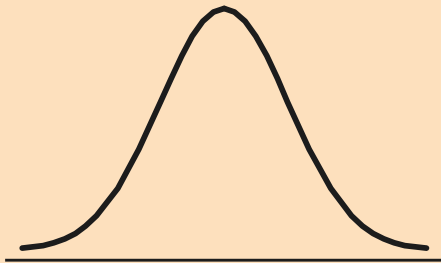
# Clicker Question 1.1

Q: In which graph, Mean < Median ?

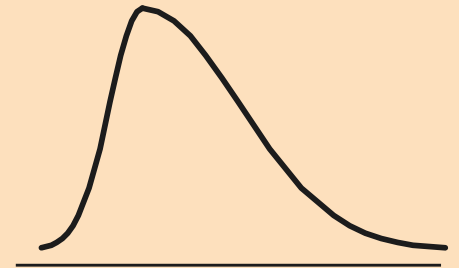
A). Left-Skewed



B). Symmetric



C). Right-Skewed



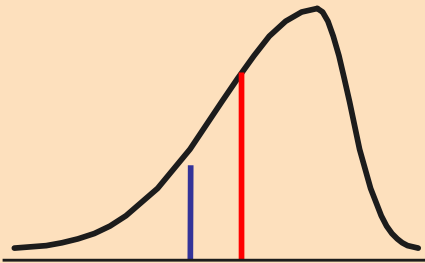


# Clicker Question 1.1

- Answer: A). “Left-Skewed” distribution.

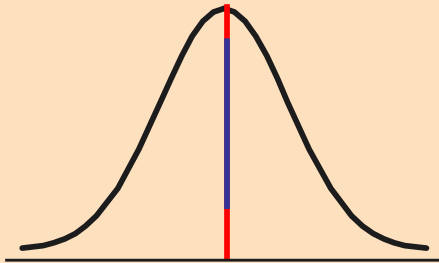
Left-Skewed

Mean < Median



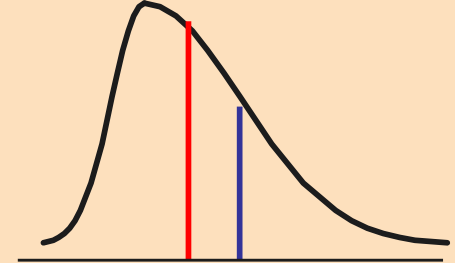
Symmetric

Mean = Median



Right-Skewed

Median < Mean





# Geometric Mean

---

- Geometric mean

$$\bar{X}_g = \sqrt[n]{(X_1 \times X_2 \times \cdots \times X_n)} = (X_1 \times X_2 \times \cdots \times X_n)^{1/n}$$

- Geometric mean rate of return

$$\bar{r}_g = (X_1 \times X_2 \times \cdots \times X_n)^{1/n} - 1$$



## Clicker Question 1.2

---

Suppose you invested \$100 in stocks and, after 5 years, the value of stocks becomes \$125 worth. What is the **average annual compound rate of returns**?

