

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. → many reasonable ways to define this
 The three most widely-used measures of center are the mean, median, and mode. Mean = "average"

The (arithmetic) **mean** of a data set is computed by adding all of the values of the variable in the data set and dividing by the number of observations.

The **population** arithmetic mean, μ , is computed using ALL of the individuals in a population. The population mean is a parameter.

$$\mu = \frac{x_1 + x_2 + \dots + x_N}{N} = \frac{\sum x}{N}$$

\sum = sum

The **sample** arithmetic mean, \bar{x} , is computed by using some of the individuals in a population. The sample mean is a statistic.

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{\sum x}{n}$$

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.

$$\bar{x} = \frac{\sum x}{n} = \frac{82 + 77 + 90 + 71 + 62 + 68 + 74 + 84 + 94 + 88}{10} = 79$$

Calc Tip: • Enter Data into List (L1)
 • STAT > Calc > 1-VAR STATS

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

Med The **median** of a data set is the value that lies in the "MIDDLE" of the data when arranged in ascending order. We use M to represent the median.

ODD number of data

1st: Arrange the data in ascending (increasing) order

2nd: The median will be the middle number

Ex: 11, 14, 16, 19, 28 (inc) (#5 odd)
 Median

EVEN number of data

1st: Arrange the data in ascending order

2nd: The median will be the

MEAN of the middle numbers

Ex: 14, 18, 20, 26, 31, 39

mean of 20 & 26 = $\frac{20+26}{2} = 23$

The **midrange** of a data set is the value midway between the minimum and maximum values.

MR

$$\text{Midrange} = \frac{\text{min value} + \text{max value}}{2}$$



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

Using 1-VAR STATS

• Med = 79.5

• MR = $\frac{62+94}{2} = 78.0$

"By Hand"

Sort: 62 68 71 74 77 82 84 88 90 94

Med = $\frac{77+82}{2} = 79.5$

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

STATS RULE OF ROUNDING

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 N/A.
 *If the data set has more than one observation that repeat the same number of time, then it is considered multimodal. (bimodal)

Ex: Find the **mode** for each example below.

a) The following data represent the number of O-ring failures on the shuttle *Columbia* for its 17 flights prior to its fatal flight:

0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 2, 3
mode is 0

b) The data of the test scores from above:

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

mode is N/A

c) Hair color of ten people in line:
Brown, Blonde, Red, Brown,
Brown, Blonde, Brown, Blonde,
Blonde, Red

(Bimodal)

mode is Brown & Blonde

4 Br
4 Bl
2 R

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)	$f \cdot x$
Class 1	3-4.9	5	$\frac{3+4.9}{2} = 3.95$	$5 \times 3.95 = 19.75$
Class 2	5-6.9	32	$\frac{5+6.9}{2} = 5.95$	$32 \times 5.95 = 190.4$
Class 3	7-8.9	40	$\frac{7+8.9}{2} = 7.95$	$40 \times 7.95 = 318$
Class 4	9-10.9	18	$\frac{9+10.9}{2} = 9.95$	$18 \times 9.95 = 179.1$
Class 5	11-12.9	5	$\frac{11+12.9}{2} = 11.95$	$5 \times 11.95 = 59.75$
	$n = \sum f = 100$		<u>data</u>	$\sum (f \cdot x) = 767$

Find the sample mean.

$$\bar{x} = \frac{\sum (f \cdot x)}{n} = \frac{767}{100} = 7.67 \text{ lbs.}$$

Sum up

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.

$$\text{Med} = \frac{2+2}{2} = \frac{4}{2} = 2 \quad \text{Mean } \bar{x} = \frac{\sum x}{n} = \frac{14}{6} = 2.3 \quad \boxed{\text{Med} = 2} \quad \boxed{\text{Mean} = 2.3}$$

