# Collecting & Organizing Univariate Data

<u>Univariate</u> – having only 1 variable

• Data can either be Qualitative or Quantitative
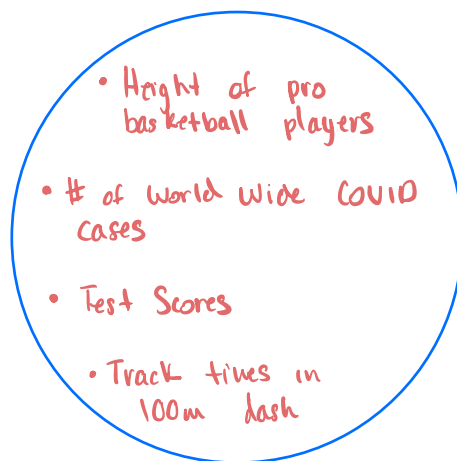
<u>Quantitative</u> – numerical data focused on <u>Quantity</u>
• Can be Discrete or Continuous

<u>Qualitative</u> – Non-numerical data focused <u>Quality</u> (i.e. Categorical)

lets think of some examples of each type of data

### Quantitative

- Height of pro basketball players

- # of world wide COVID cases

- Test Scores

- Track times in 100m dash

### Qualitative

- Ice creame flavors

- Happiness Rating

- Pass/Fail

- Eye Color

- Interview Transcript

# Central Tendency

aka "averages" (mean, mode, median)

"X bar"

$\Sigma$ - "the sum of"

**Mean ($\bar{X}$)** : average value in data set  $\bar{X} = \frac{1}{n} \sum_{i=1}^{n} X_i$  , where  n - # of values

Set: 6   8   9   10   12

$\overset{"}{X_1}$  $\overset{"}{X_2}$  $\overset{"}{X_3}$  $\overset{"}{X_4}$  $\overset{"}{X_5}$

n = 5 terms

$$\bar{X} = \frac{1}{5} \sum_{i=1}^{5} X_i \quad = \quad \frac{1}{5}\left( X_1 + X_2 + X_3 + X_4 + X_5 \right)$$

$$= \frac{1}{5}\left( 6 + 8 + 9 + 10 + 12 \right) \quad = \frac{1}{5}\left( 45 \right) = 9$$

**Mode** : The most frequently occurring value in the data set

✗ there can be 0, 1, or 2 modes (namely, bimodal)

Set₁:   1   2   3   3   4   4   5   6          mode: 3, 4

Set₂:   1   2   3   4   5   6                    mode: none

Set₃:   1   2   3   4   5   6   6              mode: 6

**Median** : The middle value in a <u>sorted</u> data set

n terms

Set₁ :   1   2   3   4   5   6

✗ If n is even (there's no exact middle term)
then median is the avg of the two middle terms  = 3.5

Set₂ :   8   10   12   13   14

median: 12

**ex)** Grades for a History test for 14 students are shown below

69   58   67   66   58   79   83   76   44   35   58   88   91   47

When a 15ᵗʰ student took the test the mean became 66.2
Calculate the grade for the 15ᵗʰ student

$$\text{Mean}_i = 66 \quad , \quad \text{Mean}_f = 66.2 \quad \Rightarrow \quad \text{Student}_G \equiv \frac{924 + x}{15} = 66.2$$

$$x = 15 \cdot 66.2 - 924 = 69$$

# Frequency Tables

The lengths, in minutes, of 20 telephone calls are shown below

| 4.2 | 6.8 | 10.4 | 8.2 | 11.5 | 1.6 | 5.8 | 7.6 | 3.1 | 21.5 |
| 13.5 | 5.8 | 4.1 | 22.8 | 13.6 | 11.2 | 9.5 | 1.8 | 12.4 | 4.9 |

Organize the data into a Frequency Table

| Length (mins) | Frequency |
|---|---|
| $0 \leq t < 5$ | 6 |
| $5 \leq t < 10$ | 6 |
| $10 \leq t < 15$ | 6 |
| $15 \leq t < 20$ | 0 |
| $20 \leq t < 25$ | 2 |

*"Grouped data" because it has continuous ranges for each class
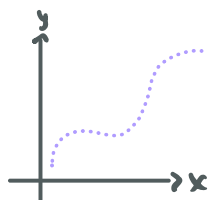
IS this data Continuous or Discrete?

