

# Chapter 3: Describing, Exploring, and Comparing Data

## Section 3.1: Measures of Center

Def **Measure of Center**: a value at the center or middle of a data set.

*will define these more precisely later*

The three most widely-used measures of center are the mean, median, and mode.

The <b>(arithmetic) mean</b> of a data set is computed by <u>adding</u> all of the values of the variable in the data set and <u>dividing</u> by the number of observations.	
The <b>population arithmetic mean</b> ( $\mu$ ) is computed using <u>ALL</u> of the individuals in a population. The population mean is a <u>parameter</u> .	The <b>sample arithmetic mean</b> ( $\bar{x}$ ) is computed by using some of the individuals in a population. The sample mean is a <u>statistic</u> .
$\mu = \frac{x_1 + x_2 + \dots + x_N}{N} = \frac{\sum x}{N}$	$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{\sum x}{n}$

*pop. mean  
↓  
parameter*

*sample mean  
x̄ → stat*

*Σ = "sum"*

Ex: Of the 42 students enrolled in an Introductory Statistics course, the data below are the first 10 exam scores. Treat the 10 students as a sample of the population, which means you use  $\bar{x}$  to find the mean.

Student	Score
Michelle	82
Ryanne	77
Bilal	90
Pam	71
Jennifer	62
Dave	68
Joel	74
Sam	84
Justine	94
Juan	88

*n = 10*

$$\bar{x} = \frac{\sum x}{n} = \frac{82 + 77 + 90 + \dots + 88}{10} \quad (n=10)$$

$$= (82 + 77 + \dots + 88) / 10 = 79$$

*calc  
2nd estimate*

The <b>median</b> of a data set is the value that lies in the <u>middle</u> of the data when arranged in <u>ascending</u> order. We use <u>M</u> to represent the median.	
<p><b>ODD</b> number of data</p> <p>1<sup>st</sup>: Arrange the data in <u>ascending</u> order</p> <p>2<sup>nd</sup>: The median will be the <u>middle number</u></p> <p>Ex: 11, 14, <u>16</u>, 19, 28 <i>already ascending</i></p> <p><i>M = 16</i></p>	<p><b>EVEN</b> number of data</p> <p>1<sup>st</sup>: Arrange the data in <u>ascending</u> order</p> <p>2<sup>nd</sup>: The median will be the <u>mean</u> of the middle numbers</p> <p>Ex: <u>14</u>, 18, <u>20</u>, 26, 31, <u>39</u> ← max</p> <p><i>min</i>      <i>M = (20+26)/2 = 23</i></p>

The <b>midrange</b> of a data set is the value midway between the minimum and maximum values.
$\text{Midrange} = \frac{\text{min value} + \text{max value}}{2}$ <p><i>Midrange = (14+39)/2 = 26.5</i></p>

Ex: Use the data from the Introductory Statistics example from above to find the median and midrange.

*Use calculator! • Enter data: (STAT → edit)  
Sort: • Sort: (STAT → SORTA(L1))*

*sorted: Median: M = (77+82)/2 = 79.5*

*62 68 71 74 77 82 84 88 90 94*

*n=10 even "middle two" = 10/2 = 5  
" + 1" = 6*

**Round-Off Rule:** Carry one more decimal place than is present in the original set of values.

*"STAT RULE OF ROUNDING"*

*Ex: Data whole #s, rounding rule is to one decimal place.*

*Midrange: (62+94)/2 = 78*

The **mode** of a variable is the **most frequent** observation of the variable that occurs in the data set.  
 \*If no observation occurs more than once, we say that the data have no mode (or mode is N/A)  
 \*If the data set has more than one observation that repeat the same number of times, then it is considered

Ex: Find the mode for each example below.

a) The following data represent the number of O-ring failures on the shuttle <i>Columbia</i> for its 17 flights prior to its fatal flight: 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 2, 3 mode = 0	b) The data of the test scores from above: 82, 77, 90, 71, 62, 68, 74, 84, 94, 88 no mode	c) Hair color of ten people in line: Brown, Blonde, Red, Brown, Brown, Blonde, Brown, Blonde, Blonde, Red qualitative Brown: 4 Blonde: 4 Mode: Brown & Blonde "bimodal"
---	---	---

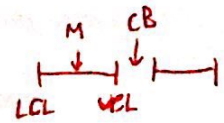
**Mean from a Frequency Distribution**

Formula:

$$\bar{x} = \frac{\sum (f \cdot x)}{n}$$

Ex: The following table gives the weights of a sample of 100 babies born at a local hospital.