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

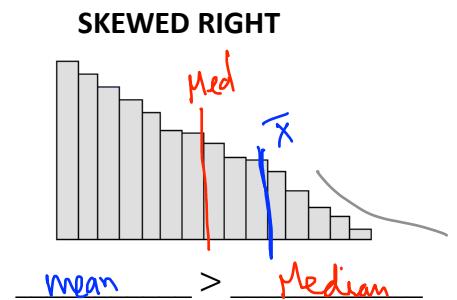
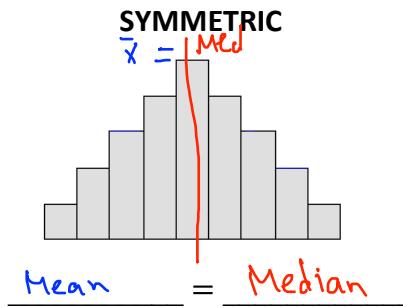
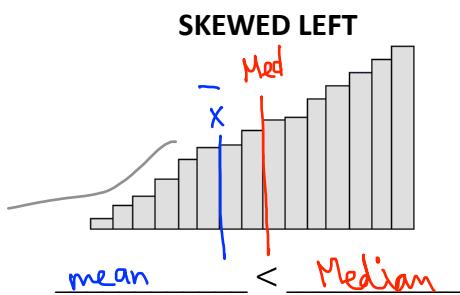
$$\bar{x} = \frac{14+102}{7} = \frac{116}{7} = 16.6$$

$$\boxed{\bar{x} = 16.6}$$

1, 1, 2, 2, 3, 5, 102
 ↑ Med ↓ Outlier
 $\boxed{\text{Med} = 2}$

- key points
- ① Med is resistant
 $\text{Med } 2 \rightarrow \text{Med } 2$
 - ② Mean is not resistant!
 $\bar{x} = 2.3 \rightarrow \bar{x} = 16.6$

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.

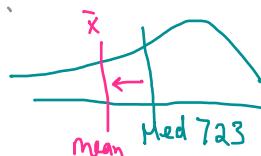
Ex: FICO scores range in value from 300 to 850, with a higher score indicating a more creditworthy individual. The distribution of FICO scores is skewed left with a median score of 723.

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

Less than

mean is less than 723.

↳ preferred ans.



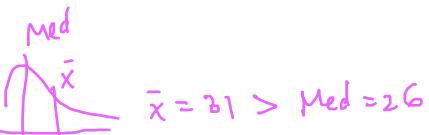
(b) What proportion of individuals have a FICO score above 723?

↳ decimal

0.50

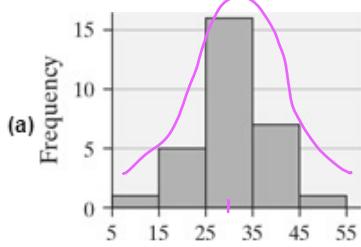
about $\frac{1}{2}$ of scores above median!

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



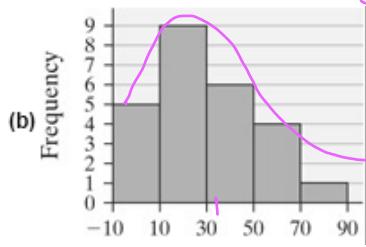
	Mean	Median
1	42	42
2	31	36
3	31	26
4	31	32

$\rightarrow Med = \bar{x}$ Symmetric
 $\rightarrow \bar{x} = 31 < Med = 36$

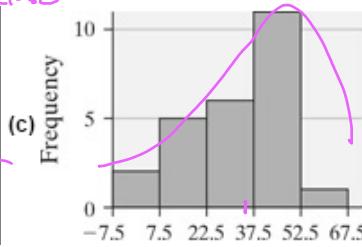


4

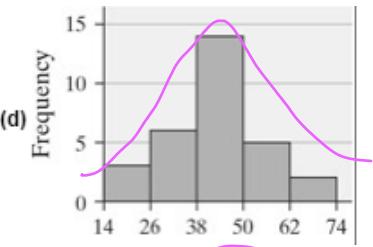
process
of elim.



