



Annexure No.:

Descriptive statistic

- When you take into account the entire population in your data analysis it is called descriptive statistic.

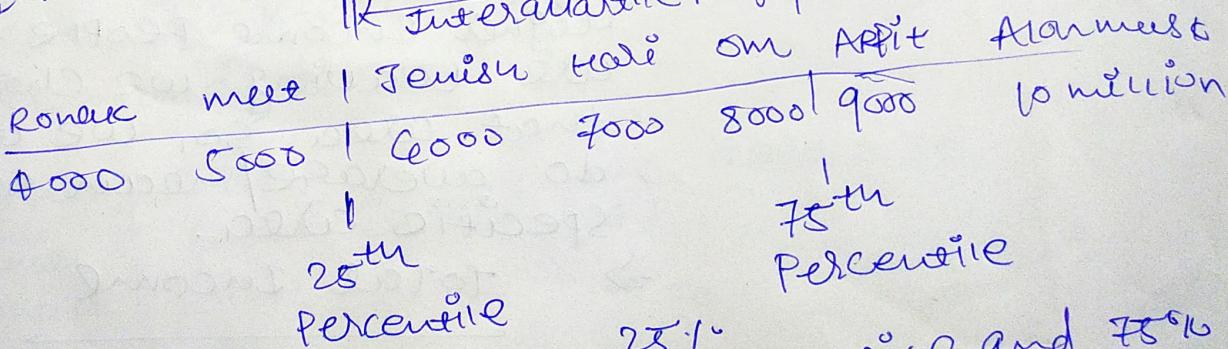
Inferential statistic

- When you take a sample for your process and then you see the analysis applies to the entire population, it is called inferential statistic.

outlier: outlier is a datapoint that is very different than rest of the data point.

IQR - Interquartile Range

1/4 Interquartile Range



Range b/w 25th and 75th Percentile
is Interquartile Range

How is Percentile used in data science

Semester

Income Data

outlier
Removal

General Data
Analysis

Ex example is
CISCE board
exam result
on percentile

iel is
mean
not

all
use
de

What is mean?

Roneak	4K
meet	5K
Keval	6K
Krunal	7K
John bengz	10 million

⇒ If, I want to open a premium car showroom in a specific area of a city, so as a Data scientist we will do analysis

⇒ So, our first step is to find area which have higher income people are leaving we choose to that area so, we will do average/mean of specific area.

⇒
$$\frac{\text{Total Income}}{\text{Total Person}}$$



Annexure No :

1) What is median.

2) In previous example one outlier is there because of which our mean value is higher so, mean is not good choice.

So, instead of this we can use median.

So, firstly we arrange in ascending order and we find middle element

⇒ What is mode?

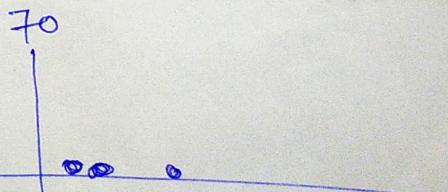
- mode here is median.

- mode means most frequently occurring value.

Q. What is mean absolute deviation and Standard deviation?

History test Avg for data ≥ 70

Name	Score	Avg (Score - Avg)
Monan	75	5
Andrea	72	2
Sofia	68	2
Joe	65	5
Vizat	67	3
Abdela	73	2



Points are nearby Average

Mean = 3.16

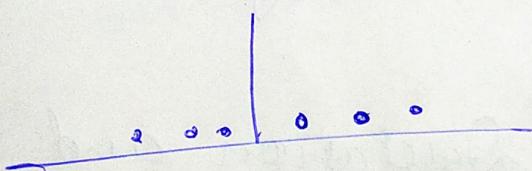
Math Test

(Score - Avg)

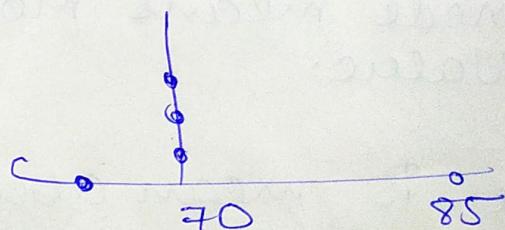
Name	Score	MAD	Average = 70
Mohan	83	23	
Andrea	96	26	
Sofia	43	27	
Joe	47	23	
Uzat	51	19 20 30	70
Abdul	50	20	
		MAD 23	date point all spread out then Average

By seeing this we can tell that Math test dataset is widespread than History test, so this is called ~~Avg~~ mean Absolute Deviation.