Weight (in lbs)	Freq (f)	class midpt (x)	f · x
3-4.9	5	$\frac{3+4.9}{2} = 3.95$	$5 \cdot 3.95 = 19.75$
5-6.9	32	5.95	190.4
7-8.9	40	7.95	318
9-10.9	18	9.95	179.1
11-12.9	5	11.95	59.75
	$n = \sum f = 100$		$\sum (f \cdot x) = 767$



Find the sample mean.

$$\bar{x}$$

thus represents a lot of data!

$$\bar{x} = \frac{767}{100} = 7.67 \text{ lbs}$$

!important! Units!

**Resistance Statistics**

A numerical summary of data is said to be resistant if extreme values (very large or small) relative to the data do not affect its value substantially.

Ex: The following are wait times (in minutes) at a dentist office: 1, 1, 2, 2, 3, 5. (a) Find the mean and median.

$$\bar{x} = \frac{\sum x}{n} = \frac{14}{6} = 2.3 \quad M = \frac{2+2}{2} = 2$$

b) Note the value of 102 minutes added to this data. Find the mean and median. Which measure is resistant to the added value?

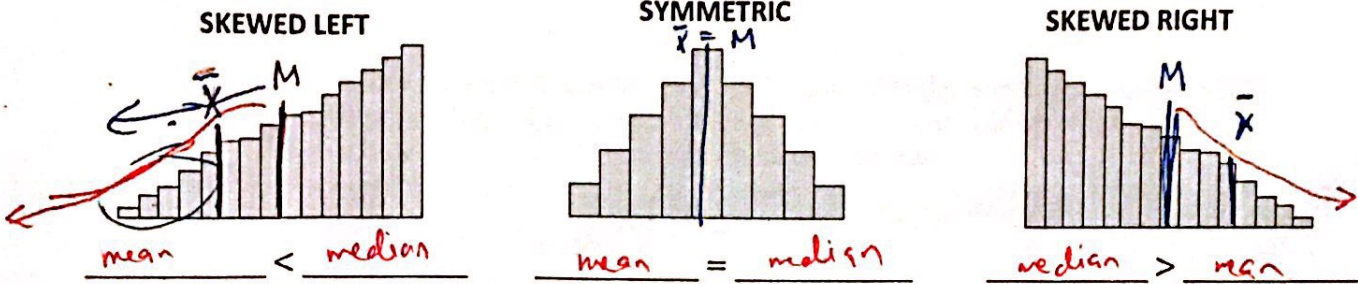
1, 1, 2, 2, 3, 5, 102

$$\bar{x} = \frac{116}{7} = 16.6 \quad M = 2$$

so mean changed from 2.3 to 16.6  
 Median is resistant!

Explanation of skewness: tail moves the mean

When data are skewed, there are extreme values in the tail, which tend to pull the mean in the direction of the tail.



General rule: If the data are symmetric use the MEAN as the best measure of center.

If the data are skewed use the MEDIAN as the best measure of center.

Ex: FICO scores range in value from 300 to 850, with a higher score indicating a more creditworthy individual. The distribution of FICO scores is skewed left with a median score of 723

(a) Do you think the mean FICO score is greater than, less than, or equal to 723? Justify your response.

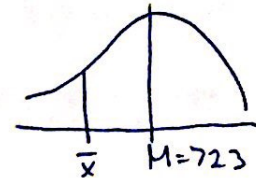
mean < 723

mean is less than 723

(b) What proportion of individuals have a FICO score above 723?

"fraction" since 723 is median = middle

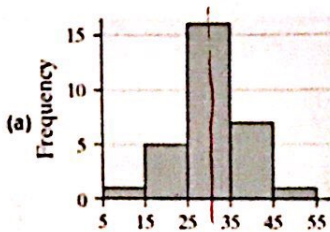
proportion above 723 is 1/2 or 0.5



Ex: Match the histograms shown to the appropriate summary statistics by writing the appropriate number under each histogram.

