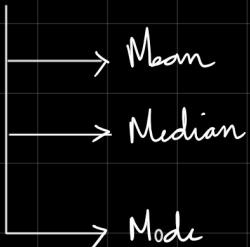


Intermediate Stats

Agenda

- ① Measure of central tendency.
- ② Measure of dispersion
- ③ Percentiles and Quantiles
- ④ 5 Number Summary (Box plots)

① Measure of Central Tendency



$$\{ 23, 24, 28, 29, 31, 32, 33 \}$$

↓

This data can be represented with a central value.

- ⇒ Measure of CT is a single value that attempts to describe a set of data by identifying the central position.

1.1 Mean (Average)

⇒ We calculate average

Population (N)

$$x_i = \{ 1, 2, 3, 4, 5 \}$$

Sample (n)

$$\text{Mean}(\mu) = \sum_{i=1}^N \frac{x_i}{N}$$

$$\text{Mean}(\bar{x}) = \sum_{i=1}^n \frac{x_i}{n}$$

$\left\{ \text{When to use mean?} \right\}$

- ⇒ Use when data is normally distributed.
- ⇒ mean is sensitive to outliers.

1.2 Median

$\left\{ \text{Used when data is skewed or contains outliers} \right\}$

$$x_i = \{1, 2, 3, 4, 5\}$$

$$\mu = 3$$

$$\begin{array}{l} \downarrow \\ \text{range} = 5 - 1 \\ = 4 \end{array}$$

$$x_j = \{1, 2, 3, 4, 5, 100\}$$

$$M = 19.16$$

$$\begin{array}{l} \downarrow \\ \text{range} = 100 - 1 \\ = 99 \end{array}$$

when range is large
there is a possibility of outlier

Finding Median

① Sort the data.

② Find the central number

→ If no. of elements are odd we take avg. of central number.

→ If no. of elements are even we take central number.

$$\underline{\text{Ex 1:}} \quad \{1, 2, \underbrace{3}_{\text{median}}, 4\}$$

$$\text{median} = \frac{2+3}{2} = 2.5$$

$$\underline{\text{Ex 2:}} \quad \{1, 2, 3, \cancel{4}, 5, 6, 7\}$$

$$\text{median} = 4$$

1.3 Mode \rightarrow Most frequently occurring data.

$$x_i = \{22, \cancel{24}, 26, 27, 100, \cancel{24}, 199, \cancel{24}, 27, 32, \cancel{24}, 33, 54, \cancel{24}\}$$

\circ Most repeated value in $x_i = 24$.

\circ Which is the Mode in this case.

$\left\{ \text{when to use mode?} \right\}$

$\left\{ \text{It is appropriate to use when data has two peaks (bimodal)} \right\}$

② Measure of Dispersion

- Variance (σ^2)
- Standard Deviation (σ)

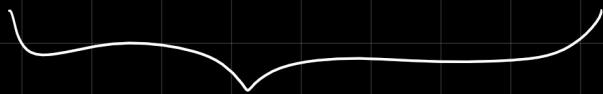
Consider samples,

$$x = \{1, 1, 2, 2, 4\}$$

$$y = \{0, 2, 2, 3, 3\}$$

$$\mu_{x_i} = \frac{1+1+2+2+4}{5} = 2$$

$$\mu_{y_i} = \frac{0+2+2+3+3}{5} = 2$$



Here, mean is same but we cannot say it follows the same distribution.

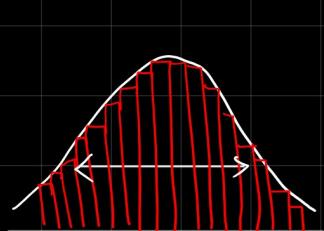
2-1 Variance (σ^2)

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

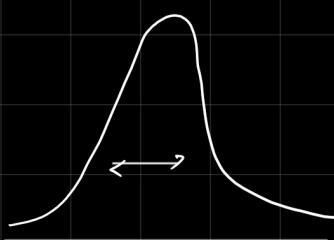
↳ tells us {spread of the distribution}

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

Bessel's correction
Degree of freedom



- Range is higher
- Variance is higher



- Range is lower.
- Variance is lower.

