

DATA: List of n numbers/measurements

$$x_1, x_2, x_3, \dots, x_n.$$

EXAMPLE: Table 3 Entrance Examination Scores of 100 Entering Freshmen

762	451	602	440	570	553	367	520	454	653
433	508	520	603	532	673	480	592	565	662
712	415	595	580	643	542	470	743	608	503
566	493	635	780	537	622	463	613	502	577
618	581	644	605	588	695	517	537	552	682
340	537	370	745	605	673	487	412	613	470
548	627	576	637	787	507	566	628	676	750
442	591	735	523	518	612	589	648	662	512
663	588	627	584	672	533	738	455	512	622
544	462	730	576	588	705	695	541	537	563

Our goal is to extract useful information from the data.

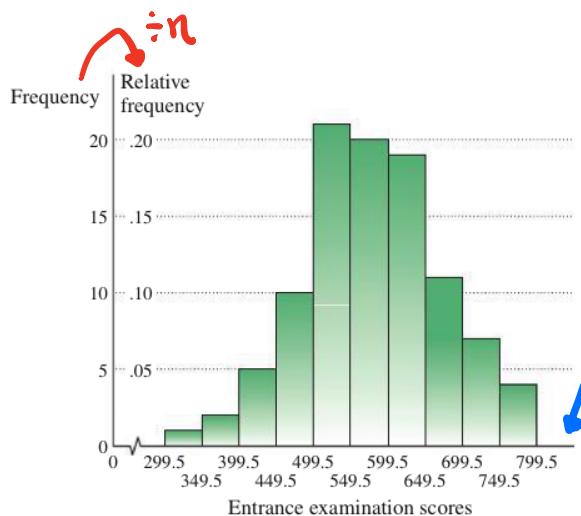
(Describe / summarize)

Two methods:

I VISUALLY

II NUMERICALLY

I. VISUALLY : FREQUENCY / RELATIVE FREQUENCY HISTOGRAM



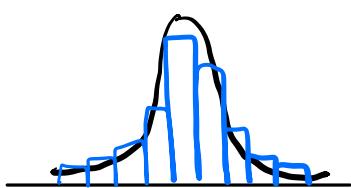
BAR CHARTS

Histograms

HORIZONTAL POTS IS # LINE
(CLASSES = INTERVALS)

Now we can see the distribution of measurements.

WORDS TO DESCRIBE DISTRIBUTIONS:



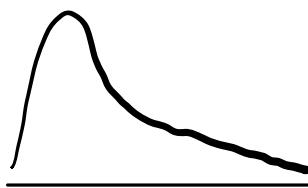
SYMMETRIC

e.g. SAT SCORES

HEIGHTS OF ADULT MALES

WEIGHTS OF APPLES

PRODUCED BY 1 TREE

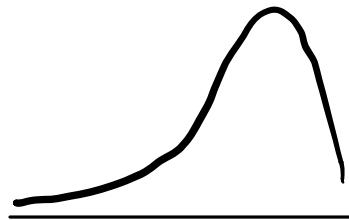


RIGHT-SKewed

e.g. HOUSEHOLD INCOME

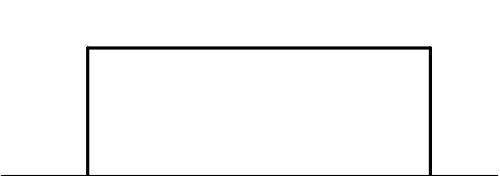
NUMBER OF SIBLINGS

SIZE OF NYC APARTMENTS



LEFT-SKewed

e.g. LIFETIME OF HUMANS



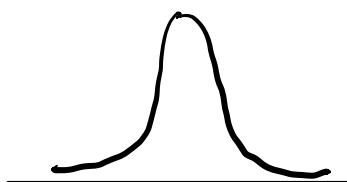
UNIFORM

e.g. # TIMES EACH FACE APPEARS ON 1000 ROLLS OF A DIE.