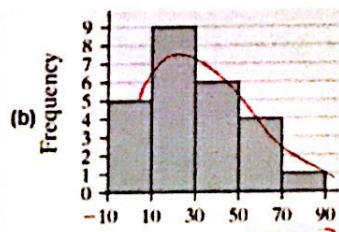
	Mean	Median
1	42	42
2	31	36
3	31	26
4	31	32

→  $\bar{x} = M$  → symmetric  
 →  $\bar{x} < M$  → skewed LEFT  
 →  $\bar{x} > M$  → skewed RIGHT  
 →  $\bar{x} \approx M$  → symmetric

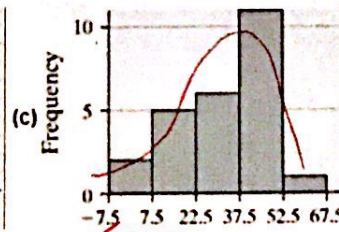


$\bar{x} = M \approx 30$

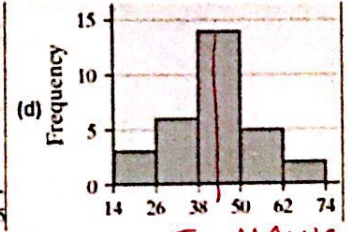
④



tail →  
 ③



tail ←  
 ②



$\bar{x} \approx M \approx 45$

①

### 3.2 Measures of Variation

#### Importance of Variation

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.

	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

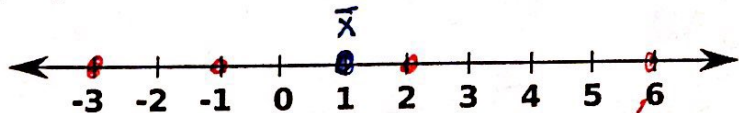
$\pm 2$  from  $\bar{x}$   
 $\pm 7$  from  $\bar{x}$

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

**Advil** Because all the values are much closer to 200 mg ( $\pm \approx 2$  mg from  $\bar{x}$ ) whereas Motrin IB has a much larger variation.... ( $\pm \approx 7$  mg from  $\bar{x}$ )

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

(a) Find the mean.  $\bar{x} = \frac{\sum x}{n} = \frac{4}{4} = 1$



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

use +/- to represent right/left side of  $\bar{x}$

add these up  
 $-4 - 2 + 1 + 5 = 0$   
not helpful!

#### Measures of variation

Def: The **range** of a data set is the difference between the maximum and minimum data values.

$$\text{range} = \text{maximum value} - \text{minimum value}$$

#### Standard Deviation of a Sample

**KEY**  $s$  - sample standard dev.  $\sigma$  - population st. dev.

Def: The **standard deviation** (denoted by  $s$ ) of a set of sample values is a measure of variation of values about the mean. It is a type of average deviation of values from the mean.

Formula: 
$$s = \sqrt{\frac{\sum (x - \bar{x})^2}{(n-1)}}$$

"Shortcut:"

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

Note: Each form can be tricky, but the alternative form tends to be easier.

#### Standard Deviation of a Population

Def: The **standard deviation** (denoted by  $\sigma$ ) of a complete set of values is a measure of variation of values about the mean. It is a type of average deviation of values from the mean.

Formula:

$$\sigma = \sqrt{\frac{\sum (x - \bar{x})^2}{N}}$$

*Note*

Note: It's rare to compute a population standard deviation. Therefore, when using technology, be sure to use the sample standard deviation unless otherwise noted.

**Variance**  $s^2$  or  $\sigma^2$

Def The variance (denoted by  $s^2$  or  $\sigma^2$ ) of a set of values is a measure of variation equal to the square of the standard deviation.

Ex: Find the range, standard deviation, and variance for the following sample of the number of chips in nine randomly sampled fun-sized bags of Doritos.

data → 25 31 28 19 24 26 29 32 20

• range  
32 - 19 = 13

range = 13

x	$x - \bar{x}$	$(x - \bar{x})^2$
19	19 - 26 = -7	(-7) <sup>2</sup> = 49
20	20 - 26 = -6	(-6) <sup>2</sup> = 36
24	24 - 26 = -2	(-2) <sup>2</sup> = 4
25	25 - 26 = -1	(-1) <sup>2</sup> = 1
26	26 - 26 = 0	(0) <sup>2</sup> = 0
28	28 - 26 = 2	(2) <sup>2</sup> = 4
29	29 - 26 = 3	(3) <sup>2</sup> = 9
31	31 - 26 = 5	(5) <sup>2</sup> = 25
32	32 - 26 = 6	(6) <sup>2</sup> = 36
$\Sigma x = 234$		$\Sigma (x - \bar{x})^2 = 164$

$$\bar{x} = \frac{\Sigma x}{n} = \frac{234}{9} = 26$$

• sample standard dev.