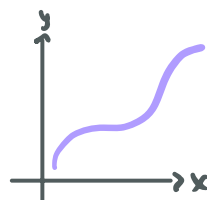Discrete data - can be counted (ex: shoe size, # of cars in a parking lot)

Continuous data - can be measured (height, weight, time)

* must contain full range of values

How can you tell if a graph is continuous or discrete?



Discrete



Continuous

ex] Men's Jeans are sized by waist measurements. Here are the
the jean sizes of 10 men:

28   30   28   34   32   30   36   28   30   30

(a) Find the mean, median, mode, Range

(b) Create a frequency table
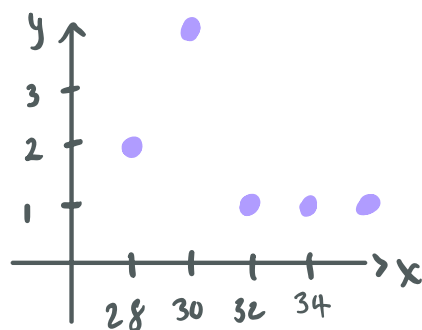
(c) Is this data Continuous or discrete? why?

(a)  Sorted: 28   28   28   30   30   30   30   32   34   36
     Mean = 30.6
     Median = 30
     Mode = 30
     Range = 8

(b)

| Jean Size | Frequency |
|-----------|-----------|
| 28 | 3 |
| 30 | 4 |
| 32 | 1 |
| 34 | 1 |
| 36 | 1 |

✗ Modal Group/Class - class or Category
   that contains the highest
   frequency
       => Jean Size 30 is modal class

(c) Discrete because jean sizes are only whole even values

# Outliers

**ex)** The ages of 15 cats are below:

10  10  11  11  11  12  12  12  12  13  13  14  14  24  25

Find the mean, mode, and median for this data

ⓐ Mean = 13.6

ⓑ Median = 12

ⓒ Mode = 12

Are there any singular data points that affect the calculation of the mean more so than the others? Remove these values and compare recalculate the mean.

✗ 24, 25 are examples of <u>Outliers</u> : extreme data values that can distort the results of statistical processes

  -They don't fit with the rest of the data

✗ Outliers are often the result of errors in collecting/reading data



Get in Loser, we're going shopping.

Outliers ex] Find the mean, median, and mode for each data set
& comment on any data values you think are
outliers:

(a) The heights of 15 sunflowers:

1.1   2.2   2.5   2.5   2.5   3.1   3.5   3.6
3.9   4.0   4.1   4.4   4.6   4.9   6.1

Mean = 3.53
Median = 3.6          Outliers: 1.1, 6.1
Mode = 2.5

(b) 20 Students' Geography test scores:

22   39   45   46   46   52   54   58   62   62
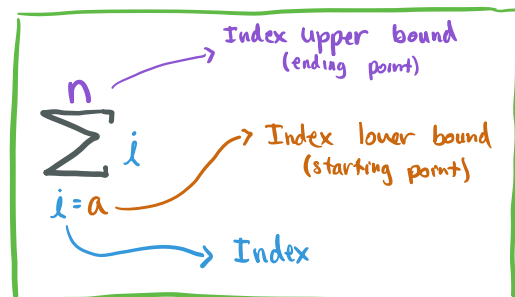62   67   70   75   78   82   89   91   95   98

Mean = 64.65
Median = 62          Outliers: 22
Mode = 62

---

Sigma Notation Warm Up

Recall: $\Sigma$ - upper case "sigma"
$\sigma$ - lower case "sigma"

$$\sum_{i=a}^{n} i$$

Index upper bound (ending point)

Index lower bound (starting point)

Index

① Find $\sum_{i=1}^{4}(i+i^2) = (1+1^2)+(2+4)+(3+9)+(4+16)$
  $2 + 6 + 12 + 20$ => 40