TIMES EACH DIGIT APPEARS IN FIRST 10 MILLION DIGITS OF π

"NORMAL #'S"

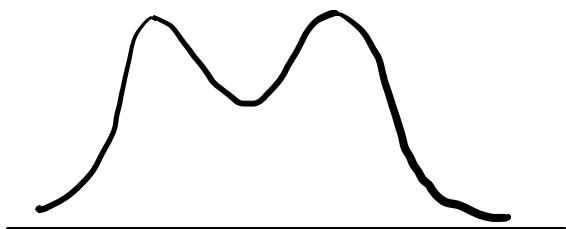
e



UNIMODAL

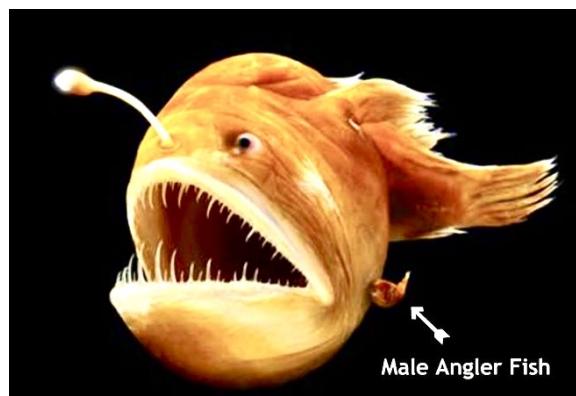
e.g. RESTING HEART RATE FOR HUMANS

PSI AT WHICH A PARTICULAR TYPE OF BICYCLE INNER TUBE POPS



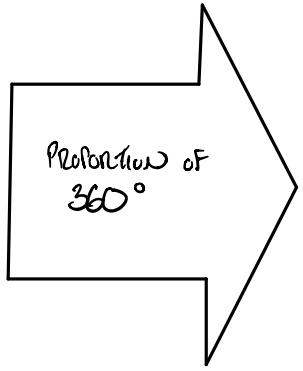
BIMODAL

e.g. SIZE OF ANGLER FISH



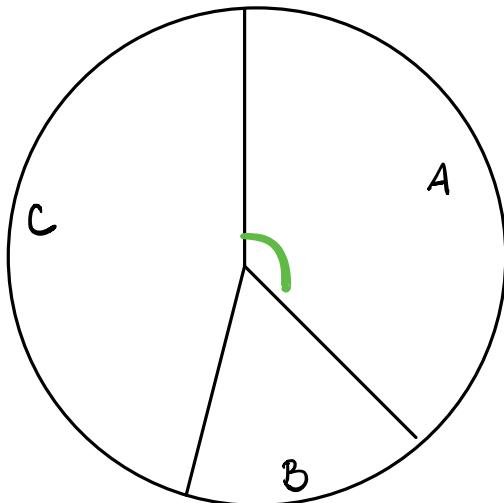
Pie Charts

	Proportion	PERCENTS
Type A	$\frac{n_1}{N}$	35%
Type B	$\frac{n_2}{N}$	18%
Type C	$\frac{n_3}{N}$	47%



CENTRAL ANGLE

$$.35 \times 360 = 126^\circ$$
$$.18 \times 360 = 64.8^\circ$$
$$.47 \times 360 = 169.2^\circ$$



