October - 28:

Measures of Dispersion:

1. Variance
2. Standard deviation
3. Range

1. Variance:

Always variance should be less if we are getting higher variance it will reduce/affect our prediction accuracy.

population variance

Sample variance

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

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

Degree of freedom.

Ex: $x = \{1, 2, 2, 3, 4, 5\}$

$\mu = \frac{17}{6} = 2.8$

$$\sigma^2 = \frac{(1-2.8)^2 + (2-2.8)^2 + (2-2.8)^2 + (3-2.8)^2 + (4-2.8)^2 + (5-2.8)^2}{6}$$

$\delta^2 = \frac{10.84}{5}$

$= 2.16$.

$= \frac{3.24 + 0.64 + 0.64 + 0.04 + 1.44 + 4.84}{6}$

$= \frac{10.84}{6}$

$\boxed{\sigma^2 = 1.8}$

* If variance is more than data spread is also more.

## 2. Standard Deviation:

| population S.D $(\sigma)$ | Sample SD $(s)$ |
|---|---|
| SD = $\sqrt{\text{Variance}}$ | SD = $\sqrt{\text{Variance}}$ |
| $\sigma = \sqrt{\sigma^2}$ | $s = \sqrt{s^2}$ |

Normalisation of the values because $\sigma^2 = 3.8$

$$\sqrt{\sigma^2} = \sqrt{3.8}$$

$$\sigma = 1.94$$

## 3. Range:

The difference between the minimum and maximum.

$$\text{Range} = \max - \min.$$

$\underline{\text{ex:}} \ \{1,2,2,3,4,5\}$

range $= 5-1$

$= 4.$

* A small range values are close to each other (less than variation).

* A large range values are widely spread out (more variable)

* Outlier can make the range misleading

    (outlier = out of the boundaries)

## Percentiles and Quartiles:

Percentile is a value below with a certain percentage of observation lie/came.

$\underline{\text{ex:}}$

$\{1,2,3,4,5, 5,5, 6,7,8,8,8,8,8, 9,9,10,11,11, 12\}$

Total = 20 = 100%

40% of to'records value is < 7.

formula:

percentile Rank of $x = \dfrac{\text{\# of values below } x}{n} \times 100.$

→In next page

Removing outliers from the data:

$[1,2,2,3,3,3,4,5,5,5,6,6,6,6,7,8,8,9,29]$

What is the % ranking of 10?

$$= \frac{16}{20} \times 100$$

$= 80\%$ of data value is below 10.

$= \frac{8}{20} \times 100 = 8 \times 5 = 40\%$ of records value is <+

What value exists at percentile ranking of 90%

$$\text{Value} = \frac{\text{percentile}}{100} \times (n+1)$$

$$= \frac{90}{100} \times 21$$

$$= 18.9$$

$$= 19 \quad \rightarrow \text{Index}$$

~~value = 19~~ $\boxed{\text{value} = 19}$

Five number Summary:

1. Minimum

2. First Quartile $(Q_1)$   $-25\%$

3. Second / Median Quartile $(Q_2)$  $-50\%$

4. Third Quartile $(75\%)$ $-Q_3$.

5. Maximum

Note: Choose these 5 numbers after removing outlier from the data.

Boundaries = [lower fence — Upper fence]

Lower fence $= Q_1 - 1.5 \ (IQR)$

Upper fence $= Q_3 + 1.5$

$$IQR = Q_3 - Q_1$$

$Q_1 = \frac{25}{100} \times 20$

$Q_1 = 5$          $s \rightarrow$ index

$Q_1 = 3$

$Q_3 = \frac{\frac{13}{75}}{\frac{100}{9}} \times 20$

$Q_{3} = \frac{15}{100} \times 20 = 15$

$Q_3 = 7$

$IQR = 7 - 3 = 4$

Lower fence $= Q_1 - Q 1.5 \leq IQR$

$\qquad = 3 - 1.5(4)$

$\qquad = -3$

Upper fence $= Q_3 + 1.56 \leq IQR$

$\qquad = 7 + 1.5(4)$

$\qquad = 7 + 6$

$\qquad = 13$

Anything $\beta -3$ & $13$ is considered as outliers.

[1, 2, 2, 2, 3,3,4; 5, 5, 5, 6, 6, 6, 6, 7, 6, 8, 9, 29]

1. min $\rightarrow 1$

2. $Q_1 \rightarrow 3$

3. $Q_2 \rightarrow 5$

4. $Q_3 \rightarrow 7$

5. $Q_4 \rightarrow 9$
   Max.