- A). 5 %
- B). 4.6 %
- C). 5.4 %



# Clicker Question 1.2

---

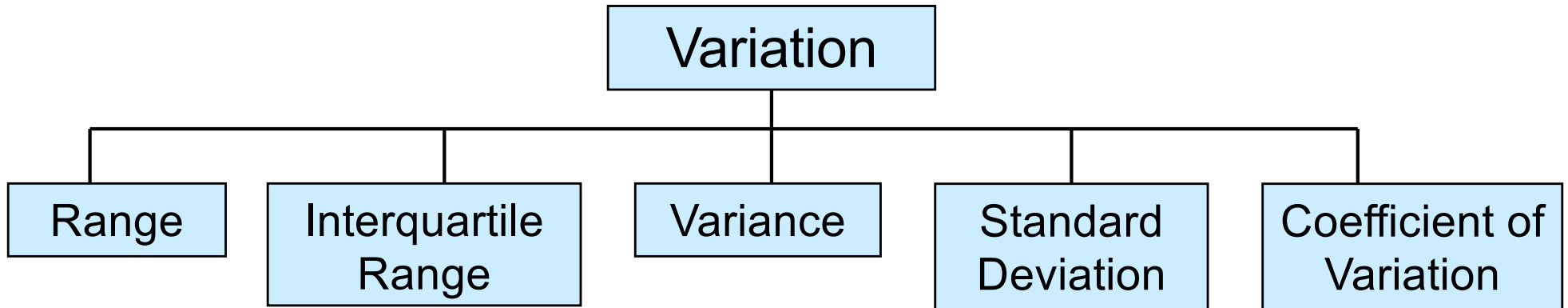
Initial Investment: \$100

After 5 years: \$125

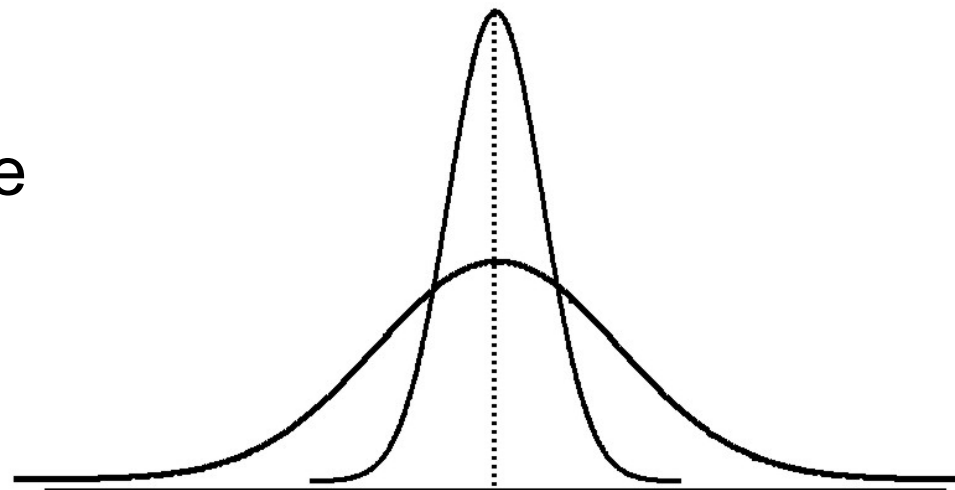
Question: the **average annual** rate of returns?

- 25 % divided by 5 years = 5%? This is Wrong!
- “Compound Interest” over 5 years  
 $\$100 \times (1+0.05)^5 = \$127.63 > \$125$  after 5 years
- **Answer is B**  
 $\$100 \times (1+r)^5 = \$125 \rightarrow r = 4.6\%$

# Measures of Variability



- Measures of variation give information on the **spread** or **variability** of the data values.



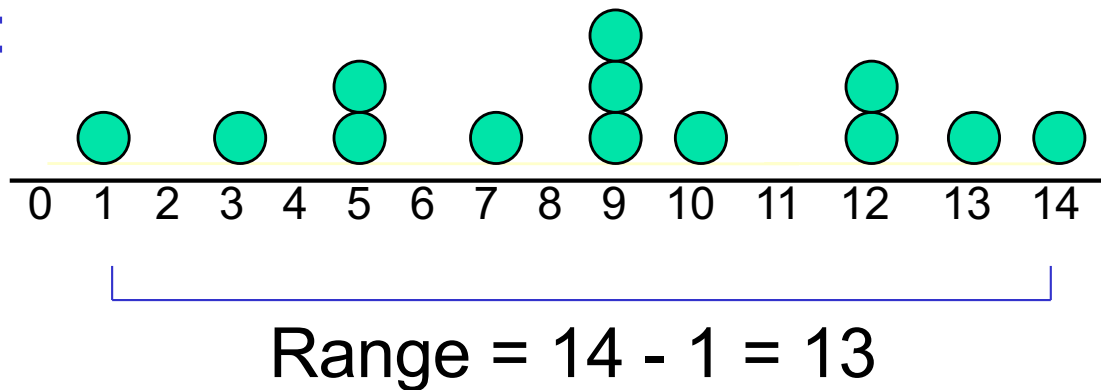
Same center,  
different variation

# Range

- Simplest measure of variation
- Difference between the largest and the smallest observations:

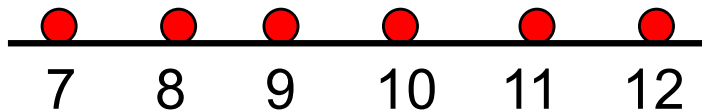
$$\text{Range} = X_{\text{largest}} - X_{\text{smallest}}$$

Example:

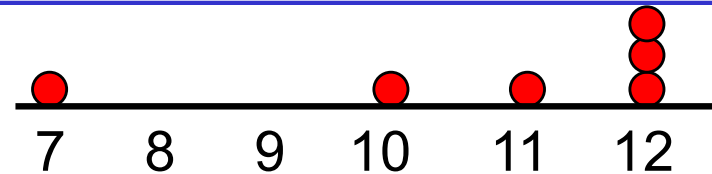


# Disadvantages of the Range

- Ignores the way in which data are distributed



$$\text{Range} = 12 - 7 = 5$$



$$\text{Range} = 12 - 7 = 5$$

- Sensitive to outliers

1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 2, 2, 2, 2, 2, 2, 2, 2, 3, 3, 3, 3, 4, 5

