# STATISTICS
## INFORMED DECISIONS USING DATA
### Fifth Edition

STATISTICS
INFORMED DECISIONS USING DATA 5e

Michael Sullivan III

# Chapter 3

Numerically Summarizing Data

Pearson

# 3.2 Measures of Dispersion (Spread)
## Learning Objectives

1. Determine the **range** of a variable from raw data

2. Determine the **standard deviation** of a variable from raw data

3. Determine the **variance** of a variable from raw data

4. Use the **Empirical Rule** to describe data that are bell shaped

5. ~~Use Chebyshev's Inequality to describe any data set~~

Pearson

The **range, *R*,** of a variable is the difference between the largest data value and the smallest data values. That is,

$$\boxed{R = max - min}$$ *Quick & easy way to measure spread.*

Range = *R* = Largest Data Value − Smallest Data Value

Units of Range: *Same as data* $\left( ex: 98° - 30° = 68° \right)$

Why study the range? What does it tell us about our data?

*Quick → spread*

*Limitations only uses 2 values of data.*

Pearson

**EXAMPLE Finding the Range of a Set of Data**

The following data represent the travel times (in minutes) to work for all seven employees of a start-up web development company.

$$23, 36, 23, 18, \underset{\text{min}}{\textcircled{5}}, 26, \underset{\text{max}}{\textcircled{43}}$$

Find the range.

$$R = 43 - 5 = \boxed{38 \text{ min}}$$

UNITS

Pearson

# 3.2 Measures of Dispersion
## 3.2.2 Introducing Standard Deviation

Ex: Advil and Motrin IB produce the same headache relief medication with the active ingredient ibuprofen. Each pill should contain 200 mg of ibuprofen. A health agency obtains a sample of ten tablets from both manufacturers and measures how much ibuprofen each pill actually contains.

$\overline{x}_1$

$\overline{x}_2$

| | Number of milligrams measured | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Advil | 199.25 | 198.50 | 200.10 | 200.75 | 201.00 | 198.00 | 200.10 | 199.00 | 201.10 | 202.20 |
| Motrin IB | 205.00 | 195.80 | 195.20 | 203.20 | 205.80 | 194.40 | 204.60 | 194.60 | 207.20 | 194.20 |

way high way low

Each sample has a mean value of 200 mg. However, based on the given sample values, which company would you prefer to buy from?

$\overline{x}_1 = 200 \, mg = \overline{x}_2$

which prefer? why?

Advil pills more close to 200 mg

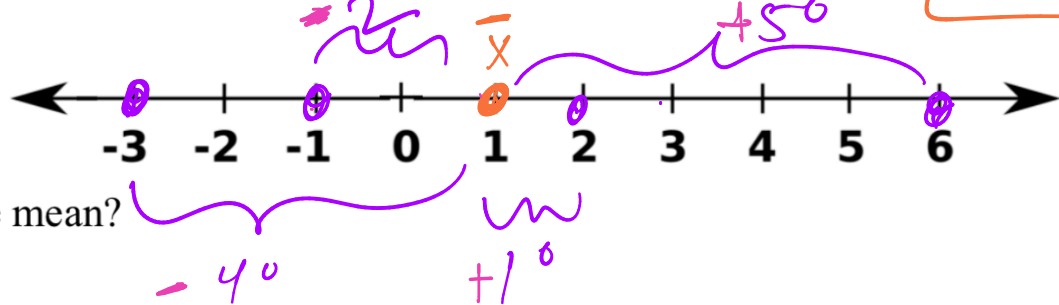why better :· high does can be harmful
· manufacture : can loose money
· low does bad b/c feel cheeted!

Pearson

# 3.2 Measures of Dispersion
## 3.2.2 Introducing Standard Deviation

Ex: The following are temperatures (in degrees) on four consecutive days in Mongolia in January: $-3, -1, 2, 6$

(a) Find the mean.

(b) How far away is each number from the mean?

*[handwritten annotations]*

$\bar{x}$

$-3 \quad -2 \quad -1 \quad 0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6$

$+5°$

$-4°$   $+1°$

a) $\bar{x} = \dfrac{-3-1+2+6}{4} = \dfrac{4}{4} = 1.0$

b) picture

why do this?

take into account how each value compares to mean.