2.2 Population Standard Deviation

Population

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

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

Sample

$$s = s^2$$

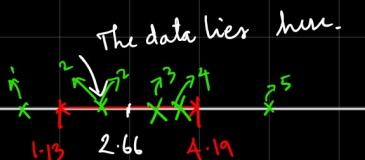
$$= \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}}$$

x	\bar{x}	$x - \bar{x}$	$(x - \bar{x})^2$
1	2.66	-1.66	2.75
2	2.66	-0.66	0.43
2	2.66	-0.66	0.43
3	2.66	0.44	0.19
4	2.66	1.44	2.07
5	2.66	2.44	5.95
2.66			11.82

$$s^2 = \frac{11.82}{6-1} = 2.36$$

$$s = \sqrt{2.36} = 1.53$$

$$\frac{2.66}{1.53} = 1.76$$



one standard deviation to right $(2.66 + 1.53)$

one s.d. to left $(2.66 - 1.53)$

$(\text{mean} - \text{variance})$

③ Percentiles and Quartiles

$$\text{Percentage} = \{1, 2, 3, 4, 5, 6, 7, 8\}$$

Ex: Percentage of even numbers = $\frac{\text{no. of even}}{\text{Total}} = \frac{4}{8} = \frac{1}{2} = 0.5$

Percentile :- A percentile is a value below which a certain percentage of observations lie.

Ex: A person's GATE score is 90 Percentile

Then it means, the person has scored better than 90% of students.

problem 1

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

What is the rank of 10?

\downarrow
[Data must be sorted before solving]

$$\text{Percentile of } x = \frac{\text{no. of values below } x}{\text{Sample}} * 100$$

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

$$= 80 \text{ percentile}$$

which means 10 is greater than entire distribution.

Problem 2

What is the value that entrts at the 25th percentile?

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

$$= \frac{25}{100} (21) \\ = \frac{21}{4} \\ = 5.25$$

$$= \frac{21}{4} \approx 5.25 \Rightarrow [5]$$

we have to choose
5th index value in data

{ When rank value is in decimal choose the corresponding
index value and index + 1 value }

$$\text{Ex: } 5.25 \Rightarrow \frac{5^{\text{th}} \text{ index} + 6^{\text{th}} \text{ index}}{2} = \frac{5+5}{2} = 5$$

(4) 5 Numbers Summary

- * Minimum
 - * First Quartile (Q_1) 25%.
 - * Median
 - * Third Quartile (Q_3) 75%.
 - * Maximum
- } helpful in understanding and removal of outliers.

$\left\{ \begin{array}{l} 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27 \end{array} \right\}$ → Outliers.

We can use 5# summary to prove / conclude that '27' is an outlier.

[Lower fence \leftrightarrow higher fence]

Lower fence $\Rightarrow Q_1 - 1.5 \text{ (IQR)}$

where $IQR = (Q_3 - Q_1)$

\downarrow \downarrow
75th 25th percentile

IQR \Rightarrow Inter Quartile Range.

higher fence $\Rightarrow Q_3 + 1.5 \text{ (IQR)}$

$$Q_1 = \left(25^{\text{th}} \text{ percentile} \right) = \frac{25}{100} \times (n+1) \quad \left\{ n \rightarrow \text{length of data} \right\}$$

$$= \frac{25}{100} \times (20) = (5^{\text{th}} \text{ index}) \Rightarrow \boxed{3}$$

$$Q_3 = \left(75^{\text{th}} \text{ Percentile} \right) = \frac{15}{100} \times (20) = \left(15^{\text{th}} \text{ index} \right) \Rightarrow 7$$

$$\text{IQR} = Q_3 - Q_1 = 7 - 3 = 4$$

$$\text{Lower fence} = 3 - 1.5(4) = -3$$

$$\text{Higher fence} = 7 + 1.5(4) = 13$$

It tells us that any value below -3 and above 13 are outliers.

{ i.e., in our case '27' is proved as outlier }

Data after removing outliers,

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

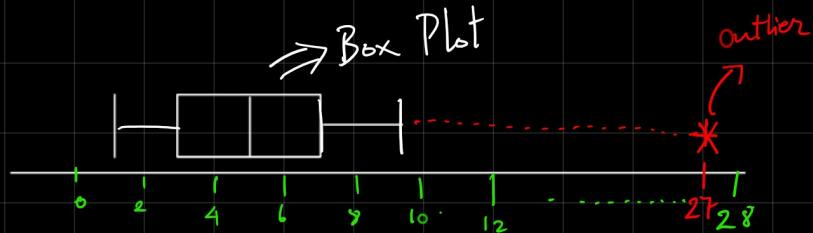
*) Minimum = 1

*) $Q_1 = 3$

*) Median = 5

*) $Q_3 = 7$

*) Maximum = 9



Q1, Median, Q3
0, 4, 10, 22