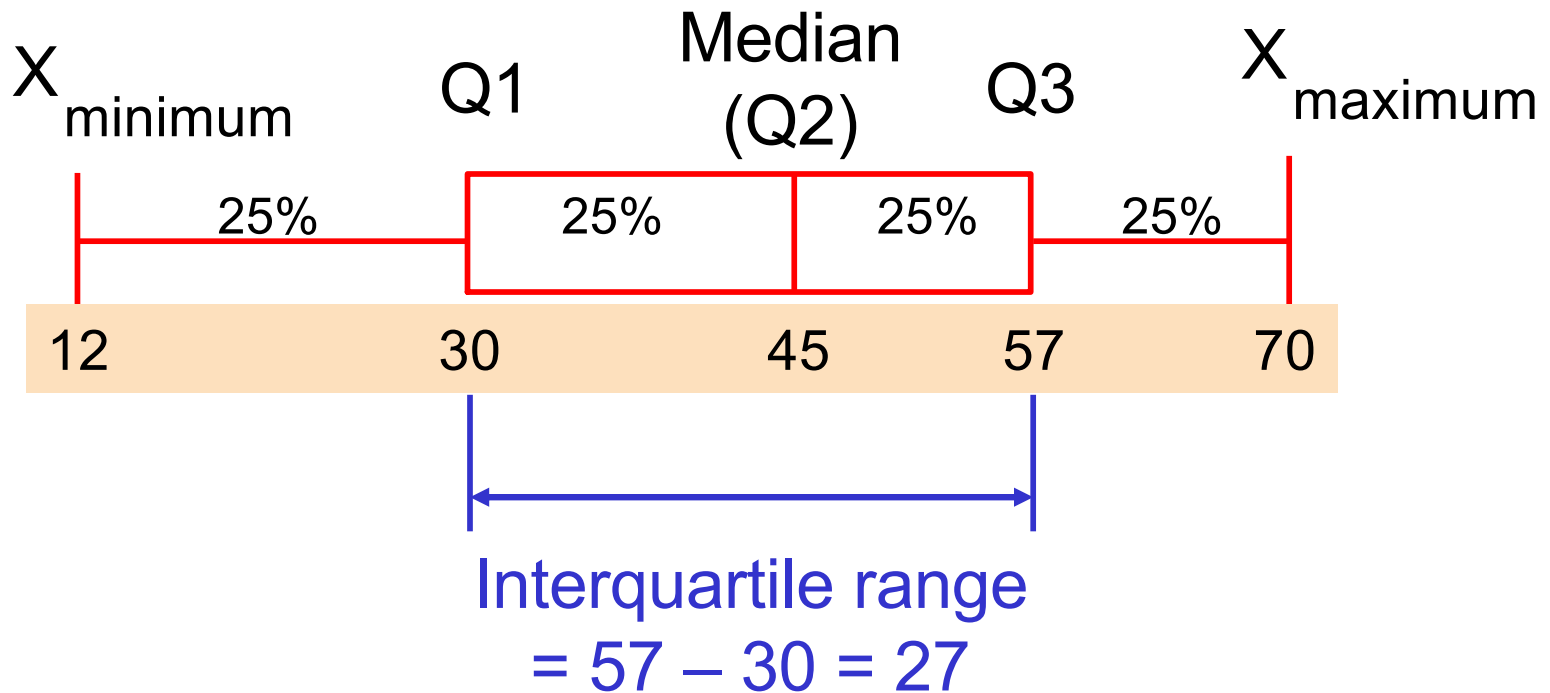
$$\text{Range} = 5 - 1 = 4$$

1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 2, 2, 2, 2, 2, 2, 2, 2, 3, 3, 3, 3, 4, 120