3



2



1

Jan 10

3.2 Measures of Variation

Measures of Spread / change

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

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 →

B/C all values are much closer together
 • Advil ± 2 mg from mean $\bar{x} = 200$ mg
 • Motrin ± 7 mg from mean $\bar{x} = 200$ mg

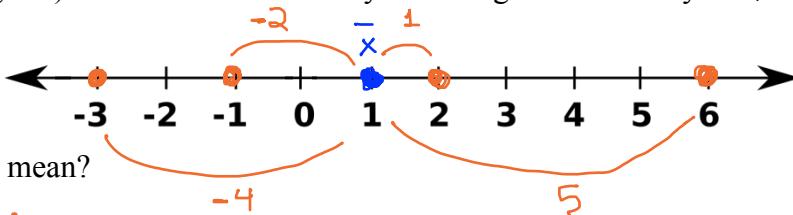
Ex: The following are temperatures (in degrees) on four consecutive days in Mongolia in January: $-3, -1, 2, 6$

$$\bar{x} = \frac{\sum x}{n}$$

(a) Find the mean.

$$\bar{x} = \frac{-3 + 1 + 2 + 6}{4} = 1$$

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



PoTIP Use $+$ or $-$ to indicate when data is above (+) or below (-) the mean.

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

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:

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

conceptual

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 σ) 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:

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

units are same as data!

Good Tip.

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

→ units are units squared of data.

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

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

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



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.

25 31 28 19 24 26 29 32 20

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

Use the alternative form to find the standard deviation.

skip ☺

x	x^2
19	
20	
24	
25	
26	
28	
29	
31	
32	
$\Sigma x =$	$\Sigma x^2 =$

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

Small! I want to know what I'm paying for & get what I'm paying for!

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)	$f \cdot x$	$f \cdot x^2$
3-4.9	5	$\frac{3+4.9}{2} = 3.95$	$5 \cdot 3.95 = 19.75$	$5 \cdot (3.95)^2 = 78.0$
5-6.9	32	5.95	190.4	$32 \cdot (5.95)^2 = 1132.88$
7-8.9	40	7.95	318	$40 \cdot (7.95)^2 = 2528.1$
9-10.9	18	9.95	179.1	$18 \cdot (9.95)^2 = 1782.045$
11-12.9	5	11.95	59.75	$5 \cdot (11.95)^2 = 714.0125$
$n = \sum f = 100$			$\sum(f \cdot x) = 767$	$\sum(f \cdot x^2) = 6235.05$

Find the sample standard deviation and variance.

$$s = \sqrt{\frac{100 \cdot [6235.05] - [767]^2}{100 \cdot 99}} = \sqrt{\frac{35216}{9900}} = 1.886046 \approx 1.886 \text{ lbs}$$

Range Rule of Thumb

$$\text{Var } s^2 = \frac{35216}{9900}$$

$$3.5571717 \dots \approx 3.557 \text{ lbs}^2$$

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	AIRPORT CODE	WIFI SPEED
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

<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?

Calc gives: $\bar{x} = 42.73 \text{ Mbps}$, $\sigma = 15.98 \text{ Mbps}$ (WARNING: actually use Sample St.dev Sx)

" $x \leq \mu - 2\sigma$ means" the data value x is 2 standard deviations less than the mean at least

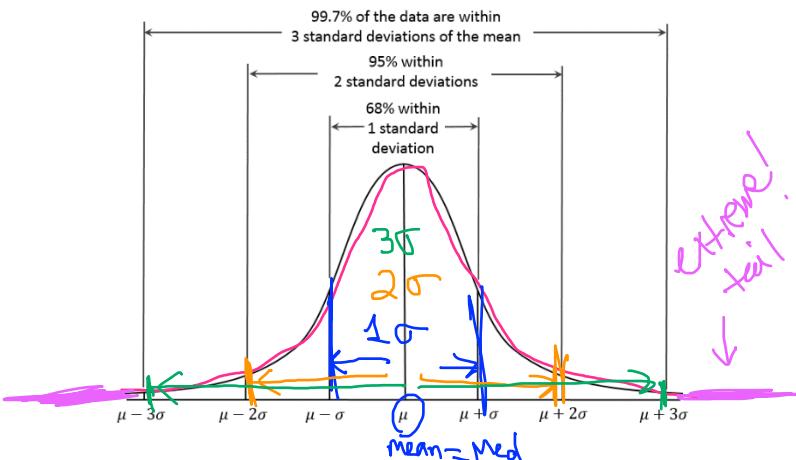
→ Sig High: $\mu + 2\sigma = 42.73 + 2(15.98) = 74.69 \text{ Mbps}$. Conclusion "60.8 Mbps is NOT significantly fast!"