But If we have



$$MAD = 3.33$$



$$MAD = 3.33$$

so in this scenario we can't predict.

Annexure No.:

What is Standard Deviation?

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

$N = \text{size of population}$

Standard deviation is a square root of Variance

Name	Score	Abs Score (Score - Avg)	$(\text{Score} - \text{Avg})^2$
Mohan	75	5	25
Andrea	72	2	4
Sofia	68	2	4
Joe	65	5	25
Uday	67	3	9
Abdel	73	3	9

$$\text{Avg} = 12.66$$

$$\sqrt{\text{Avg}} = 3.55$$

standard deviation.

Name	test score	Abs (Score - Avg)	$(\text{Score} - \text{Avg})^2$
Mohan	83	13	169
Andrea	70	0	0
Sofia	70	0	0
Joe	63	7	49
Uday	70	0	0
Abdel	70	0	0

$$\text{Avg} = 36.33$$

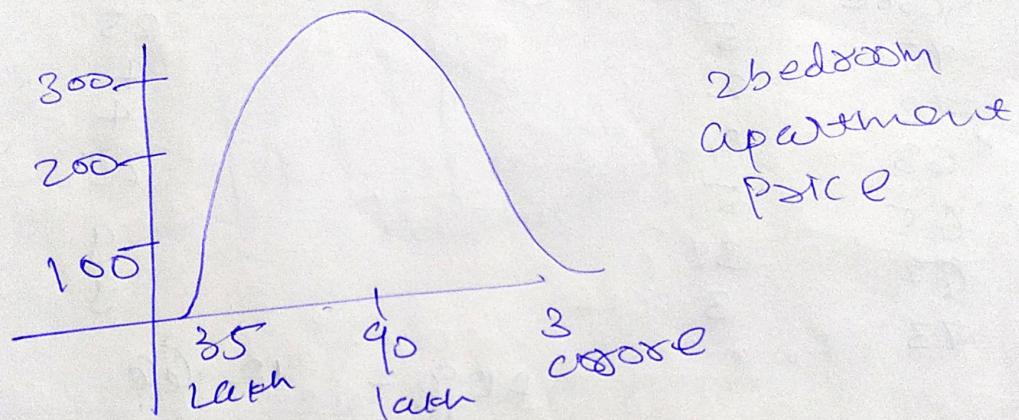
$$\sqrt{\text{Avg}} = 6.02$$

Page No.:

so, using standard deviation we can say that train test data are wide spread.

What is Normal Distribution

→ You will have most of your data sample around average value and you have some data sample that is far away from average. ~~left~~

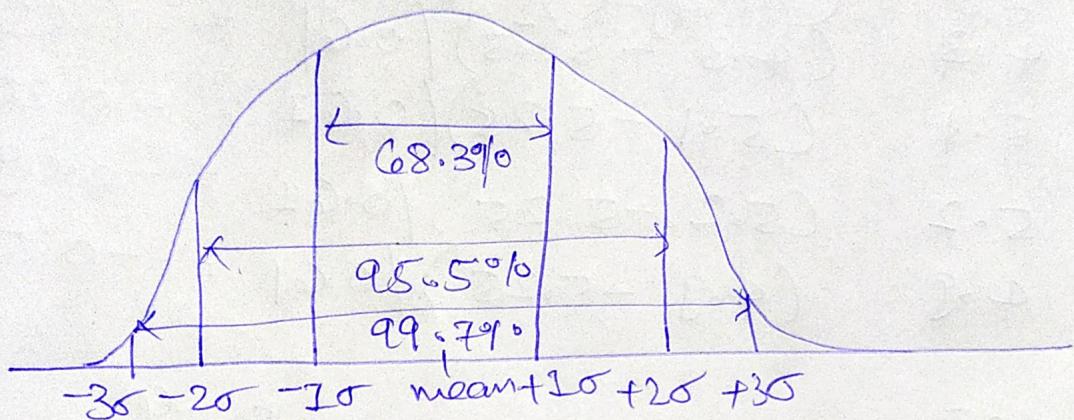


How can I use this in data analysis.

- outlier Removal.

⇒ What formulae do we use to remove outliers.

How standard deviation is used to remove standard deviation.



Any datapoint that is below -3σ and above 3σ is treated as a outlier.