② Find $\sum_{i=1}^{4}i + \sum_{i=1}^{4}i^2 = (1+2+3+4)+(1^2+2^2+3^2+4^2)$ => 40

③ Find $\sum_{s=1}^{3}y = y+y+y = 3y$

# Measures of Dispersion

- Dispersion- the action/process of distributing things over a wide area
- Measures how spread out the data is
- We already learned a measure of Dispersion: **Range**

ex] Consider the two data sets:

$S_1 = $ -10, 0, 10, 20, 30          $S_2 = $ 8, 9, 10, 11, 12

mean: $\frac{50}{5} = 10$          mean: $\frac{50}{5} = 10$

Range: $30 + 10 = 40$          Range: $12 - 8 = 4$

Thus $S_1$ is more dispersed (more spread out)

## Standard Deviation ($\sigma$) : measure of how data values are spread out
in relation to the mean

aka "root-mean-squared deviation"

$\cancel{*} \sigma$ - "sigma"

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

$\cancel{*}$ mean of sample

$$\sigma_1^2 = \frac{(-10-10)^2 + (0-10)^2 + (10-10)^2 + (20-10)^2 + (30-10)^2}{5} = 200$$

$$\sigma_2^2 = \frac{(8-10)^2 + (9-10)^2 + (10-10)^2 + (11-10)^2 + (12-10)^2}{5} = 2$$

$\Rightarrow$   $\sigma_1 = \sqrt{200} = 10\sqrt{2} \longrightarrow$ more dispersed than $\sigma_2$

$\sigma_2 = \sqrt{2}$

**ex)** The # of Ice creams sold over a period of 13 weeks:

146  151  158  158  161  149  160  147  158  160  216  225  238

Find the Standard Deviation

$$\bar{X} = \frac{2227}{13} = 171.3$$

$$\sigma^2 = \frac{(146-171.3)^2 + (149-171.3)^2 + (151-171.3)^2 + 3(158-171.3)^2 + 2(160-171.3)^2 + (161-171.3)^2 + (216-171.3)^2 + (225-171.3)^2 + (238-171.3)^2 + (147-171.3)^2}{13} = 950.98$$

$$\Rightarrow \boxed{\sigma = \sqrt{950.98} = 30.8}$$

This is EXHAUSTING

We can do some algebraic manipulation to Rewrite this formula and make it simpler

$$\sigma^2 = \frac{1}{n}\sum_{i=1}^{n}(x_i - \bar{x})^2 = \frac{1}{n}\sum_{i=1}^{n}\left(x_i^2 - 2\bar{x}\cdot x_i + \bar{x}^2\right)$$

$$= \frac{1}{n}\sum_{i=1}^{n}x_i^2 - 2\bar{x}\underbrace{\left[\frac{1}{n}\sum_{i=1}^{n}x_i\right]}_{=\bar{x}} + \underbrace{\left[\frac{1}{n}\sum_{i=1}^{n}\bar{x}^2\right]}_{=\frac{1}{n}(n\bar{x}^2)}$$

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

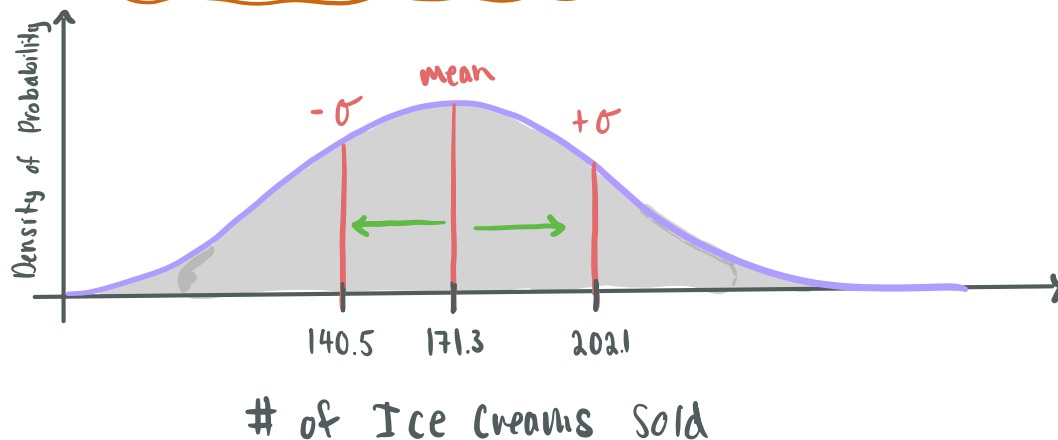$$\sigma = \boxed{= \sqrt{\frac{1}{n}\sum_{i=1}^{n}(x_i^2) - \bar{x}^2}}$$   ✗ Easier and should be used