Empirical Rule

68-95-99.7 Rule

Empirical Rule says that a normal or Bell-Shaped or Symmetric distribution has...

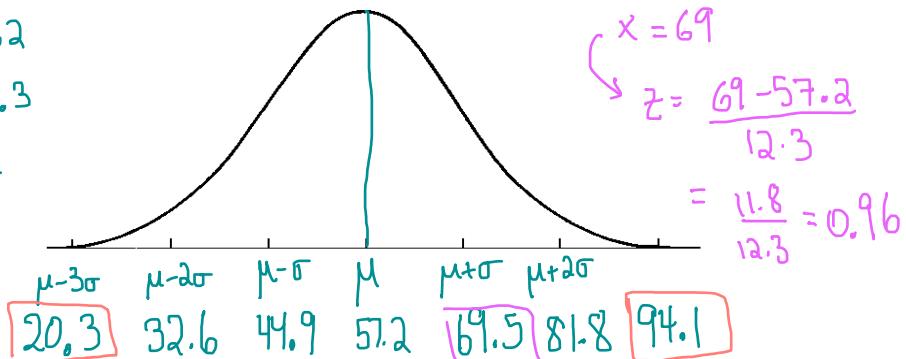
- Approximately 68 % of the data will lie within 1 standard deviation of the mean.
- Approximately 95 % of the data will lie within 2 standard deviation of the mean.
- Approximately 99.7 % of the data will lie within 3 standard deviation of the mean.



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, μ , is 57.2 years and the population standard deviation, σ , 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

$$\begin{aligned} \mu &= 57.2 \\ \sigma &= 12.3 \\ N &= 42 \end{aligned}$$



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

By Empirical Rule, approximately 99.7% of patients.

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

Between 20.3 and 94.1 years.

$\frac{69.5 - 57.2}{12.3} = \frac{12.3}{12.3} = 1$

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

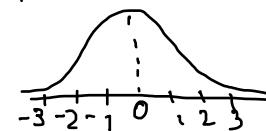
(b) STAT \Rightarrow CALC \Rightarrow 1: 1-Var Stats \Rightarrow List: (Choose list) \Rightarrow Calculate



Jan 13

3.3 Measures of Relative Standing and Boxplots

Standard Normal Dist



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.

own rule!

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

so excellent to help compare different data sets.

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:

$$z < -2 \text{ or } z > 2$$

$$z = -3.5 \quad z = 2.67$$

Key



3. The z score allows us to compare data values drawn from different samples or populations.

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?

Jacob's class:

82

$$\bar{x} = 78, s = 6$$

\bar{x} mean
 s sample st. dev.

$$z = \frac{82 - 78}{6} = 0.67 \text{ Jacob}$$

Michael's class:

$$\bar{x} = 55, s = 12$$

Compute z-scores

70

$$z = \frac{70 - 55}{12} = 1.25 \text{ Michael}$$

(M → E)

Michael wins the bet since he did better in his class relative to his classmates

Percentiles

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^{th} percentile

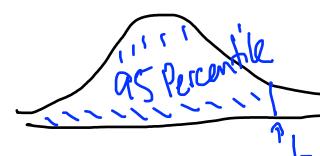
Finding a Percentile Associated to a Given Score

Formula:

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

L or P_k

Careful! not including!



Finding the Score Associated to a Given Percentile

Formula:

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

Score Value

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^{th} score and the next higher score.

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	

(a)

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

$$\frac{\text{# homes less than } \$188,000}{\text{total}} = \frac{23}{38} \approx 0.605 = 60.5\% \text{ or } 61 \text{ percentile}$$

(b)

Find the home price corresponding to the 85th percentile.

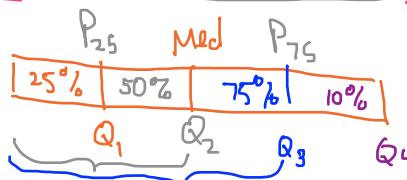
$$L = \frac{k}{100} * n = \left(\frac{85}{100} \right) * 38 = 32.3 \rightarrow \text{so round up to } 33^{\text{rd}} \text{ home}$$

home price of \$233,000 corresponds to 85th 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 $Q_3 = P_{75}$



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.