CENTRAL ANGLE =
Proportion \times 360°

DATA VISUALIZATION

STATISTICS
GRAPHIC DESIGN

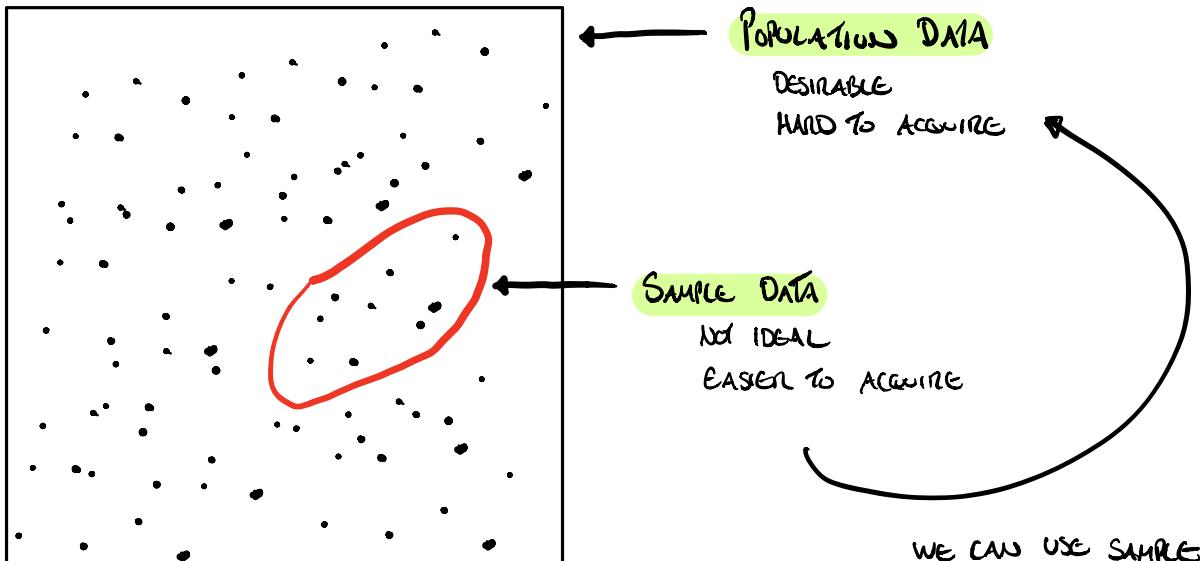
REDDIT.COM/r/DATAISBEAUTIFUL

II. NUMERICALLY : MEASURES OF CENTER & VARIATION

DATA: $x_1, x_2, x_3, \dots, x_n$

2 SOURCES

We Want to
Study
The Population



WE CAN USE SAMPLE
DATA TO MAKE ESTIMATES
FOR THE POPULATION
(INFERENTIAL STATISTICS)

3 MEASURES OF CENTER

"X-Bar"

1. SAMPLE MEAN \bar{x}

Population Mean μ

$$\left\{ \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{x_1 + x_2 + x_3 + \dots + x_n}{n} \right.$$

↑
SIGMA NOTATION

Σ GREEKS
from Sum

END VALUE FOR INDEX

INDEX

INDEX INCREASES BY 1

e.g. $\sum_{i=1}^5 \frac{1}{2^i} = \frac{1}{2^1} + \frac{1}{2^2} + \frac{1}{2^3} + \frac{1}{2^4} + \frac{1}{2^5} = \frac{63}{64}$

START VALUE FOR INDEX

TEMPLATE FOR OBJECTS/TERMS TO BE ADDED

$$\sum_{i=2}^5 2i(i+1)^2$$

$$2 \cdot 2(2+1)^2 + 2 \cdot 3(3+1)^2 + 2 \cdot 4(4+1)^2 + 2 \cdot 5(5+1)^2$$

$$4 \times 9 + 6 \times 16 + 8 \times 25 + 10 \times 36$$

The mean of 4 numbers is 90. If the mean of the first three numbers is 88, find the fourth number.

$$\text{Wzano: } \frac{88 + 92}{2} = 90 \Rightarrow x_4 = 92$$

$$\frac{x_1 + x_2 + x_3 + x_4}{4} = 90$$

$$x_1 + x_2 + x_3 + x_4 = 4 \cdot 90 = 360$$

$$\underbrace{\qquad\qquad\qquad}_{\text{Sum}}$$

$$\frac{x_1 + x_2 + x_3}{3} = 88 \quad \text{Ave}$$

$$x_1 + x_2 + x_3 = 3 \cdot 88 = 264$$

$$264 + x_4 = 360$$

$$x_4 = 96$$

Example. Suppose I buy 20 gallons of gas at an average price of \$2.40/gallon, and you buy 10 gallons of gas at an average price of \$2.10/gallon. Together, what is the average price per gallon that we've paid for gas?

WEIGHTED AVERAGES:

GIVEN n 'S THE AVERAGE (MEAN) IS

$$\frac{x_1 + x_2 + x_3 + \dots + x_n}{n} = \frac{x_1}{n} + \frac{x_2}{n} + \frac{x_3}{n} + \dots + \frac{x_n}{n}$$

$$= \frac{1}{n}x_1 + \frac{1}{n}x_2 + \frac{1}{n}x_3 + \dots + \frac{1}{n}x_n$$

ADD UP TO 1. ALL THE SAME.