- signed variation
- look total variation

problem: total signed variation = 0.

Pearson

# 3.2 Measures of Dispersion

The **population standard deviation** of a variable is the square root of the sum of squared deviations about the population mean divided by the number of observations in the population, $N$.

That is, it is the square root of the mean of the squared deviations about the population mean.

**Notation:** population standard deviation is symbolically represented by $\sigma$ (lowercase Greek "sigma").

Units of Standard Deviation:

*Same units as data.*

Pearson

The **population standard deviation** of a variable is the square root of the mean of the squared deviations about the population mean.

$\mu$ — pop. mean

$$\sigma = \sqrt{\frac{(x_1 - \mu)^2 + (x_2 - \mu)^2 + \cdots + (x_N - \mu)^2}{N}}$$

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

where $x_1$, $x_2$, . . . , $x_N$ are the $N$ observations in the population and $\mu$ is the population mean.

P Pearson

# 3.2 Measures of Dispersion

Notation:

- **population standard deviation:** $\sigma$

- **sample standard deviation:** $s$ (small s)

the actual # not that important

simple   is it big   or   small ?

Pearson

**IMPORTANT**

Why study the standard deviation?

*tells us about how spread out our data is*

What does it tell us about our data?

*Better than range b/c it incorporates*

*every value of data!*

**EXAMPLE Computing a Population Standard Deviation**

The following data represent the travel times (in minutes) to work for all seven employees of a start-up web development company.

$$23, 36, 23, 18, 5, 26, 43$$

Compute the population standard deviation of this data.

Mean: $\mu = 24.9$ min per employee

Pearson

# 3.2 Measures of Dispersion $\mu = 24.9$

## 3.2.2 Determine the Standard Deviation of a Variable from Raw Data (5 of 16)

$L1$     $L2$     $\boxed{L3 = L2 - L1}$     $L4 = L3^2$

| $X_i$ | $\mu$ | $X_i - \mu$ | $(X_i - \mu)^2$ |
|---|---|---|---|
| 23 | 24.9 | $23 - 24.9 = -1.9$ | 3.61 |
| 36 | 24.9 | $36 - 24.9 = 11.1$ | 123.21 |
| 23 | 24.9 | $23 - 24.9 = -1.9$ | 3.61 |
| 18 | 24.9 | $18 - 24.9 = -6.9$ | 47.61 |
| 5 | 24.9 | $5 - 24.9 = -19.9$ | 396.01 |
| 26 | 24.9 | $26 - 24.9 = 1.1$ | 1.21 |
| 43 | 24.9 | $43 - 24.9 = 18.1$ | 327.61 |

$$\boxed{\Sigma(X_i - \mu)^2 = 902.87}$$

$$\sigma = \sqrt{\frac{\Sigma(X_i - \mu)^2}{N}} = \sqrt{\frac{902.87}{7}} \approx 11.4 \text{ minutes}$$

↖ round w/ 1

"Stats Law"

Pearson

# 3.2 Measures of Dispersion

| $x_i$ | $\mu$ | $x_i - \mu$ | $(x_i - \mu)^2$ |
|---|---|---|---|
| 23 | 24.85714 | −1.85714 | 3.44898 |
| 36 | 24.85714 | 11.14286 | 124.1633 |
| 23 | 24.85714 | −1.85714 | 3.44898 |
| 18 | 24.85714 | −6.85714 | 47.02041 |
| 5 | 24.85714 | −19.8571 | 394.3061 |
| 26 | 24.85714 | 1.142857 | 1.306122 |
| 43 | 24.85714 | 18.14286 | 329.1633 |

$$\Sigma(x_i - \mu)^2 = 902.8571$$

$$\sigma = \sqrt{\frac{\Sigma(x_i - \mu)^2}{N}} = \sqrt{\frac{902.8571}{7}} \approx 11.36 \text{ minutes}$$

*[handwritten notes: check w) "1 VAR STATS" "σx" =]*

Pearson

The **sample standard deviation**, *s*, of a variable is the square root of the sum of squared deviations about the sample mean **divided by *n* − 1**, where *n* is the sample size.