$$\text{Range} = 120 - 1 = 119$$

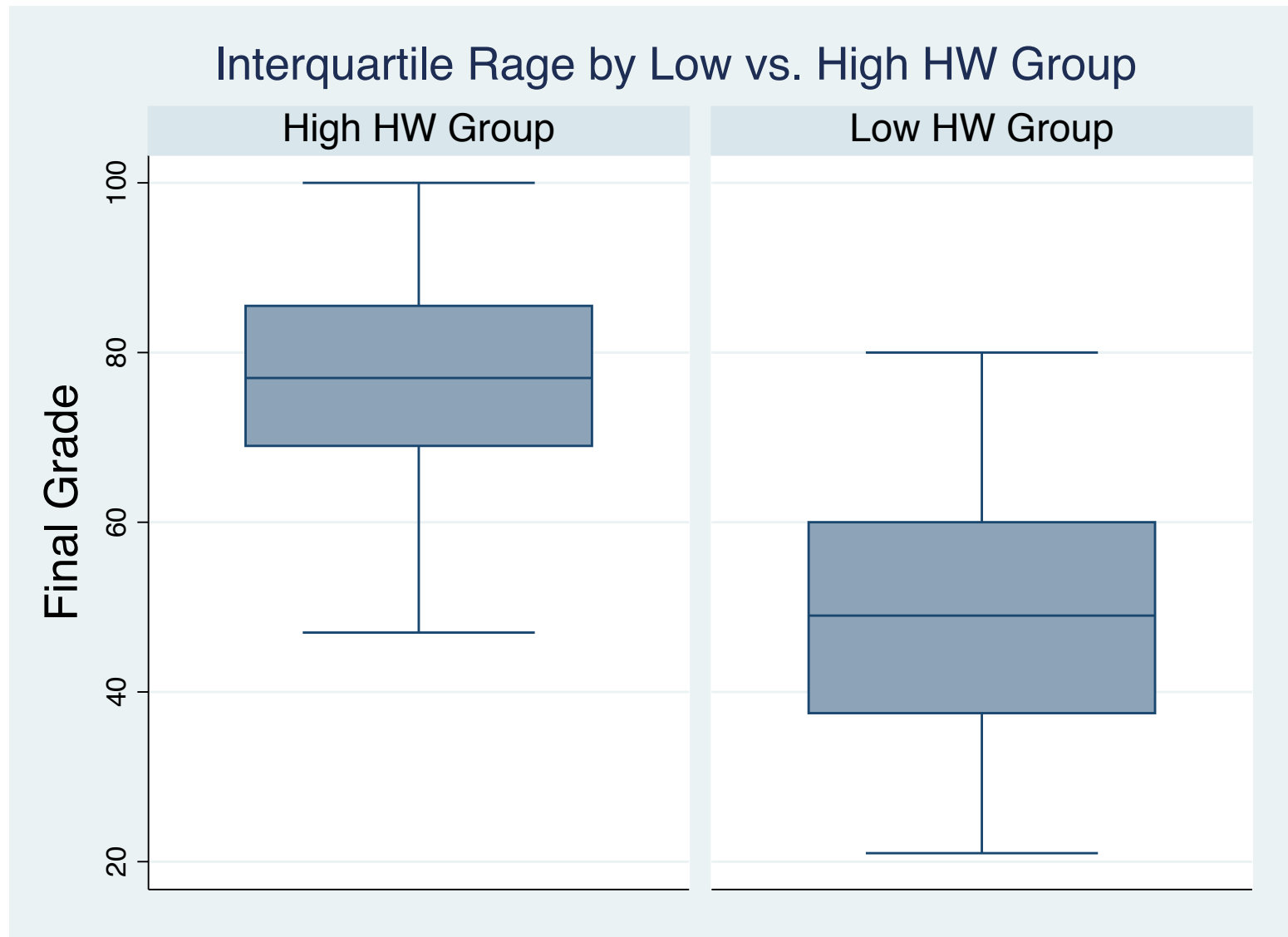
# Interquartile Range

Example:





# STATA Example





# Population Variance

- Average of squared deviations of values from the mean

- Population variance:

$$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$$

Where

$\mu$  = population mean

$N$  = population size

$x_i$  =  $i^{\text{th}}$  value of the variable  $x$



# Sample Variance

- Average (approximately) of squared deviations of values from the mean

- Sample variance:

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

Where  $\bar{X}$  = arithmetic mean

$n$  = sample size

$X_i$  =  $i^{\text{th}}$  value of the variable  $X$



# Population Standard Deviation

---

- Most commonly used measure of variation
- Shows variation about the mean
- Has the **same units as the original data**
- Population standard deviation:

$$\sigma = \sqrt{\sigma^2} = \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{N}}$$



# Sample Standard Deviation

- Most commonly used measure of variation
- Shows variation about the mean
- Has the **same units as the original data**

- Sample standard deviation:

$$s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}}$$



# Calculation Example: Sample Standard Deviation

Sample

Data ( $x_i$ ):

10 12 14 15 17 18 18 24

$n = 8$

Mean =  $\bar{x} = 16$

$$s = \sqrt{\frac{(10 - \bar{x})^2 + (12 - \bar{x})^2 + (14 - \bar{x})^2 + \dots + (24 - \bar{x})^2}{n - 1}}$$

$$= \sqrt{\frac{(10 - 16)^2 + (12 - 16)^2 + (14 - 16)^2 + \dots + (24 - 16)^2}{8 - 1}}$$