More Generally, A WEIGHTED AVERAGE OF n 'S IS

$$\alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_3 + \dots + \alpha_n x_n$$

ADD UP TO 1. ALL POSITIVE.

Example. Suppose you have a homework grade of 90, and quiz grade of 85, and an exam grade of 80. If homework counts for 20% of your grade, the quiz counts for 35% of your grade, and the exam counts for 45% of your grade, calculate your average for the class.

$$.20 \times 90 + .35 \times 85 + .45 \times 80 = 83.75$$

Example. It costs a shipping company \$8.75 to ship a small package overnight. Suppose the shipping company charges a flat rate of \$19 to ship a small package overnight, and if the package is late they refund the full amount. If 93% of all packages are delivered on time, what is the average profit per package that the shipping company earns?

$$2 \text{ TYPES OF TRANSACTIONS} \quad \text{ON-TIME : } \text{Profit} = +19 - 8.75 = 10.25$$

$$\text{LATE : } \text{Profit} = +19 - 8.75 - 19 = -8.75$$

$$100 \text{ PACKAGES : } \frac{93(10.25) + 7(-8.75)}{100}$$



$$.93(10.25) + .07(-8.75) = \boxed{\$8.92}$$

2. MEDIAN

Arrange the data in order from least to greatest.

If the number of measurements is odd, then the median is the number in the middle position.

e.g. 1 1 3 3 4
 ↑
 MEDIAN

If the number of measurements is even, then the median is the mean of the two numbers that share/straddle the middle position.

e.g. 1 1 3 5 5 6
 $\frac{3+5}{2} = 4$
 MEDIAN

Note: An equal number of measurements lie to the right/left of the median when arranged in order.

"HALF THE DATA IS LESS THAN THE MEDIAN, & HALF THE DATA IS GREATER THAN THE MEDIAN."

Good when you don't want the "center" to be influenced by extreme values

35,000 42,000 46,000 48,000 51,000 63,000

MEAN 48,833

MEDIAN 47,000 (HALF ABOVE, HALF BELOW)

35,000 42,000 46,000 48,000 51,000 63,000,000

MEAN 10,538,333

MEDIAN 47,000 (HALF ABOVE, HALF BELOW)

Blood cholesterol levels. Find the mean and median for the data in the following table.

Midpoint	Interval	Frequency
159.5	149.5–169.5	4
179.5	169.5–189.5	11
199.5	189.5–209.5	15
219.5	209.5–229.5	25
239.5	229.5–249.5	13
259.5	249.5–269.5	7
279.5	269.5–289.5	3
299.5	289.5–309.5	2

(TOTAL 80)

WE ASSUME ALL MEASUREMENTS IN EACH CLASS/BIN ARE EQUAL TO THE MIDPOINT OF THAT CLASS/BIN'S INTERVAL

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \rightarrow \frac{1}{n} \sum_{k=1}^m m_k f_k$$

MIDPOINT OF EACH CLASS
FREQUENCY OF k^{th} CLASS

$$= \frac{1}{80} \left(4(159.5) + 11(179.5) + 15(199.5) + 25(219.5) + 13(239.5) + 7(259.5) + 3(279.5) + 2(299.5) \right)$$

$$= \frac{4}{80}(159.5) + \frac{11}{80}(179.5) + \frac{15}{80}(199.5) + \frac{25}{80}(219.5) + \frac{13}{80}(239.5) + \frac{7}{80}(259.5) + \frac{3}{80}(279.5) + \frac{2}{80}(299.5)$$

{

WEIGHTED AVERAGE OF MIDPOINTS

WEIGHTED BY RELATIVE FREQUENCY

$$\frac{\text{FREQUENCY}}{n}$$

3. Mode : Most frequently occurring measurement(s).

Yes, there may be multiple modes (ties allowed).

Blood cholesterol levels. Find the mean and median for the data in the following table.