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

$$= \sqrt{\frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \cdots + (x_n - \bar{x})^2}{n - 1}}$$

where $x_1, x_2, \ldots, x_n$ are the *n* observations in the sample and $\bar{X}$ is the sample mean.

Pearson

We call $n - 1$ the **degrees of freedom** because the first $n - 1$ observations have freedom to be whatever value they wish, but the $n^{\text{th}}$ value has no freedom.

It must be whatever value forces the sum of the deviations about the mean to equal zero.

Why do we divide by the degrees of freedom in the sample standard deviation?

This has to do with the idea of "biased" vs "unbiased" statistic. The standard deviation is biased if we divide by n in a sample, so to correct for this, we divide by n-1 which "unbiases" the sample standard deviation.

Pearson

## EXAMPLE Computing a Sample Standard Deviation

Here are the results of a random sample taken from the travel times (in minutes) to work for all seven employees of a start-up web development company:

5, 26, 36

$n = 3$

Find the sample standard deviation.

*use stats law of rounding*

*calc enter into list
1 VAR Stats
"Sx" sample st.dev*

$S = 15.8$ min

**P** Pearson

| $x_i$ | $\overline{x}$ | $x_i - \overline{x}$ | $(x_i - \overline{x})^2$ |
|---|---|---|---|
| 5 | 22.33333 | −17.333 | 300.432889 |
| 26 | 22.33333 | 3.667 | 13.446889 |
| 36 | 22.33333 | 13.667 | 186.786889 |

$$\Sigma(x_i - \overline{x})^2 = 500.66667$$

$$s = \sqrt{\frac{\Sigma(x_i - \overline{x})^2}{n-1}} = \sqrt{\frac{500.66667}{2}} \approx 15.82 \text{ minutes}$$

**IN CLASS ACTIVITY**

**The Sample Standard Deviation**

Using the pulse data we collected from Section 3.1, do the following:

a) Obtain a simple random sample of $n = 4$ students and compute the sample standard deviation.

b) Obtain a second simple random sample of $n = 4$ students and compute the sample standard deviation.

c) Are the sample standard deviations the same? Why?

Pearson

**EXAMPLE Comparing Standard Deviations**

Determine the standard deviation waiting time for Wendy's and McDonald's. Which is larger? Why?

(Use Calc)

Pearson

# 3.2 Measures of Dispersion

*σ* — spreed small/values close together

*from 0* — more spread

## Wait Time at Wendy's

| | | | | | |
|------|------|------|------|------|------|
| 1.50 | 0.79 | 1.01 | 1.66 | 0.94 | 0.67 |
| 2.53 | 1.20 | 1.46 | 0.89 | 0.95 | 0.90 |
| 1.88 | 2.94 | 1.40 | 1.33 | 1.20 | 0.84 |
| 3.99 | 1.90 | 1.00 | 1.54 | 0.99 | 0.35 |
| 0.90 | 1.23 | 0.92 | 1.09 | 1.72 | 2.00 |

$S = 0.738$

## Wait Time at McDonald's

data is more spread out!

larger!

| | | | | | |
|------|------|------|------|------|------|
| 3.50 | 0.00 | 0.38 | 0.43 | 1.82 | 3.04 |
| 0.00 | 0.26 | 0.14 | 0.60 | 2.33 | 2.54 |
| 1.97 | 0.71 | 2.22 | 4.54 | 0.80 | 0.50 |
| 0.00 | 0.28 | 0.44 | 1.38 | 0.92 | 1.17 |
| 3.08 | 2.75 | 0.36 | 3.10 | 2.19 | 0.23 |

$S = 1.265$

Pearson

**EXAMPLE Comparing Standard Deviations**

Sample standard deviation for Wendy's:

<div align="center">0.738 minutes</div>

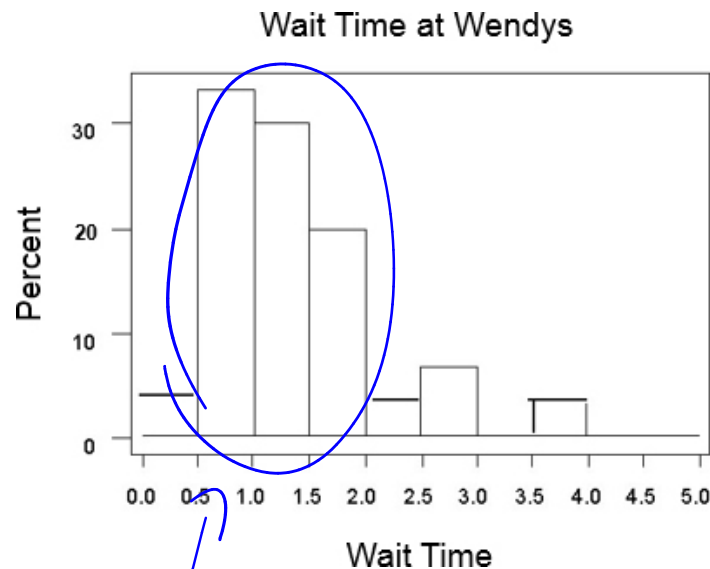Sample standard deviation for McDonald's:

<div align="center">1.265 minutes</div>

Recall from earlier that the data is more dispersed for McDonald's resulting in a larger standard deviation.
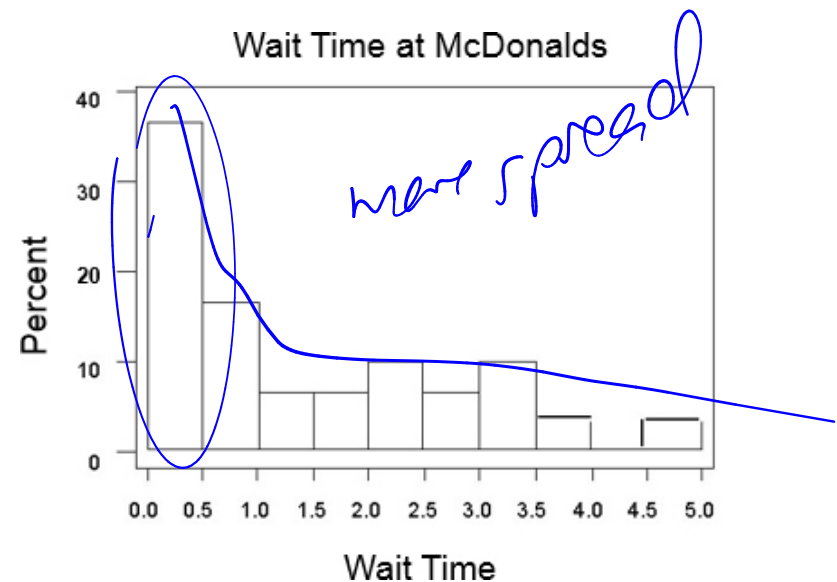
Pearson

The mean wait time in each line is 1.39 minutes.

Histograms for wait time data.



Wait Time at Wendys

close together

Wait Time at McDonalds

more spread

Pearson

# 3.2 Measures of Dispersion

The **variance** of a variable is the square of the standard deviation.

$$\sigma^2 = \frac{\sum(x_i - \mu)^2}{N} \qquad \text{vs}$$

"possible
frey Q"

## NOTATION:

- The **population variance** is $\sigma^2$

- The **sample variance** is $s^2$

## Units:

- The units of population variance are units squared!

- The units of sample variance are units squared!

Pearson

**EXAMPLE Computing a Population Variance**

The following data represent the travel times (in minutes) to work for all seven employees of a start-up web development company.

23, 36, 23, 18, 5, 26, 43

Compute the population and sample variance of this data.