$$S = \sqrt{\frac{\Sigma (x - \bar{x})^2}{n - 1}} \quad \begin{matrix} n = 9 \\ n - 1 = 8 \end{matrix}$$

$$S = \sqrt{\frac{164}{8}} = \sqrt{20.5}$$

$S \approx 4.5$  "stat rule rounding"

Use the alternative form to find the standard deviation.

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

$$S = \sqrt{\frac{9(6248) - (234)^2}{9(9-1)}}$$

$$S = \sqrt{20.5}$$

$$S \approx 4.5$$

SAME!  
☺

x	$x^2$
19	361
20	400
24	576
25	625
26	676
28	784
29	841
31	961
32	1024
$\Sigma x = 234$	$\Sigma x^2 = 6248$

Question: If you bought a bag of chips everyday, would you prefer to have a small or large standard deviation between bags?

Small! so bag of chips can cant on the size of chips being consistent!

## Standard Deviation from a Frequency Distribution

Formula:

$$s = \sqrt{\frac{n[\sum(f \cdot x^2)] - [\sum(f \cdot x)]^2}{n(n-1)}}$$

Ex: Recall: The following table gives the weights of a sample of 100 babies born at a local hospital.

Weight (in lbs)	Freq (f)	class midpt (x)	f · x	f · x <sup>2</sup>
3-4.9	5	3.95	19.75	5 · (3.95) <sup>2</sup> = 78.0125
5-6.9	32	5.95	190.4	32 · (5.95) <sup>2</sup> = 1132.88
7-8.9	40	7.95	318	40 · (7.95) <sup>2</sup> = 2529.1
9-10.9	18	9.95	179.1	18 · (9.95) <sup>2</sup> = 1782.045
11-12.9	5	11.95	59.75	5 · (11.95) <sup>2</sup> = 714.025
	n = ∑f = 100		∑(f · x) = 767	∑(f · x <sup>2</sup> ) = 6235.05

Find the sample standard deviation and variance. →  $s^2 = (\sqrt{3.5571})^2 = 3.5571$  → ~~3.5571~~

st. dev

$$s = \sqrt{\frac{100[6235.05] - [767]^2}{100(99)}} = \sqrt{3.5571} \approx 1.8860$$

$$s = 1.89$$

↑ stats for rounding

## Range Rule of Thumb

Variance:  $s^2 = 3.56$

Data are significantly LOW if the value is  $\mu - 2\sigma$  or lower.

Data are not significant if the value is between  $\mu - 2\sigma$  and  $\mu + 2\sigma$ .

Data are significantly HIGH if the value is  $\mu + 2\sigma$  or higher.

Ex: The data below are free download wifi speeds (in Mbps) from ten of the busiest international airports.

AIRPORT CODE	WIFI SPEED (x)	AIRPORT CODE	WIFI SPEED (x)
DEN	78.2	YYC	41.8
YVR	55.1	BOS	32.2
PHL	48.4	DFW	32.0
SFO	45.3	MEX	27.7
SEA	43.7	DTW	22.9

$$\sum x = 427.3$$

$$\sum x^2 = 20555.37$$

<https://www.speedtest.net/insights/blog/fastest-airports-north-america-2017/>

If an international airport began to provide new 60.8 Mbps free wifi, and claimed that their speed is "miles above" many others, would you agree with their claim?

$$\bar{x} = \frac{\sum x}{n} = \frac{427.3}{10} = 42.73$$

check if significantly high if "60.8 >  $\bar{x} + 2s$ "

$$s = \sqrt{\frac{\sum(x - \bar{x})^2}{n-1}} = \sqrt{\frac{n(\sum x^2) - (\sum x)^2}{n(n-1)}} = \sqrt{\frac{10(20555.37) - (427.3)^2}{10 \cdot 9}} = 15.98 = s$$

$$\bar{x} + 2s = 42.73 + 2(15.98) = 74.69$$

Conclusion "60.8 is not significantly high!"