What is Z score

How many standard deviation away a datapoint is from mean.

$$Z = \frac{x - \mu}{\sigma}$$

μ = mean

σ = standard deviation

Name	Height (ft)	z-score
Ronak	6.2	$(6.2 - 5.25) / 0.61 = 1.53$
meet	5.7	$(5.7 - 5.25) / 0.61 = 0.72$
Kevul	4.6	$(4.6 - 5.25) / 0.61 = -1.06$
Adpit	5.4	$(5.4 - 5.25) / 0.61 = 0.25$
Abhay	5.9	$(5.9 - 5.25) / 0.61 = 1.04$
Harshom	4.3	$(4.3 - 5.25) / 0.61 = -1.55$
Jenish	5.1	$(5.1 - 5.25) / 0.61 = -0.25$
minu	5.2	$(5.2 - 5.25) / 0.61 = -0.09$
maheen	4.9	$(4.9 - 5.25) / 0.61 = -0.58$

$$\text{Average} = 5.25$$

$$z = \frac{x - \bar{x}}{\sigma}$$

$$\text{Standard deviation} = 0.61$$

Now let's apply z-score

What is logarithm?

I have invested 5 \$
After some point of time it gives me
125 \$ so, how many years it
take

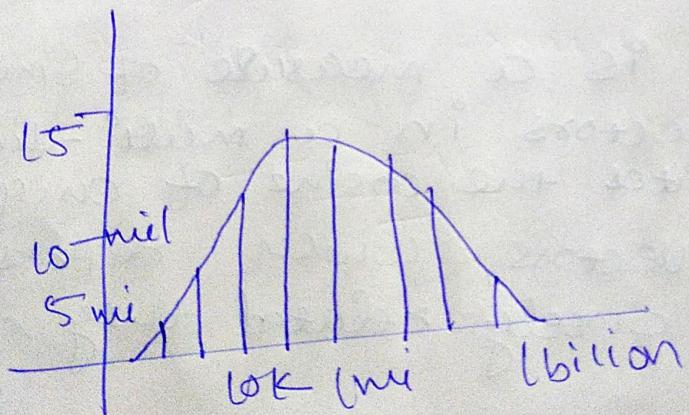
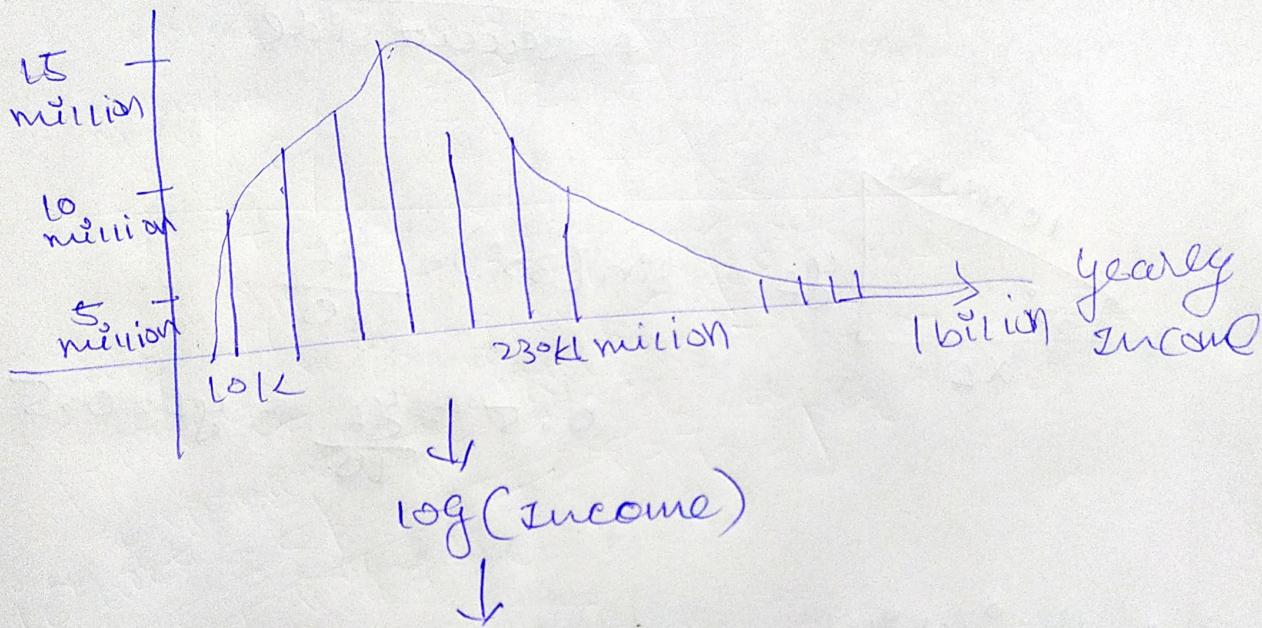
$$\log_5 125 \rightarrow 3$$