**Note: calculator doesn't compute variance. You need to compute it from the standard deviation by squaring it.**

Pearson

# 3.2 Measures of Dispersion

## 3.2.3 Determine the Variance of a Variable from Raw Data (3 of 3)

**EXAMPLE Computing a Population Variance**

Recall that the population standard deviation (from previous slide) is $\sigma = 11.36$

so the population variance is…

$$\sigma^2 = 11.36^2 = 129.0496 = \boxed{129.1 \text{ minutes squared}}$$

Recall that the sample standard deviation is $s = 15.82$,
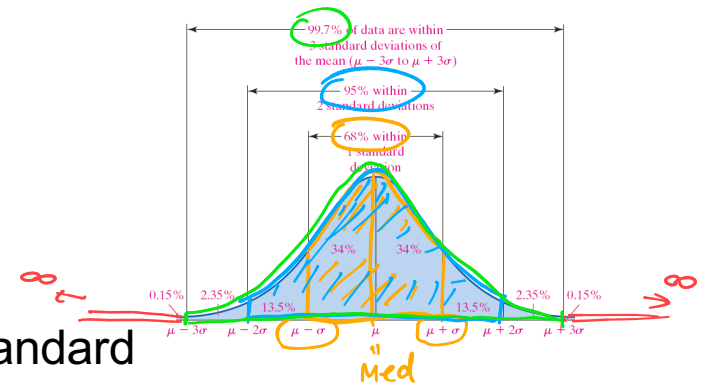
so the sample variance is…

Pearson

68-95-99.7% Rule

**The Empirical Rule**

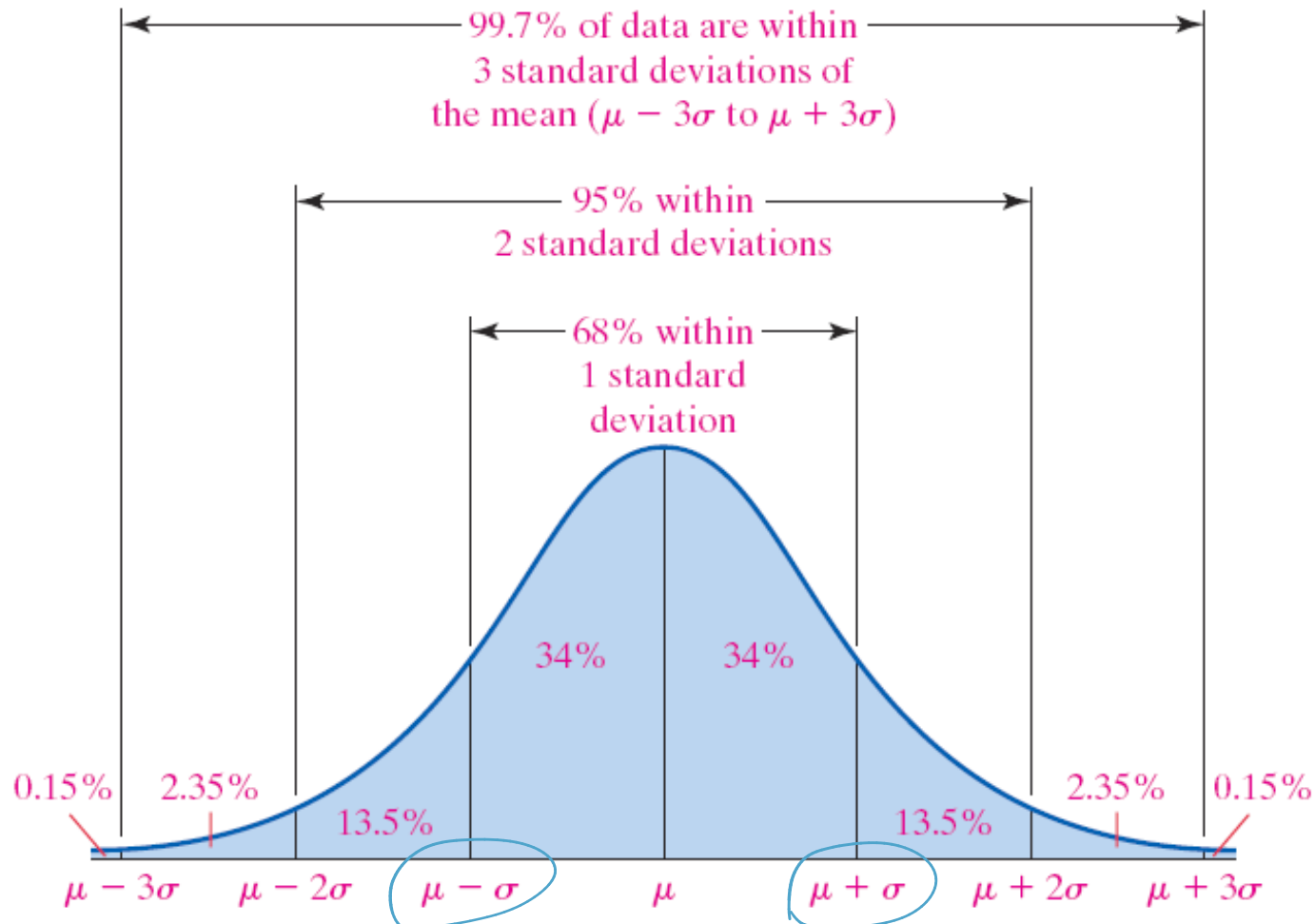If a distribution is roughly bell shaped, then