Remember: statistical significance is within 95%

don't forget thin tails 'go forever'

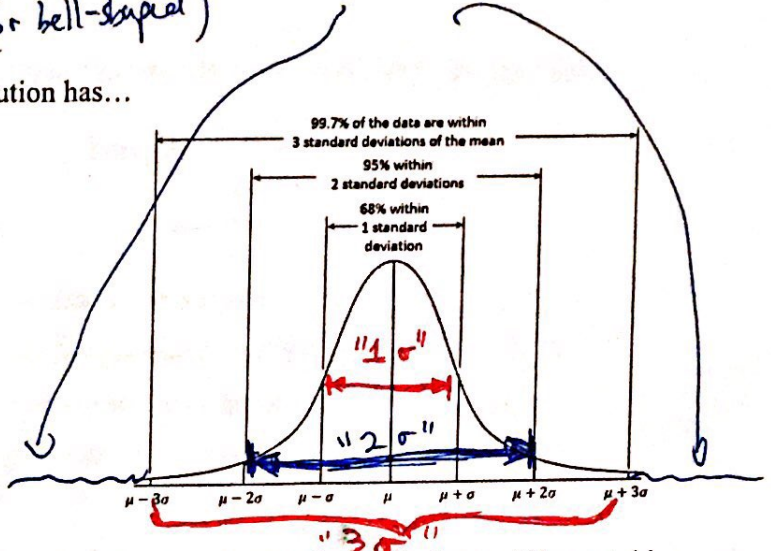
**"68-95-99.7 Rule"**

**Empirical Rule**

(or normal) (or bell-shaped)

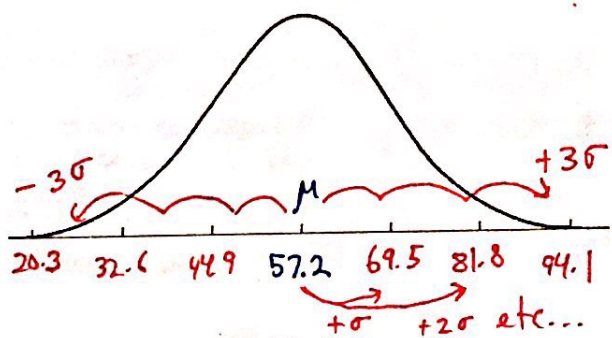
Empirical Rule says that a symmetric distribution has...

- Approximately 68 % of the data will lie within 1 standard deviation of the mean.
- Approximately 95 % of the data will lie within 2 standard deviation of the mean.
- Approximately 99.7 % of the data will lie within 3 standard deviation of the mean.



Ex: The following data represent the ages of all of the forty-two female patients of a family doctor. We are told that the data has a bell-shaped distribution and that the population mean,  $\mu$ , is 57.2 years and the population standard deviation,  $\sigma$ , is 12.3 years.

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



$\mu = 57.2$   
 $\sigma = 12.3$

(a) Determine the percentage of all patients whose age is within 3 standard deviations of the mean.

Approximately 99.7% by Empirical Rule

(b) Between what two values will this percentage lie?

Between 20.3 years and 94.1 years

units!!

**Finding summary statistics for a data set**

Use Graphing Calculator (TI-84 Plus)

Instructions: (a) STAT  $\Rightarrow$  1: Edit...  $\Rightarrow$  Enter data into a list  
(b) STAT  $\Rightarrow$  CALC  $\Rightarrow$  1: 1-Var Stats  $\Rightarrow$  List: (Choose list)  $\Rightarrow$  Calculate



### 3.3 Measures of Relative Standing and Boxplots

#### z Scores

Def: A **z score** is the number of standard deviations that a given value  $x$  is above or below the mean.

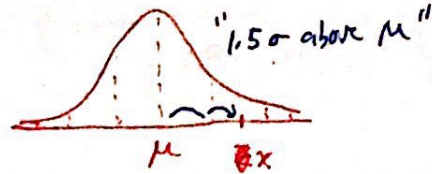
Formula:

Sample :

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

Population:

$$z = \frac{x - \mu}{\sigma}$$



**Round-Off Rule:** Round z scores to two decimal places.

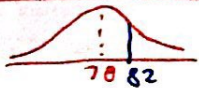
The z score is a standardized value that describes a data value's relative standing.