Recalculate the Standard Deviation using our newly formulated definition

$$\sigma^2 = \frac{146^2 + 151^2 + 158^2 + 158^2 + 161^2 + 149^2 + 160^2 + 147^2 + 158^2 + 160^2 + 216^2 + 225^2 + 238^2}{13} - 171.3$$
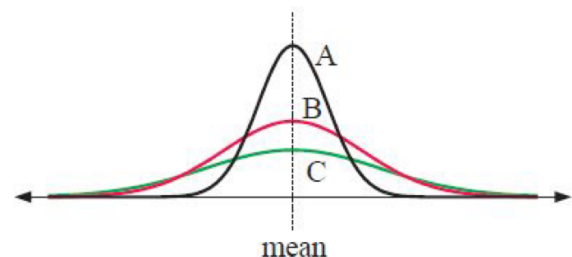
$$\sigma = \sqrt{953.6} = 30.8$$

## Normal Distribution Curve



Density of Probability (y-axis)

mean

$-\sigma$     $+\sigma$

140.5     171.3     202.1

# of Ice Creams Sold

★ Density of Probability - likelihood of obtaining a value corresponding to the x-axis

For Example

There's a high probability that the ice cream shop will sell 171.3 ice cream cones in 13 weeks because it has the highest y-value
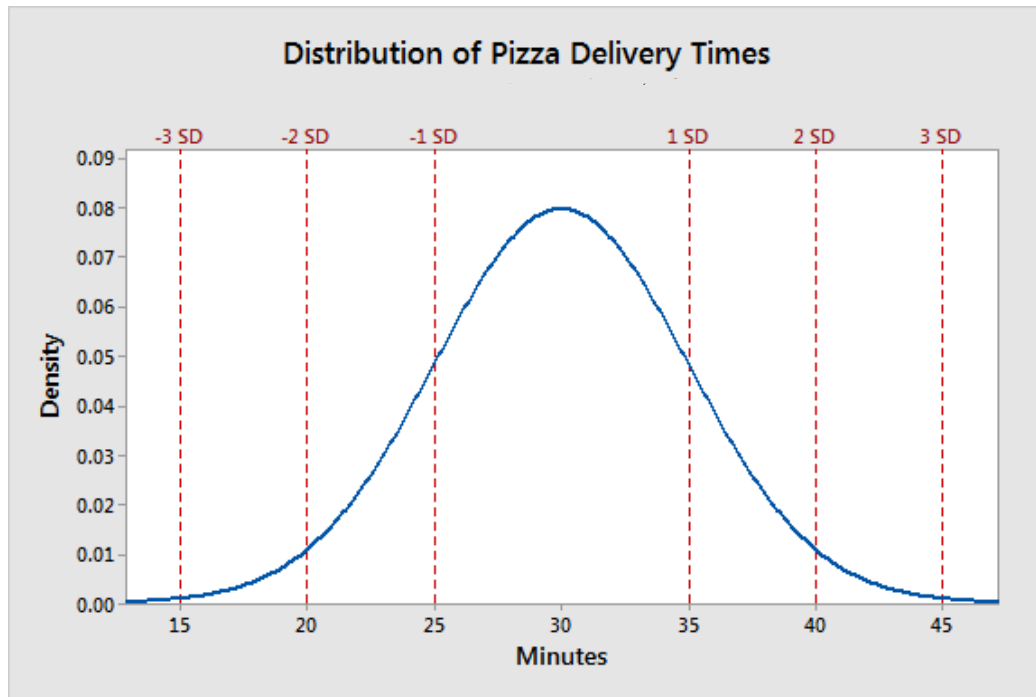
The given distributions have the same mean, but clearly they have different spreads. The A distribution has most scores close to the mean whereas the C distribution has the greatest spread.