- **Approximately 68%** of the data will lie <u>within 1</u> standard deviation of the mean. That is, approximately 68% of the data lie between ___$\mu - \sigma$___ and ___$\mu + \sigma$___

- **Approximately 95%** of the data will lie <u>within 2</u> standard deviations of the mean. That is, approximately 95% of the data lie between ___$\mu - 2\sigma$___ and ___$\mu + 2\sigma$___

- **Approximately 99.7%** of the data will lie <u>within 3</u> standard deviations of the mean. So, approx. 99.7% of the data lie between ___$\mu - 3\sigma$___ and ___$\mu + 3\sigma$___

**Note**: We can also use the Empirical Rule based on sample data with $\bar{x}$ used in place of $\mu$ and $s$ used in place of $\sigma$.

P Pearson

99.7% of data are within
3 standard deviations of
the mean ($\mu - 3\sigma$ to $\mu + 3\sigma$)

95% within
2 standard deviations

68% within
1 standard
deviation

34%    34%

0.15%    2.35%    13.5%    13.5%    2.35%    0.15%

$\mu - 3\sigma$    $\mu - 2\sigma$    $\mu - \sigma$    $\mu$    $\mu + \sigma$    $\mu + 2\sigma$    $\mu + 3\sigma$

Pearson

# 3.2 Measures of Dispersion

$\mu = 150$

$\sigma = 8$

Use the Empirical Rule to fill out the normal distribution:

mean of 150 and standard deviation of 8

99.7%

68%

95%

126    134    142    150    158    166    174

Pearson

*population*

## EXAMPLE Using the Empirical Rule

The following data represent the serum HDL cholesterol of the 54 female patients of a family doctor.

| 41 | 48 | 43 | 38 | 35 | 37 | 44 | 44 | 44 |
|----|----|----|----|----|----|----|----|----|
| 62 | 75 | 77 | 58 | 82 | 39 | 85 | 55 | 54 |
| 67 | 69 | 69 | 70 | 65 | 72 | 74 | 74 | 74 |
| 60 | 60 | 60 | 61 | 62 | 63 | 64 | 64 | 64 |
| 54 | 54 | 55 | 56 | 56 | 56 | 57 | 58 | 59 |
| 45 | 47 | 47 | 48 | 48 | 50 | 52 | 52 | 53 |