$$= \sqrt{\frac{126}{7}}$$

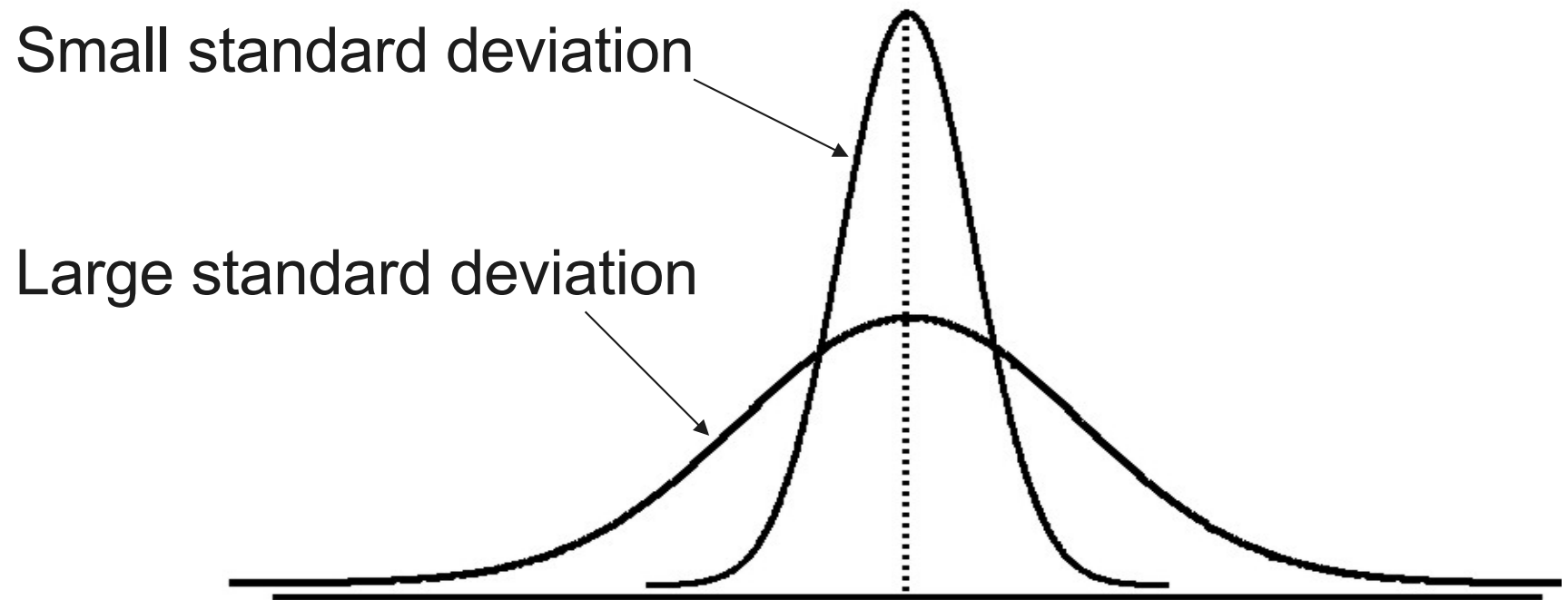
=

4.2426



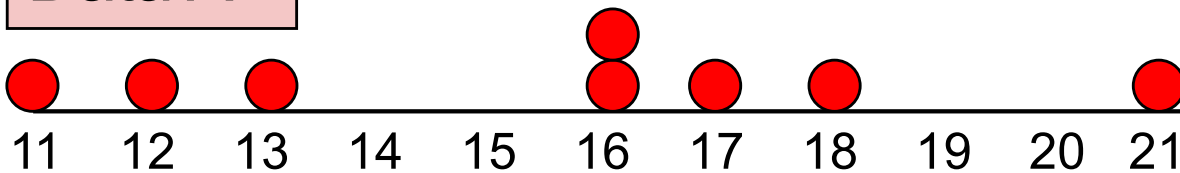
A measure of the “average”  
scatter around the mean

# Measuring variation



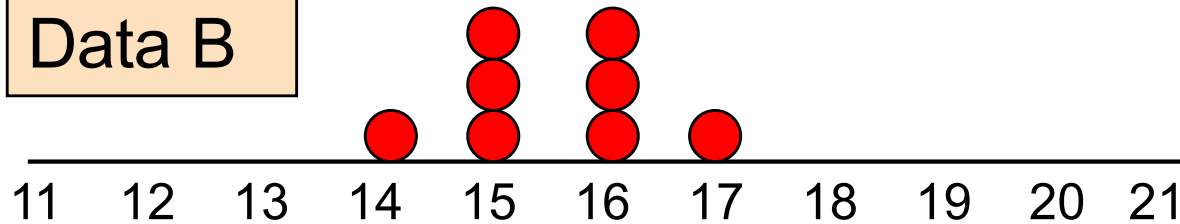
# Comparing Standard Deviations

Data A



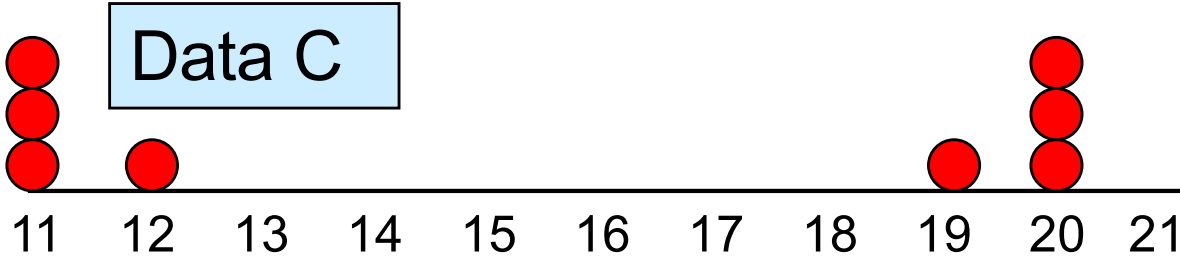
Mean = 15.5  
s = 3.338

Data B



Mean = 15.5  
s = 0.926

Data C



Mean = 15.5  
s = 4.570





# Advantages of Variance and Standard Deviation

- Each value in the data set is used in the calculation
- Values far from the mean are given extra weight  
(because deviations from the mean are squared)



