

Statistics – 5

- 1) **Data types**
- 2) **Categorical and Numerical**
- 3) **Qualitative and Quantitative**
- 4) **Continuous and discrete type**
- 5) **Levels of data**
- 6) **Nominal level**
- 7) **Ordinal level**
- 8) **interval level: It does not have true zero point**
- 9) **Ratio level: It has zero point**
- 10) **Population and sample**
- 11) **Inferential statistics: Will work on sample and estimate on population**
- 12) **Descriptive statistics: Analyse the data (analyse the population)**
- 13) **Frequency table**
- 14) **Bar chart**
- 15) **Relative frequency table**
- 16) **pie chart**
- 17) **frequency distribution table**
- 18) **histogram**
- 19) **Distribution plot**
- 20) **How to find the interval**
- 21) **How to choose class width**
- 22) **Central tendency**
- 22) **Mean – mode – Median**
- 23) **Median vs Mean**
- 24) **Positive skew : Right side skew: $\text{Mean} > \text{Median} > \text{Mode}$**
- 25) **Negative skew : Left side skew: $\text{Mean} < \text{Median} < \text{Mode}$**
- 26) **No skew: Normal distribution: $\text{Mean} = \text{Median} = \text{Mode}$**

In order to understand about data dispersions

we can approach in two ways

1) Center point analysis

2) Data flow analysis

1) Range :

- *in a data the marks values consider as raw data*
- *in that raw data we have a lowest value and Highest value*
- $\text{Range} = H - L$
- *The data lowest value has 1 mark , highest value = 100 marks*
- $\text{Range} = 100 - 1 = 99$

draw back:

Range will not tell about the middle points of the data

2) Mean deviation :

mean deviation tells about how a data point is deviated from mean

for example consider 5 data points : 1, 2, 3, 4, 5

The mean is = 3

How the data point (observation) 1 is deviated from mean value 3: $1 - 3 = -2$

$x_1: 1, x_2: 2, x_3: 3, x_4: 4, x_5: 5$

Mean is denoted with : \bar{x} (sample mean)

Mean is denoted with : $\bar{\mu}$ (population mean)

$$\bar{x} = \frac{x_1 + x_2 + x_3 + x_4 + x_5}{5}$$

$$\bar{x} = \frac{1 + 2 + 3 + 4 + 5}{5} = \frac{15}{5} = 3$$

How the data point (observation) $x_1: 1$ is deviated from mean value 3:

$$x_1 - \bar{x} = 1 - 3 = -2 \text{ (-indicates back side of mean)}$$

How the data point (observation) $x_2: 2$ is deviated from mean value 3:

$$x_2 - \bar{x} = 2 - 3 = -1$$

How the data point (observation) x_3 : 3 is deviated from mean value 3:

$$x_3 - \bar{x} = 3 - 3 = 0$$

How the data point (observation) x_1 : 1 is deviated from mean value 3:

$$x_4 - \bar{x} = 4 - 3 = 1 \text{ (+indicates ahead of mean)}$$

How the data point (observation) x_1 : 1 is deviated from mean value 3:

$$x_5 - \bar{x} = 5 - 3 = 2$$

Total deviation:

$$= -2 - 1 + 0 + 1 + 2$$

$$= (x_1 - \bar{x}) + (x_2 - \bar{x}) + (x_3 - \bar{x}) + (x_4 - \bar{x}) + (x_5 - \bar{x})$$

$$= \sum_{i=1}^5 (x_i - \bar{x})$$

$$\text{Mean deviation: } \frac{1}{5} * \sum_{i=1}^5 (x_i - \bar{x})$$

$$= (x_1 - \bar{x}) + (x_2 - \bar{x}) + (x_3 - \bar{x}) + \dots + (x_n - \bar{x})$$

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

$$\text{Mean deviation: } \frac{1}{N} * \sum_{i=1}^N (x_i - \bar{x})$$

Draw back:

What is the Mean deviation from above example:

$$= -2 - 1 + 0 + 1 + 2 = 0$$

But actually we are seeing the deviation, but the maths says No deviation

Chiru and vignesh are discussing why this is happend

they identified because of the neagtive values

we will convert negative to postive

3) Absolute Mean deviation :