A
B
C
mean

ex] Use the graph below to find:

    ⓐ Mean ($\bar{x}$): 30 mins

    ⓑ Standard Deviation ($\sigma$): 5 mins

**Distribution of Pizza Delivery Times**



# Interquartile Range (IQR)

better represents the dispersion of data compared to Range

ex] The following data set shows the # of animal crackers in each kid's lunch box

4   4   10   11   15   7   14   12   6

Step 1) Find the Median: 10

Sort list: 4  4  6  7  10  11  12  14  15

    lower Quartile $Q_1$        Upper Quartile $Q_3$

$\boxed{\text{Step 2}}$ Find median of $Q_1$ and $Q_3$

$$Q_1 = \frac{4+6}{2} = 5 \qquad\qquad Q_3 = \frac{12+14}{2} = 13$$

$\boxed{\text{Step 3}}$ Find the difference btw the two Ranges

$$Q_3 - Q_1 = 13 - 5 = 8$$

$$\Rightarrow \boxed{IQR = 8}$$

IQR = Upper Quartile ($Q_3$) — Lower Quartile ($Q_1$)

data point at the          data point at the
75% percentile            25% percentile

[EX] Find the IQR of the data shown in the dot plot below

Songs on Each Album in Shane's Collection



Number of Songs

| 7 | 9 | 9 | 10 | 10 | 10 | 11 | 12 | 12 | 14 |

Median: $\frac{10+10}{2} = 10$

$Q_1$ median: 9

$Q_2$ median: 12

$IQR = Q_3 - Q_1 = 12 - 9 = \boxed{3}$

A **percentile** is the score below which a certain percentage of the data lies.

For example:
- the 85th percentile is the score below which 85% of the data lies.
- If your score in a test is the 95th percentile, then 95% of the class have scored less than you.

Notice that:
- the **lower quartile** ($Q_1$) is the 25th percentile
- the **median** ($Q_2$) is the 50th percentile
- the **upper quartile** ($Q_3$) is the 75th percentile.

A cumulative frequency graph provides a convenient way to find percentiles.

# Revisiting Outliers
— an extremely high or extremely low value in our data

A data value ($z$) is an outlier if:

$$z > Q_3 + 1.5(IQR)$$

or

$$z < Q_1 - 1.5(IQR)$$

ex) Phone Calls Recieved in a day

Find any outliers

| 10 | 12 | 11 | 15 |
| 11 | 14 | 13 | 17 |
| 12 | 22 | 14 | 11 |

Sort: | 10 | 11 | 11 | 11 | 12 | 12 | 13 | 14 | 14 | 15 | 17 | 22 |

Median: 12.5
$Q_1$: 11
$Q_3$: 14.5
IQR: 3.5

$\Rightarrow$

Outliers $<$ $11 - 1.5(3.5)$

$14.5 + 1.5(3.5) >$ Outliers

$\Downarrow$

Outliers $< 5.75$ $\longrightarrow$ None

$19.75 >$ Outliers $\longrightarrow$ 22

# Central Tendencies of Grouped Data

☆ Note if only given a grouped frequency table then you cannot find the exact values for mean, median, and mode but we can approximate them

ex) Consider the grouped data of student test scores:

| Score | Frequency | Midpoint | (Midpoint)(Frequency) |
|---|---|---|---|
| $0 \le x < 60$ | 5 | 30 | 150 |
| $60 \le x < 70$ | 4 | 65 | 260 |
| $70 \le x < 80$ | 10 | 75 | 750 |
| $80 \le x < 90$ | 15 | 85 | 1275 |
| $90 \le x \le 100$ | 4 | 95 | 380 |

Find

ⓐ Mean

Step 1: Find Midpoint of each class

Step 2: Find product of Midpoints and Frequencies

Step 3: Find n - total # data points $= \sum \text{Frequencies} = 38$

Step 4: Find sum of all (Midpoint)(Frequency) products $= 2815$

Step 5: $\overline{X} = \dfrac{\sum (\text{Midpoint})(\text{Frequency})}{n} = \dfrac{2815}{38} = \boxed{76.1}$