Blood Cholesterol Levels
(milligrams per deciliter)

MODAL CLASS [209.5, 229.5]

MIDPOINT	Interval	Frequency
159.5	149.5–169.5	4
179.5	169.5–189.5	11
199.5	189.5–209.5	15
219.5	209.5–229.5	25
239.5	229.5–249.5	13
259.5	249.5–269.5	7
279.5	269.5–289.5	3
299.5	289.5–309.5	2

MODE 219.5

Blood cholesterol levels. Find the mean and median for the data in the following table.

Blood Cholesterol Levels
(milligrams per deciliter)

MIDPOINT	Interval	Frequency
159.5	149.5–169.5	4
179.5	169.5–189.5	11
199.5	189.5–209.5	15
219.5	209.5–229.5	25
239.5	229.5–249.5	13
259.5	249.5–269.5	7
279.5	269.5–289.5	3
299.5	289.5–309.5	2

MEDIAN

$$\overbrace{x_1, x_2, \dots, x_{39}}^{39}, x_{40}, x_{41}, \dots, \overbrace{x_{79}, x_{80}}^{39}$$

$$\text{MEDIAN} = \frac{x_{40} + x_{41}}{2}$$

$$= \frac{219.5 + 219.5}{2}$$

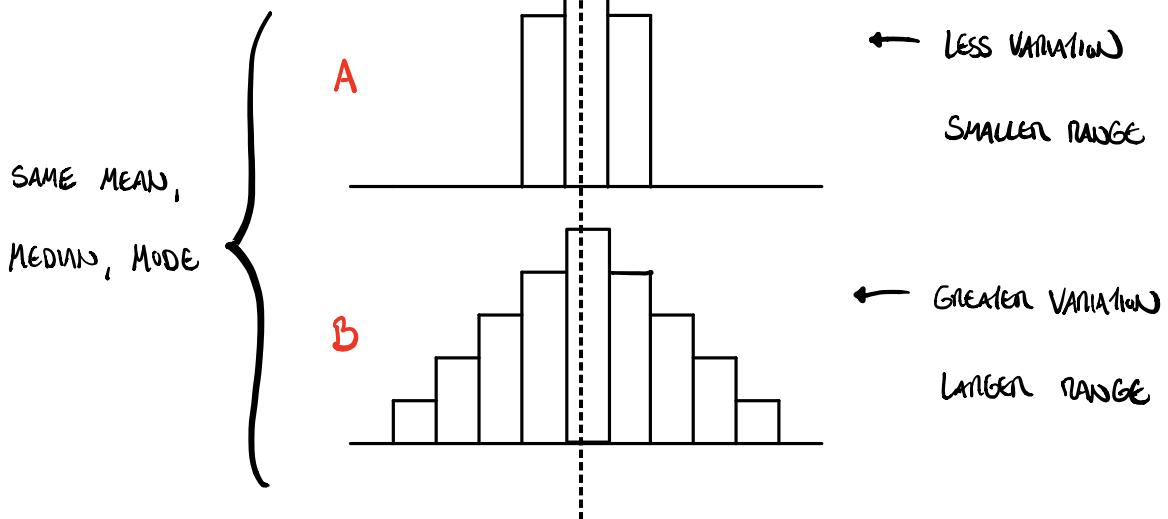
$$= 219.5$$

x_{40}, x_{41} are in this interval

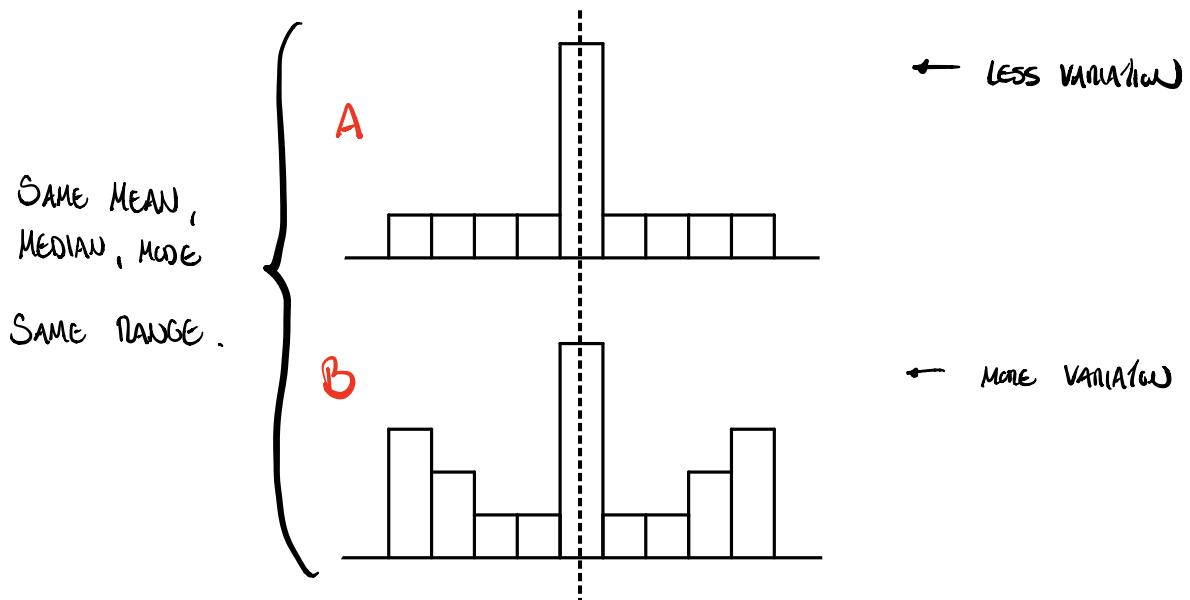
3 MEASURES OF VARIATION

DISTRIBUTION OF SIZE OF FISH IN LAKES A & B

$$1. \text{ RANGE} = \text{MAX} - \text{MIN}$$



WHAT ABOUT THE FOLLOWING DISTRIBUTIONS:



WE CAN MEASURE THE DIFFERENCE IN VARIATION HERE WITH
VARIANCE & STANDARD DEVIATION.

Variance σ^2 & Standard Deviation σ For Population