- mean deviation has drawback of total deviation becomes zero due to positive and negative values
- will convert negative values to positive
- we choose a method $\text{mod} = |-5| = 5$ also $|5| = 5$

for example consider 5 data points : 1, 2, 3, 4, 5

The mean is = 3

How the data point (observation) 1 is deviated from mean value 3: $1 - 3 = -2$

$x_1: 1, x_2: 2, x_3: 3, x_4: 4, x_5: 5$

Mean is denoted with : \bar{x} (sample mean)

Mean is denoted with : $\bar{\mu}$ (population mean)

$$\bar{x} = \frac{x_1 + x_2 + x_3 + x_4 + x_5}{5}$$
$$\bar{x} = \frac{1 + 2 + 3 + 4 + 5}{5} = \frac{15}{5} = 3$$

How the data point (observation) $x_1: 1$ is deviated from mean value 3:

$$|x_1 - \bar{x}| = |1 - 3| = |-2| = 2$$

How the data point (observation) $x_2: 2$ is deviated from mean value 3:

$$|x_2 - \bar{x}| = |2 - 3| = |-1| = 1$$

How the data point (observation) $x_3: 3$ is deviated from mean value 3:

$$|x_3 - \bar{x}| = |3 - 3| = 0$$

How the data point (observation) $x_4: 4$ is deviated from mean value 3:

$$|x_4 - \bar{x}| = |4 - 3| = 1 \text{ (+indicates ahead of mean)}$$

How the data point (observation) $x_5: 5$ is deviated from mean value 3:

$$|x_5 - \bar{x}| = |5 - 3| = 2$$

Total deviation:

$$= |-2| + |-1| + |0| + |1| + |2|$$

$$= |(x_1 - \bar{x})| + |(x_2 - \bar{x})| + |(x_3 - \bar{x})| + |(x_4 - \bar{x})| + |(x_5 - \bar{x})|$$

$$= \sum_{i=1}^5 |(x_i - \bar{x})|$$

Absolute Mean deviation: $\frac{1}{5} * \sum_{i=1}^5 |(x_i - \bar{x})|$

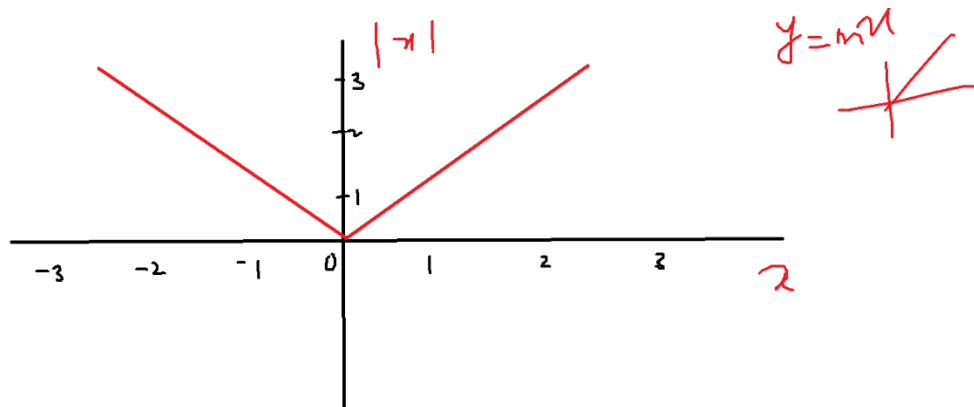
Absolute Mean deviation: $\frac{1}{N} * \sum_{i=1}^N |(x_i - \bar{x})|$

draw back:

$|x|$ graph

X	$ x $
-1	1
-2	2
-3	3
0	0
1	1
2	2
3	3

$(-1, 1), (-2, 2), (-3, 3), (0, 0), (1, 1), (2, 2), (3, 3)$



$|x|$ graph is not continuous, it does not have sharp point at 0

Maths says differentiation fails for not continuous equations

$|x|$ differentiation = $\frac{1}{|x|}$ or $\frac{x}{|x|}$

$\frac{1}{0}$ or $\frac{0}{0}$ = zero division error or undefined

4) Variance :