# Coefficient of Variation

---

- Measures **relative variation**
- Always in percentage (%)
- Shows **variation relative to mean**
- Can be used to compare two or more sets of data measured in different units

$$CV = \left( \frac{s}{\bar{x}} \right) \cdot 100\%$$

# Comparing Coefficient of Variation

- Stock A:

- Average price last year = \$50
- Standard deviation = \$5

$$CV_A = \left( \frac{s}{\bar{x}} \right) \cdot 100\% = \frac{\$5}{\$50} \cdot 100\% = 10\%$$

- Stock B:

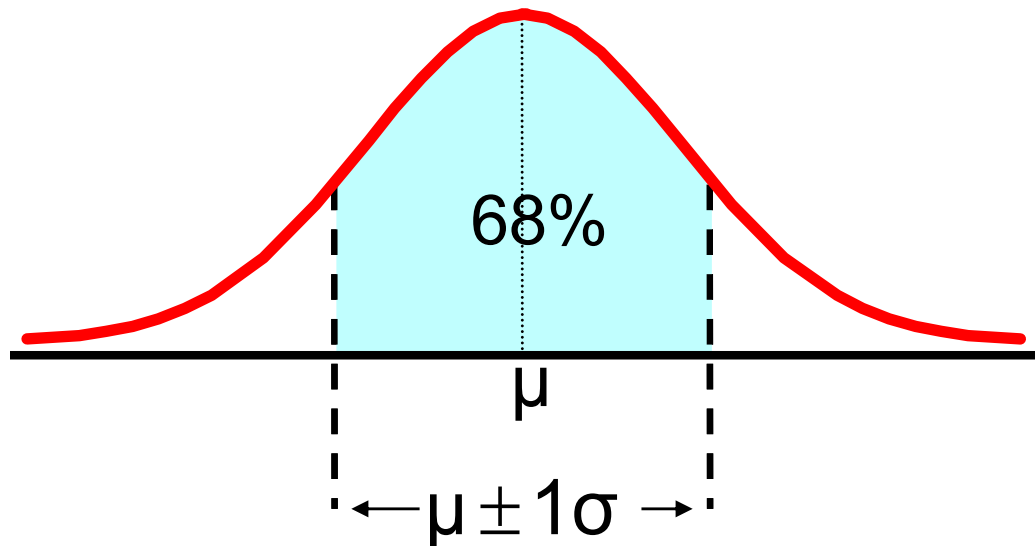
- Average price last year = \$100
- Standard deviation = \$5

$$CV_B = \left( \frac{s}{\bar{x}} \right) \cdot 100\% = \frac{\$5}{\$100} \cdot 100\% = 5\%$$

Both stocks have the same standard deviation, but stock B is less variable relative to its price

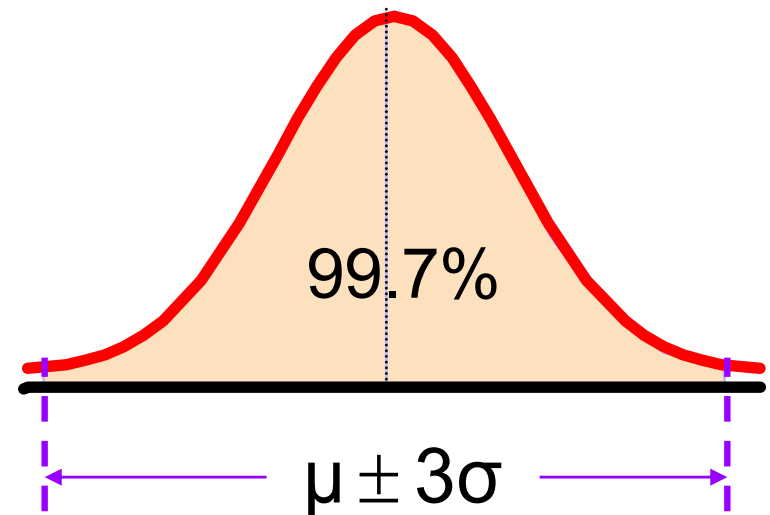
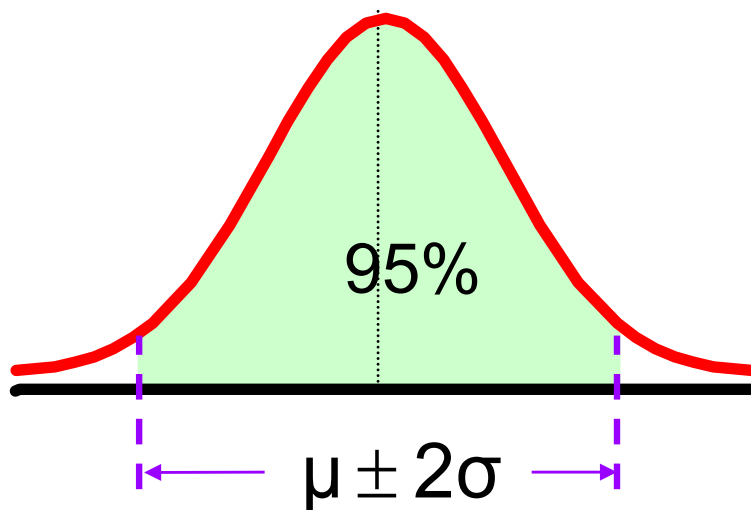
# The Empirical Rule