Logarithm is an inverse of an
exponent.

$$\text{Popular log is } \log_{10} 10 = 1$$

Annexure No :

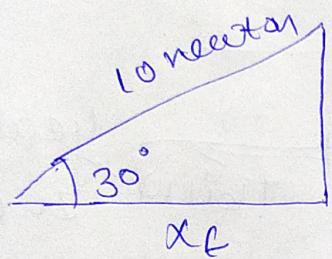
If you get a normal distribution by applying a log function to a dataset then dataset is log normally distributed.



$$\sin(\theta) = \frac{\text{opposite}}{\text{hypotenuse}}$$

$$\cos(\theta) = \frac{\text{adjacent}}{\text{hypotenuse}}$$

$$\tan(\theta) = \frac{\text{opposite side}}{\text{adjacent side}}$$



$$\sin(30^\circ) = \frac{y_f}{10}$$

$$0.5 = \frac{y_f}{10} = y_f = 0.5$$

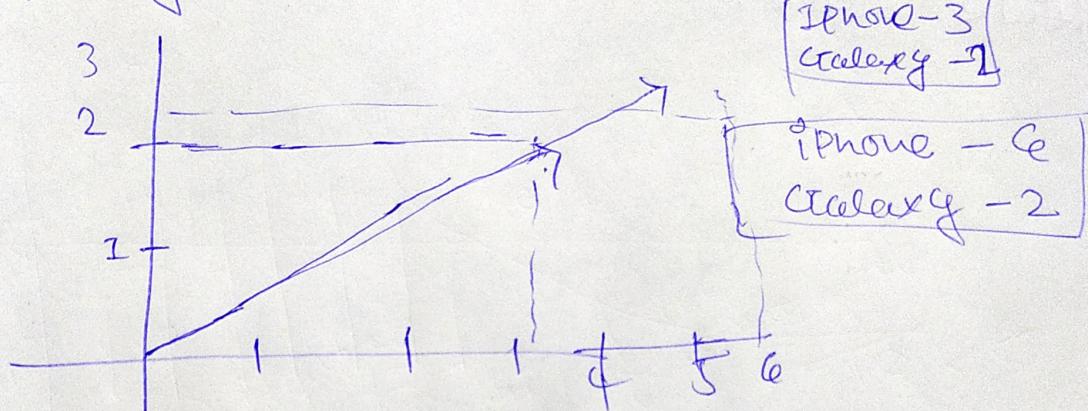
Cosine Similarity

Cosine similarity is a measure of similarity between two vectors in a multi-dimensional space. It calculates the cosine of angle between the vectors, which represents how similar or closely related they are.

Annexure No. :

Example

- document similarity of iPhone and Galaxy.



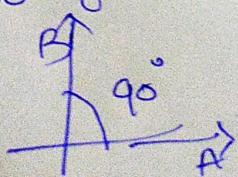
so here angle is 0 so document is similar.

$$\text{Cosine Similarity} = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|}$$

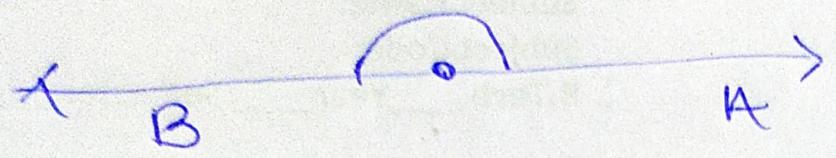
$$\text{Cosine Similarity} = \frac{\|\mathbf{A}\| \|\mathbf{B}\| \cos(\theta)}{\|\mathbf{A}\| \|\mathbf{B}\|}$$

$$\text{Cosine Similarity} = \frac{1}{\|\mathbf{A}\| \|\mathbf{B}\|} = \cos(\theta)$$

$$\text{Cosine Similarity} = 0$$



3.) Cosimilarity - 1



$\text{cosine distance} = 1 - \text{cosimilarity}$
cosine distance always represented in
positive space.