ⓑ Median Class

$\dfrac{n}{2} = \dfrac{38}{2} = 19$

$\Rightarrow$ 19th and 20th values

Either $\begin{array}{l} 70 \le x < 80 \\ 80 \le x < 90 \end{array}$

ⓒ Modal Class

$\boxed{80 \le x < 90}$ — class with highest frequency

# Summary

<u>Quantitative</u> - numerical data focused on <u>Quantity</u>

<u>Qualitative</u> - non-numerical data focused <u>Quality</u> (i.e. Categorical)
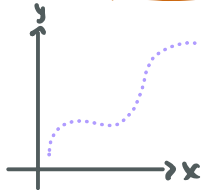
<u>Univariate</u> - having only 1 variable

<u>Central Tendency</u> - values that describe the middle of a data set

$$\left( \begin{array}{c} \text{mean, median, and mode are all measures of} \\ \text{Central tendency} \end{array} \right)$$

<u>Mean</u> $(\overline{x})$ : average value in data set. $\overline{x} = \frac{1}{n} \sum\limits_{i=1}^{n} (x_i)$ , where n- # of values
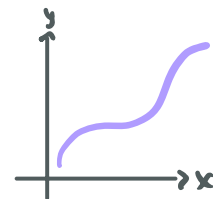
<u>Mode</u> : The most frequently occurring value in the data set

<u>Median</u> : The middle value in a <u>sorted</u> data set

<u>Discrete data</u> - can be <u>counted</u> (ex: shoe size, # of cars in a parking lot)



<u>Continuous data</u> - can be <u>measured</u> (height, weight, time)



<u>Outliers</u> : extreme data values that can distort the results of statistical processes

A data value (z) is an outlier if:

$$z > Q_3 + 1.5(IQR)$$
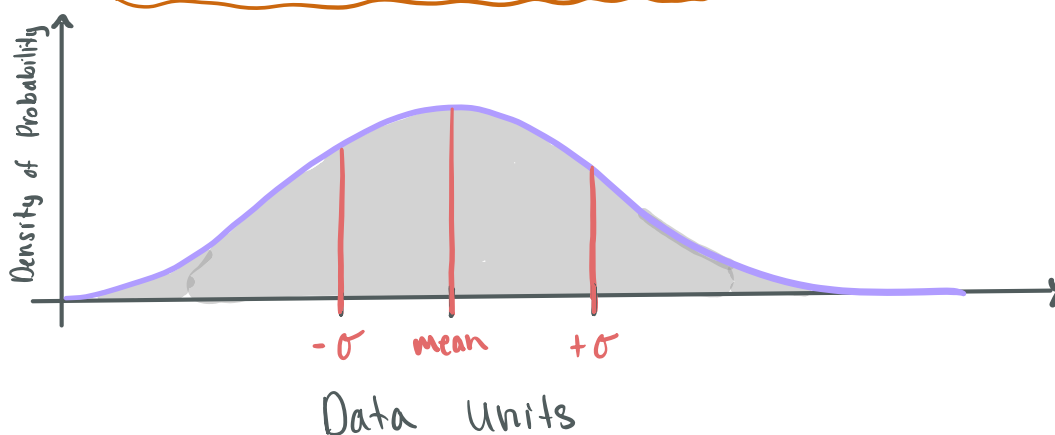$$\text{or}$$
$$z < Q_1 - 1.5(IQR)$$

<u>Standard Deviation $(\sigma)$</u> : measure of how data values are spread out in relation to the mean

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

## Normal Distribution Curve



<u>Measures of Dispersion</u> : values that describe how spread out the data is ( Range, and IQR are measures of Dispersion)

<u>Interquartile Range (IQR)</u> $= Q_3 - Q_1$
(better represents dispersion)

<u>Range</u> : (highest value) $-$ (Lowest value)

### Practice Problems

Pg 99 , 3A, Q1

Pg 100, Ex2, Q.2

Pg 107, Ex5, Q.b,c

Pg 109, Ex6, Q. a,b,c

Pg 110, Investigation 6