- If the data distribution is approximated by normal distribution, then the interval:
- $\mu \pm 1\sigma$  contains about 68% of the values in the population or the sample



# The Empirical Rule

- $\mu \pm 2\sigma$  contains about **95%** of the values in the population or the sample
- $\mu \pm 3\sigma$  contains **almost all** (about **99.7%**) of the values in the population or the sample





## Clicker Question 1.3

---

- In one year, the average stock price of Google Inc. was \$800 with the standard deviation equal to \$100. In what interval, approximately 95% of the stock price of Google Inc. will be?

A). Between \$400 and \$1200

B). Between \$600 and \$1000

C). Between \$700 and \$900



## Clicker Question 1.3

- Because  $\mu \pm 2\sigma$  contains about 95% of the values,

$$800 \pm 2 \times 100 = 800 - 200 \text{ and } 800 + 200$$

contains about 95% of the Google stock price.

**Answer: B). Between \$600 and \$1000**

# Weighted Mean

- The **weighted mean** of a set of data is

$$\bar{X} = \sum_{i=1}^n w_i X_i = w_1 X_1 + w_2 X_2 + \cdots + w_n X_n$$

- Where  $w_i$  is the weight of the  $i^{\text{th}}$  observation  
and  $\sum w_i = 1$
- Use when data is already grouped into  $n$  classes, with  $w_i$  values in the  $i^{\text{th}}$  class



# Example

- Consider a student with the scores of assignment ( $x_1$ ), clicker ( $x_2$ ), midterm ( $x_3$ ), and final exam ( $x_4$ ) given by

$$x_1 = 100, x_2 = 0, x_3 = 90, x_4 = 70$$

- The weights:

$$w_1 = 0.2, w_2 = 0.05, w_3 = 0.30, w_4 = 0.45.$$

- The final grade for this student is

$$\sum_{i=1}^4 w_i x_i = 20 + 0 + 27 + 31.5 = 78.5$$

# The Sample Covariance

- The covariance measures the strength of the linear relationship between **two variables**
- The **population covariance**:

$$\text{Cov}(x, y) = \sigma_{xy} = \frac{\sum_{i=1}^N (x_i - \mu_x)(y_i - \mu_y)}{N}$$

- The **sample covariance**:

$$\text{Cov}(x, y) = s_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

- Only concerned with the strength of the relationship
- No causal effect is implied
- Depends on the unit of measurement



# Interpreting Covariance

---

- **Covariance** between two variables:

$\text{Cov}(x,y) > 0 \rightarrow$  x and y tend to move in the **same** direction

$\text{Cov}(x,y) < 0 \rightarrow$  x and y tend to move in **opposite** directions

$\text{Cov}(x,y) = 0 \rightarrow$  x and y are independent



# Coefficient of Correlation

---

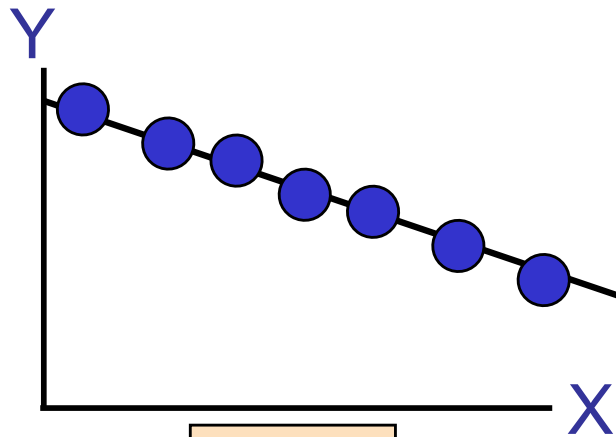
- Measures the relative strength of the linear relationship between two variables
- Population correlation coefficient:

$$\rho = \frac{\text{Cov}(x, y)}{\sigma_X \sigma_Y}$$

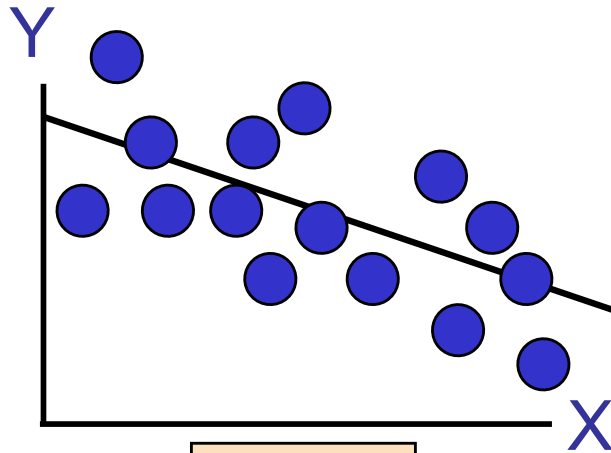
- Sample correlation coefficient:

$$r = \frac{\text{Cov}(x, y)}{s_X s_Y}$$

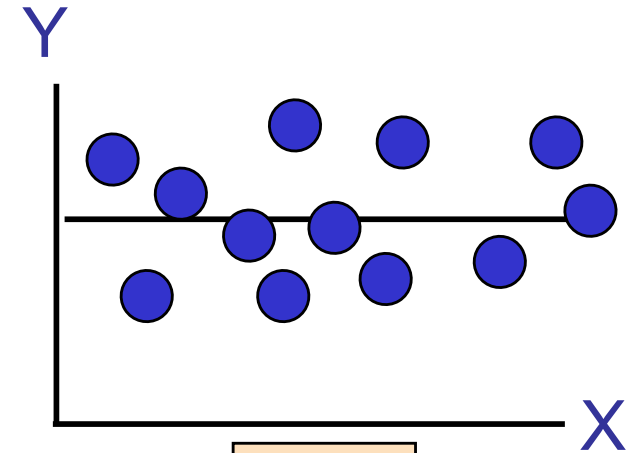
# Scatter Plots of Data with Various Correlation Coefficients



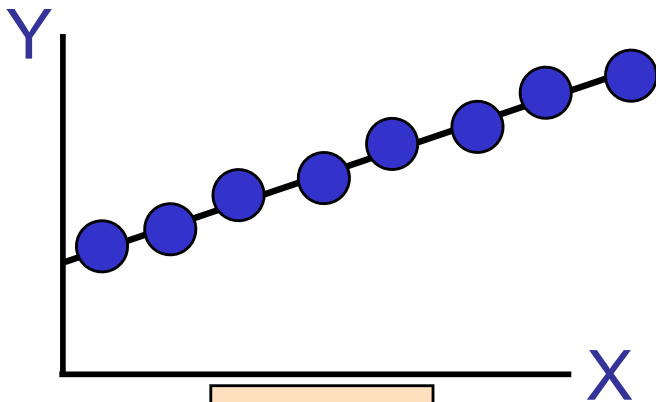
$r = -1$



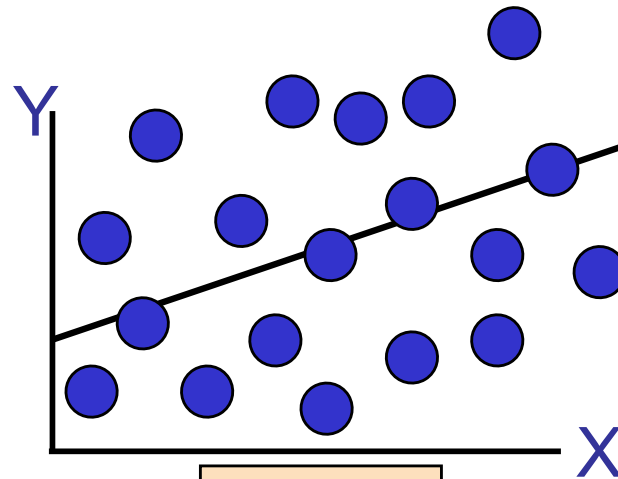
$r = -.6$



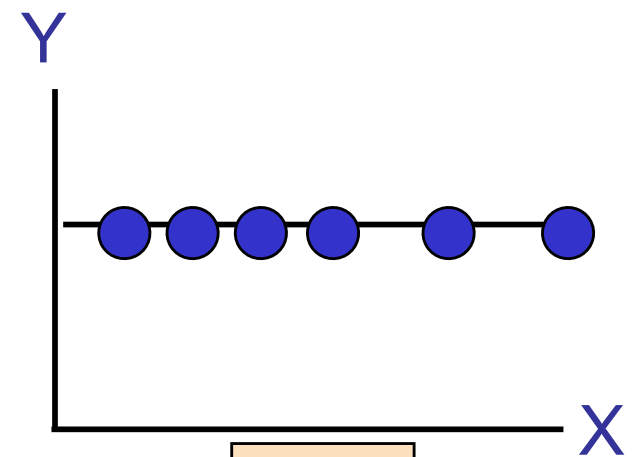
$r = 0$



$r = +1$



$r = +.3$



$r = 0$

# Example: HW and Final Grade

- $r = 0.499$
- There is a **relatively strong positive linear relationship** between HW scores and Final Grades
- Students who scored high on HW assignment tended to have high final grades