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

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

$$\min = \$128,000 \quad \max = \$525,000$$

(b) Find Q_2 (the median).

$$\text{Med} = \$176,500$$

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

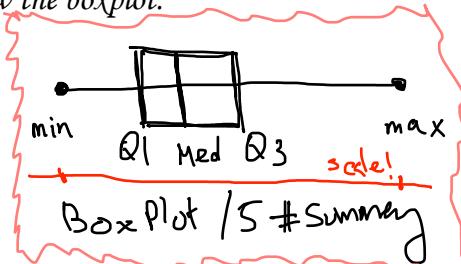
$$Q_1 = \$163,000$$

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

$$Q_3 = \$210,000$$

(e) List the five-number summary.

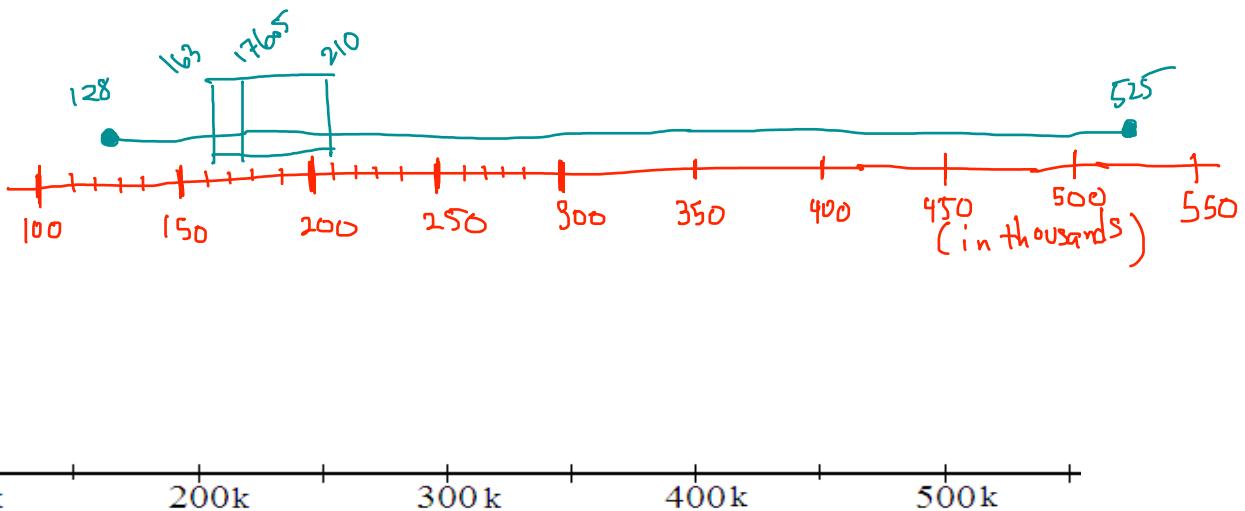
$$\$128,000, \$163,000, \$176,500, \$210,000, \$525,000$$



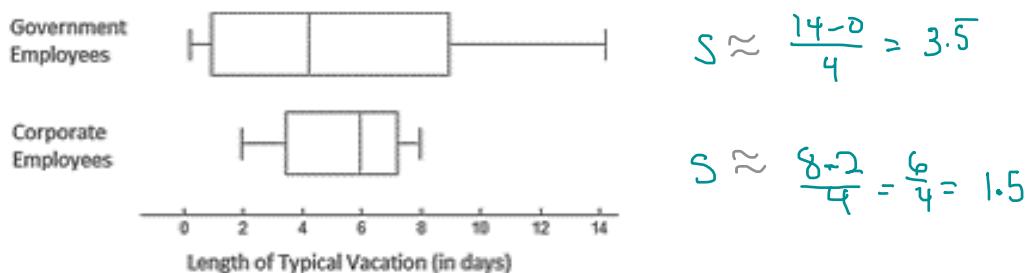
Ex: Construct a boxplot for the previous data set regarding home prices.

Graphing a boxplot:

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



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.



→ Answer Vary.

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.

Bonus Range Rule of Thumb Approximation for Standard Deviation (Sample)

$$S \approx \frac{\text{Range}}{4} = \frac{\max - \min}{4}$$