DATA: $x_1, x_2, x_3, \dots, x_n$

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

$$\text{STANDARD DEVIATION } \sigma = \sqrt{\sigma^2} \quad (\text{SQR OF VARIANCE})$$

Variance s^2 & Standard Deviation s For Sample

DATA: $x_1, x_2, x_3, \dots, x_n$

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

$$\text{STANDARD DEVIATION } s = \sqrt{s^2} \quad (\text{SQR OF VARIANCE})$$

s^2 & s ARE
GOOD ESTIMATES FOR
 σ^2 & σ

SMALLER DENOMINATOR \Rightarrow LARGER NUMBER

$$\frac{1}{2} > \frac{1}{3}$$

e.g. DATA : 13 14 17 25 26 (mean = 19)

CALCULATE THE VARIANCE & STANDARD DEVIATION

ASSUMING THE DATA COMES FROM A (a) POPULATION

(b) SAMPLE

(c) HOW MANY MEASUREMENTS LIE WITHIN 1 STD. DEV. OF MEAN?

(d) HOW MANY MEASUREMENTS LIE WITHIN 2 STD. DEV. OF MEAN?

x_i	$x_i - \bar{x}$	$(x_i - \bar{x})^2$
13	$13 - 19 = -6$	36
14	$14 - 19 = -5$	25
17	$17 - 19 = -2$	4
25	$25 - 19 = 6$	36
26	$26 - 19 = 7$	49

$$\begin{aligned}
 (a) \sigma^2 &= \frac{1}{5} \sum (x_i - \bar{x})^2 \\
 &= \frac{1}{5} (36 + 25 + 4 + 36 + 49) \\
 &= \frac{1}{5} (150) = 30
 \end{aligned}$$

$$\begin{aligned}
 (b) s^2 &= \frac{1}{5-1} \sum (x_i - \bar{x})^2 \\
 &= \frac{1}{4} (36 + 25 + 4 + 36 + 49) \\
 &= \frac{1}{4} (150) = 37.5
 \end{aligned}$$