1. A **negative z score** corresponds to a data value **below the mean**.
2. **Unusual data values** are more than two standard deviations from the mean.
  - Ordinary values:  $-2 \leq z \text{ score} \leq 2$
  - **Unusual data values:**  $z < -2$  or  $z > 2$
3. The z score allows us to compare data values drawn from different samples or populations.

} purpose of a z-score!

Ex: Two college roommates are taking different physics courses at a university. They agree to a wager regarding their midterm scores, whereby the loser must do the dishes for a month. After scoring an **82**, Jacob insists that **Michael** lost since he earned a **70** on his exam. However, Michael argues that he performed better relative to the rest of his class than did Jacob. Use the given class results to determine who won the bet? → use z-score to compare

Jacob's class:  $\bar{x} = 78, s = 6$



Michael's class:  $\bar{x} = 55, s = 12$



more variation in Michael's class

$$z = \frac{x - \bar{x}}{s} = \frac{82 - 78}{6} \approx 0.67$$

$$z = \frac{x - \bar{x}}{s} = \frac{70 - 55}{12} = 1.25$$

#### Percentiles

**Michael scored better!**

b/c Michael's z-score was  $z = 1.25$  which is larger than Jacob's z-score of  $z = 0.67$

Def: **Percentiles (denoted  $P_k$ )** are measures of location in relation to all the other data values.

#### Notation

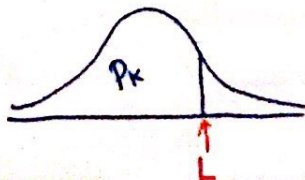
Symbol	Represents
$L$	locator that gives the rank of a value
$P_k$	$k^{\text{th}}$ percentile

•  $k$  percentage (blw 0% & 100%)  
•  $P_k = \frac{k}{100}$  (decimal)

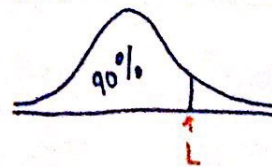
#### Finding a Percentile Associated to a Given Score

Formula:

$$P_k \text{ percentile of score} = \frac{\# \text{ of scores less than given score}}{\text{total number of scores}}$$



Ex: the 90%





Also:  $k = \frac{L}{n} \times 100 \leftrightarrow P_k = \frac{L}{n}$

**Finding the Score Associated to a Given Percentile**

Formula:

$$L = \frac{k}{100} \cdot n$$

Note: If  $L$  is a decimal, then always round up to find the score with that specific rank.

If  $L$  is a whole number, then average the  $k^{\text{th}}$  score and the next higher score.

*ie the one in that ordered spot*

*so 128 is really 128,000*

Ex: The following data set represents the selling price (in thousands) of 38 randomly selected homes.

128	135	138	145	149	152	155	158	159	163	163	165	167
168	170	170	172	173	176	177	180	180	185	188	191	193
199	205	210	212	215	229	233	250	325	450	500	525	

*Annotations: min at 128, max at 525, 10th spot at 163, 29th spot at 210, 33rd spot at 233.*

(a) Find the percentile corresponding to a selling price of \$188,000.

$$\% \text{ score} = \frac{\# \text{ less than } 188}{\text{total}} = \frac{23}{38} \approx 60.5\% = 61\% \text{ or } 61^{\text{st}} \text{ percentile}$$

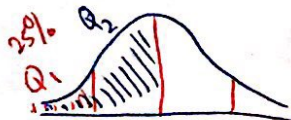
(b) Find the home price corresponding to the 85th percentile.

$$L = \frac{k}{100} \cdot n = \frac{85}{100} \cdot (38) = 32.3 = 33^{\text{rd}} \text{ spot} \rightarrow \$233,000 \text{ is the price of home at the } 85^{\text{th}} \text{ percentile.}$$

**Quartiles**

Def **Quartiles** (denoted  $Q_1$ ,  $Q_2$ , and  $Q_3$ ) are measures of location which divide a data set into four groups with about 25% of the values in each group.

Note:  $Q_1 = P_{25}$ ,  $Q_2 = \text{the median}$ , and  $Q_3 = P_{75}$



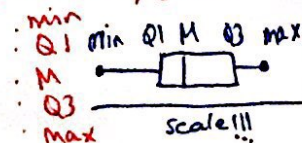
**Boxplot**

$Q_4 = P_{100} = \text{everything!}$

Def A **boxplot** is a graph of a data set that consists of a line extending from the minimum value to the maximum value, and a box with lines drawn at the first quartile, the median, and the third quartile.