- mean deviation has drawback of total deviation becomes zero due to positive and negative values
- Absolute mean deviation has drawback of no smooth curve at point '0' so that the differentiation fails
ALL ML algorithms developed by Maths only
If any maths equation not holds the properties, we will not consider
- Our main goal is : To convert negative values to positive
- we choose square method = $(-5)^2 = 25$ also $(5)^2 = 25$

for example consider 5 data points : 1,2,3,4,5

The mean is = 3

$x_1: 1, x_2: 2, x_3: 3, x_4: 4, x_5: 5$

Mean is denoted with : \bar{x} (sample mean)

Mean is denoted with : $\bar{\mu}$ (population mean)

$$\bar{x} = \frac{x_1 + x_2 + x_3 + x_4 + x_5}{5}$$

$$\bar{x} = \frac{1 + 2 + 3 + 4 + 5}{5} = \frac{15}{5} = 3$$

How the data point (observation) $x_1: 1$ is deviated from mean value 3:

$$(x_1 - \bar{x})^2 = (1 - 3)^2 = (-2)^2 = 4$$

How the data point (observation) $x_2: 2$ is deviated from mean value 3:

$$(x_2 - \bar{x})^2 = (2 - 3)^2 = (-1)^2 = 1$$

How the data point (observation) $x_3: 3$ is deviated from mean value 3:

$$(x_3 - \bar{x})^2 = (3 - 3)^2 = (0)^2 = 0$$

How the data point (observation) $x_4: 4$ is deviated from mean value 3:

$$(x_4 - \bar{x})^2 = (4 - 3)^2 = (1)^2 = 1$$

How the data point (observation) $x_5: 5$ is deviated from mean value 3:

$$(x_5 - \bar{x})^2 = (5 - 3)^2 = (2)^2 = 4$$

Total deviation:

$$= 4 + 1 + 0 + 1 + 2$$

$$= (x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + (x_3 - \bar{x})^2 + (x_4 - \bar{x})^2 + (x_5 - \bar{x})^2$$

$$= \sum_{i=1}^5 (x_i - \bar{x})^2$$

$$\text{Variance: } \frac{1}{5} * \sum_{i=1}^5 (x_i - \bar{x})^2$$

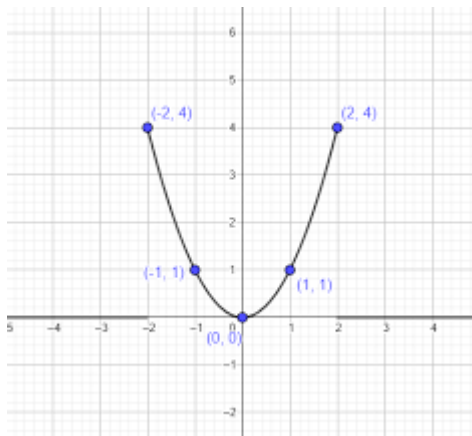
$$\text{Variance: } \frac{1}{N} * \sum_{i=1}^N (x_i - \bar{x})^2$$

draw back:

x^2 graph

X	X ²
-3	9
-2	4
-1	1
0	0
1	1
2	4
3	9

$(-3, 9), (-2, 4), (-1, 1), (0, 0), (1, 1), (2, 4), (3, 9)$



$$y = x^2 \quad \text{power} = 2 \quad \text{Non linear}$$

$$y = x \quad \text{power} = 1 \quad \text{Linear}$$

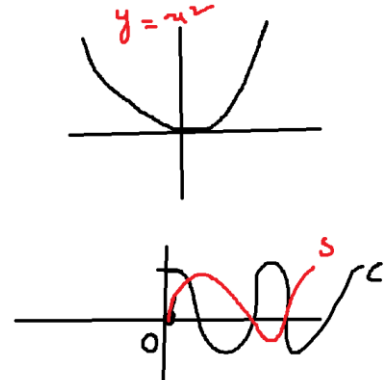
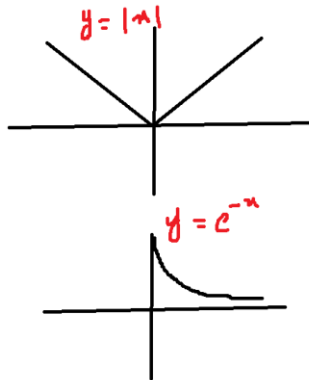
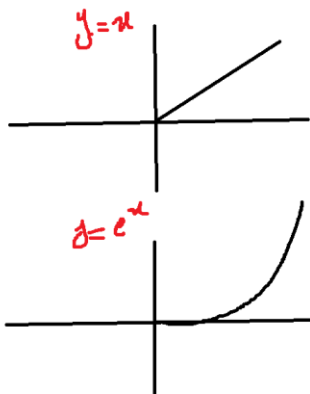
$$y = x^2$$

$$\frac{dy}{dx} = 2x$$

The curve has smooth every where $x = 0, \frac{dy}{dx} = 2 * 0 = 0$

So variance is a valid metric

Type equation here.



Drawback:

Variance: $\frac{1}{N} * \sum_{i=1}^N (x_i - \bar{x})^2$

Step – 1: Calculate mean \bar{x}

step – 2: calculate individual deviation: $x_i - \bar{x}$

Step – 3: Calculate square of the deviation = $(x_i - \bar{x})^2$

Step – 4: Calcuatue total deviation

Step – 5: Calculate the variance by dividing N

$x = km$	$x_i - \bar{x}$	$(x_i - \bar{x})^2$
1km	-2	4km ² : (1km-3km) ² = (-2km) ²
2km	-1	1km ²
3	0	0km ²
4	1	1km ²
5	2	4km ²
$\bar{x} = 3$		

- 1) When we do the variance the values are increasing
- 2) when we do the variance the units also mention as square units
- 3) This leads interpretation problem
- 4) In the above example the variance $= \frac{10}{5} km^2 = 2km^2$
- 5) In the realtime we does not have km^2 units , we can not explain in the proper way
- 6) Inorder to avoid the draw back we need to apply root on variance

Standard Deviation

It is denoted with σ

$$\sigma = \sqrt{\text{variance}}$$

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

Interpretation:

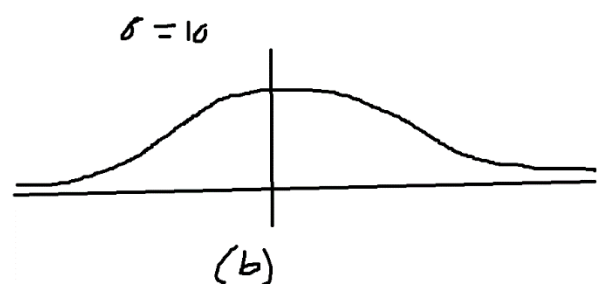
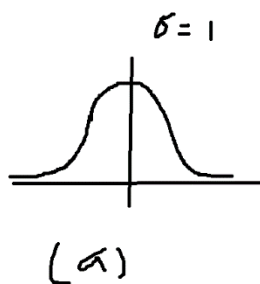
standard deviation will explain on of average a data point how much deviated from mean

for example: 4,5 var = $2km^2$

$$std = \sqrt{2km^2} = 1km$$

Each and every data point on of average 1km deviated from mean

on of average How much a data point is deviated from mean



Standard deviation is low: The data points are close to mean

Standard deviation is hight: The data points are far from mean

How much hyd is far from Nagpur: 400km

How much a data point is far from mean = standard deviation

- *Range* = $H - L$
- **Mean deviation:** $\frac{1}{N} * \sum_{i=1}^N (x_i - \bar{x})$
- **Absolute Mean deviation:** $\frac{1}{N} * \sum_{i=1}^N |x_i - \bar{x}|$
- **Variance:** $\frac{1}{N} * \sum_{i=1}^N (x_i - \bar{x})^2$
- $\sigma = sd = \sqrt{\frac{1}{N} * \sum_{i=1}^N (x_i - \bar{x})^2}$