$$s = \sqrt{s^2} = \sqrt{30}$$

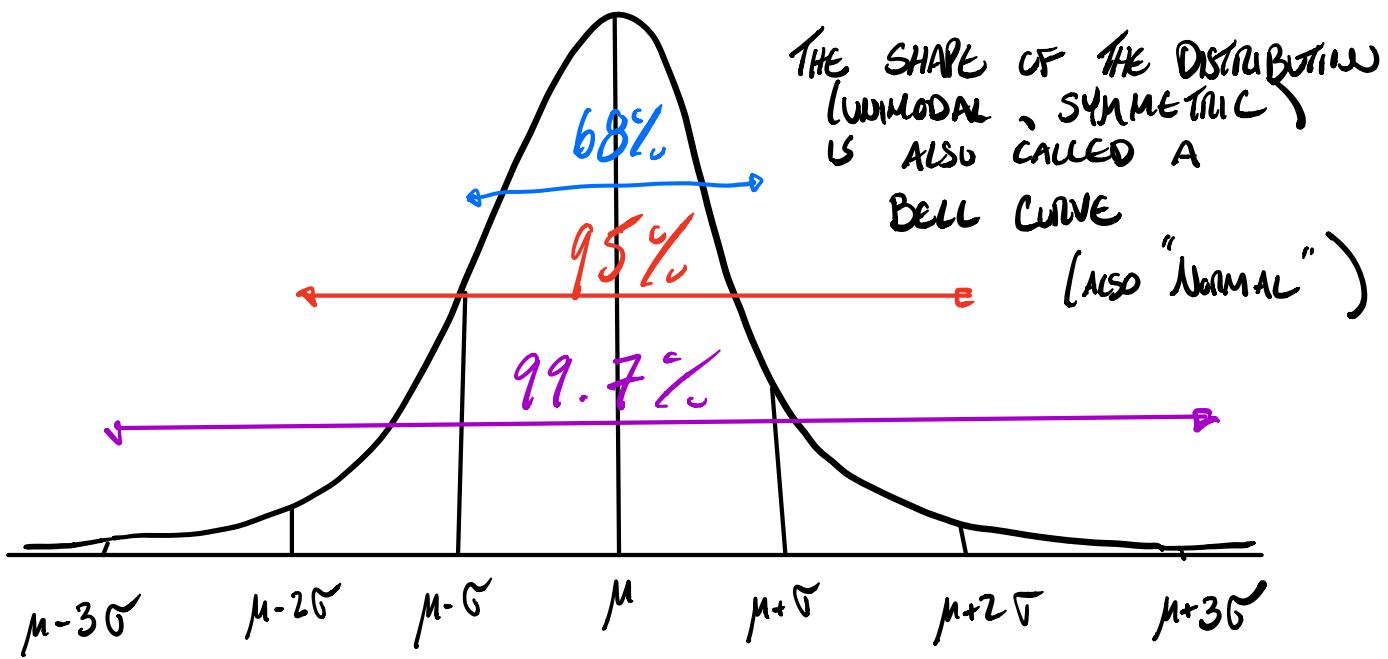
$$s = \sqrt{s^2} = \sqrt{37.5}$$

THE EMPIRICAL RULE :

APPROX. 68% OF DATA LIES WITHIN 1 STAND. DEV. OF MEAN

APPROX. 95% OF DATA LIES WITHIN 2 STAND. DEV. OF MEAN

APPROX. 99.7% OF DATA LIES WITHIN 3 STAND. DEV. OF MEAN



In Problems 11 and 12, find the standard deviation for each set of grouped sample data using formula (5) on page 525.

11. Interval	Frequency
2	2
5	5
8	7
11	1

$$\bar{x} = 6.4$$

x_i	$x_i - \bar{x}$	$(x_i - \bar{x})^2$	
2	-4.4	19.36	$\times 2$
5	-1.4	1.96	$\times 5$
8	1.6	2.56	$\times 7$
11	4.6	21.16	$\times 1$

$$\sigma^2 = \frac{2(19.36) + 5(1.96) + 7(2.56) + 1(21.16)}{15}$$

$$\sigma = \sqrt{\frac{2(19.36) + 5(1.96) + 7(2.56) + 1(21.16)}{15}}$$