Note: A **five-number summary** refers to the five values used to draw the boxplot.

**BOX PLOT / 5 # Summary**



Ex: Find the five-number summary for the previous data regarding home prices.

(a) Find the minimum and maximum values in the data set.

$\text{min} = \$128,000$ ,  $\text{max} = \$525,000$

(b) Find  $Q_2$  (the median).

$$Q_2 = M = \frac{176 + 177}{2} = 176.5 \rightarrow Q_2 = \$176,500$$

(c) Find  $Q_1 = P_{25}$ .

(d) Find  $Q_3 = P_{75}$ .

$$Q_1 = P_{25} = \frac{25}{100} \cdot (38) = 9.5 \rightarrow 10^{\text{th}} \text{ spot}$$

$$Q_3 = P_{75} = \frac{75}{100} \cdot (38) = 28.5 \rightarrow 29^{\text{th}} \text{ spot}$$

$Q_1 = \$163,000$

$Q_3 = \$210,000$

(e) List the **five-number summary**

$\$128,000, \$163,000, \$176,500, \$210,000, \$525,000$

Ex: Construct a boxplot for the previous data set regarding home prices.

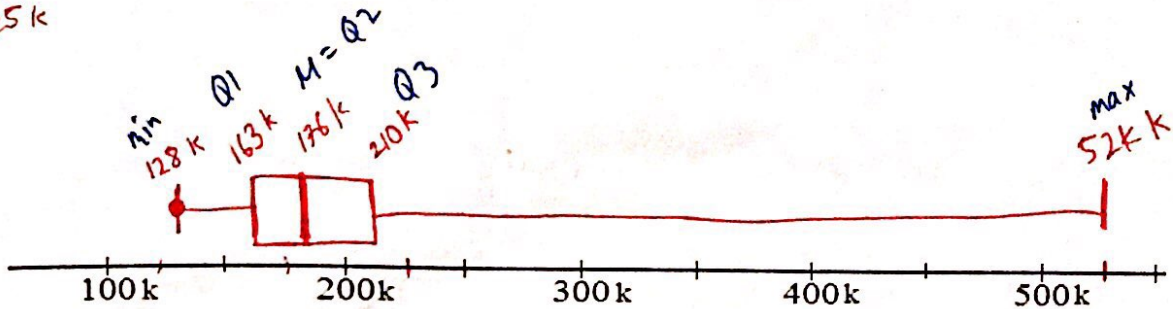
**Graphing a boxplot:**

- (a) Place the five number summary for a data set on a number line, and draw a straight line connecting them.
- (b) Draw vertical lines at each of the five number summary values.
- (c) Draw a rectangle connecting  $Q_1$  to  $Q_3$ .

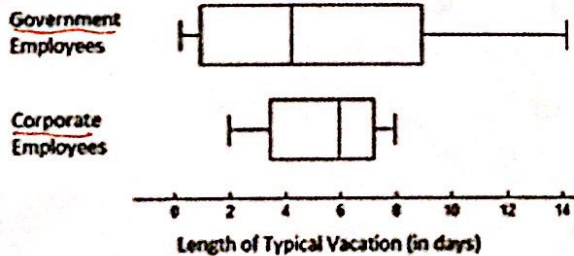
5# Summary

- 128 k
- 163 k
- 176 k
- 210 k
- 525 k

important: draw an accurate scale (horizontal axis) first, then draw box-plot.



Ex: Below, boxplots are shown for the length of a typical vacation for California residents who work for the government in some capacity and for those who work for a private company.



Based on these graphs, would you prefer a government job or corporate job based solely on the length of their vacations? There is no one right answer, but please consider center, spread, and any other relevant statistics or values from the boxplots shown to support your data.

Pro Government

- \* larger spread so can sometimes have short vacay (0 days) but sometimes have long vacay (14 days)!
- \* median = 4 days, so 50% of employees have ~~any~~ vacay length 4 to 14 days

Pro-Corporate

- \* less spread so pretty consistent vacay length of about 6 days (better to plan!)
- \* at least 2 days off guaranteed where govt can have 0 days off
- \* median = 6 days so median higher 50% have vacay 2 to 6 days 50% have vacay 6 to 8 days

10