The population mean and the standard deviation is:
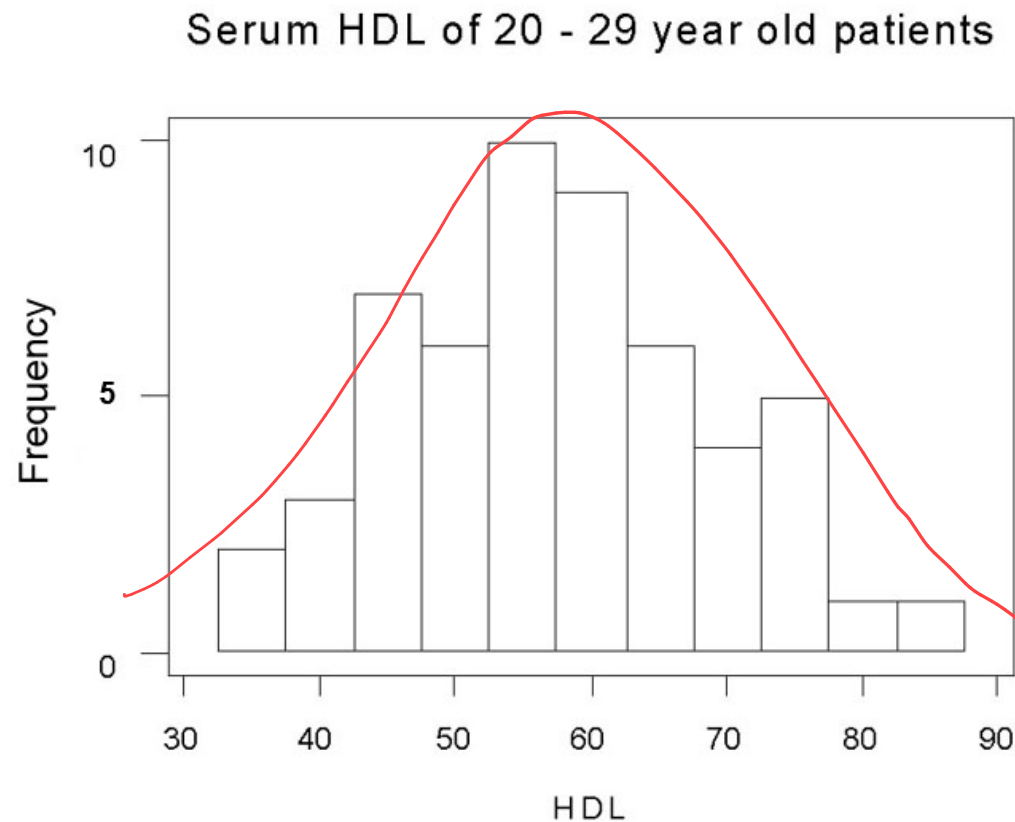
$$\mu = 57.4 \text{ and } \sigma = 11.8$$

Pearson

a) Draw a histogram to verify the data is approximately bell-shaped.

b) Determine the percentage of all patients that have serum HDL within 3 standard deviations of the mean according to the Empirical Rule.   *Ans   99.7%*

c) Determine the percentage of all patients that have serum HDL between 34 and 69.1 according to the Empirical Rule.
*Ans   68 + 13.5 = 81.5%*

d) Determine the actual percentage of patients that have serum HDL between 34 and 69.1.   *#*

*54  * 0.815  =  44.01  ⟶  45 female patients have HDL levels between 34 & 69.1*

*patients*

P Pearson

# 3.2 Measures of Dispersion

## 3.2.4 Use the Empirical Rule to Describe Data that are Bell Shaped (6 of 7)



Serum HDL of 20 - 29 year old patients

*approximately bell shaped*

*can use Emp. Rule.*

Pearson

# 3.2 Measures of Dispersion

Use the Empirical Rule to fill out the normal distribution:

$$\mu = 57.4$$
$$\sigma = 11.8$$

$$95 - 68 = 27$$
$$27/2 = 13.5$$

13.5%

68%

34

69.1

22    33.8    45.6    57.4    69.2    81    92.8

Pearson

(c) According to the Empirical Rule, 99.7% of the all patients that have serum HDL within 3 standard deviations of the mean.

(d) 13.5% + 34% + 34% = 81.5% of all patients will have a serum HDL between 34.0 and 69.1 according to the Empirical Rule.

(e) 45 out of the 54 or 83.3% of the patients have a serum HDL between 34.0 and 69.1.

## Rounding Rules:

- Stats Law of Rounding
  - When rounding statistics based on data, round (final answers) to one more significant figure than the original data.
  - Example: mean, median, standard deviation, etc

- Miscellaneous
  - When rounding people, always round UP
    - Ex: if we estimate that 23.2 people, then round up to 24 people