

Date: / /

## Introduction to Statistics

### Descriptive Stats



- ① Measure of central tendency
- ② measure of dispersion

Anything related to

summarizing the data.

Mean, Median, mode, SD, Variance

③ Gaussian Distribution

CHI-SQUARE

④ Normal ..

Hypothesis testing  
(P Value)

⑤ Binomial ..

z-table, t-table

### Statistics

↳ collecting, organising & analyzing data

↳ {Better Decision Making}

Data → Facts or pieces of information that can be measured.

### Types of Statistics

① Descriptive → consists of organizing & summarizing data

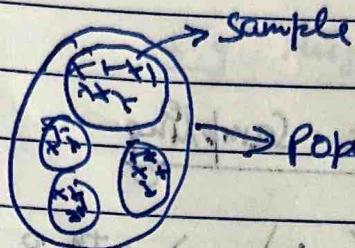
② Inferential → technique where we used the data that we have measured to form conclusion

Date: / /

## Population & Sample

Elections  $\rightarrow$  (n), UP

Exit Poll

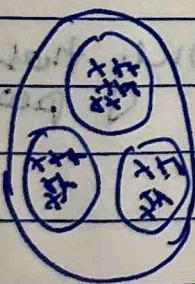


Population  $\rightarrow (N)$   
Sample  $\rightarrow (n)$

## Sampling Techniques

### ① Simple Random Sampling

Every member of the population has an equal chance of being selected for your sample ( $n$ ).



② Stratified Sampling :- where the population ( $N$ ) is split into non-overlapping groups (strata)

E.g. Gender

$\rightarrow$  Male

$\rightarrow$  Female

Date: / /

Age

(0-10)

(10-20)

(20-40)

(40-100)

Profession

### ③ Systematic Sampling

$(N) \rightarrow n^{\text{th}}$  Individual

Eg. Mall  $\rightarrow$  Survey (Covid)

$\hookrightarrow 8^{\text{th}}$  person  $\rightarrow$  survey

It means every 8th person

### ④ Convenience Sampling

Survey

Only those people

$\hookrightarrow$  Data Science

### Variables

A variable? P: a property that can take on any value.

- a) Quantitative Variables  $\rightarrow$  Measured Numerically
- b) Qualitative / Categorical Variable.

Date: / /

## Quantitative

↓  
Discrete Variables

E.g. Whole Number

Total child in family  
2, 3, 4, 5

↓  
Continuous Variables

E.g. Height: 172.5, 162.2 cm

Variable: Measurement for Scales

4 types of Measured Variable

- ① Nominal → Class → Gender, colours, type of flowers
- ② Ordinal → Order of the data matters, not values.
- ③ Interval → Order matters, values also matters, natural zero is not part
- ④ Ratio

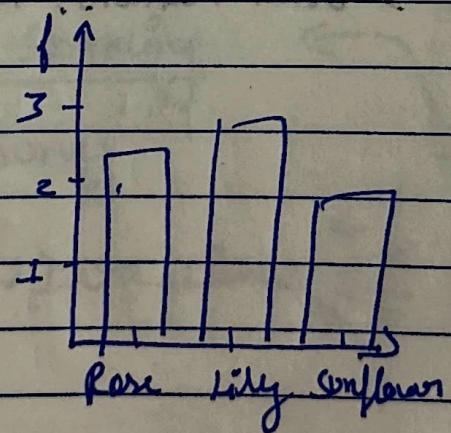
Date: / /

## Frequency Distribution

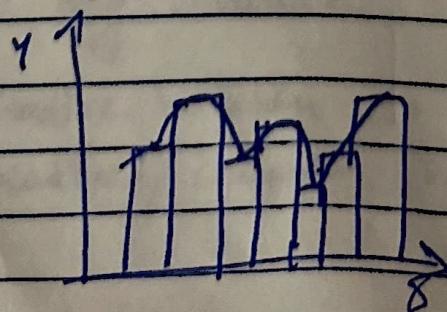
Sample Dataset := Rose, Lilly, Sunflower, Rose, Lilly

flower	frequency	relative frequency
Rose	3	3
Lilly	5	7
Sunflower	2	4

### Bar Graph



Histogram = Continuous



## Arithmetic Mean for population & Sample

Mean (Average)

Population ( $N$ )

$$K = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\}$$

$$\mu = \frac{\sum_{i=1}^N n_i}{N}$$

Sample ( $n$ )

$$\bar{x} = \frac{\sum_{i=1}^n n_i}{n}$$

$$= \underline{\underline{3.2}}$$

$$= \frac{1+1+2+2+3+3+4+5+5+6}{10}$$

$$= \frac{28}{10} = \underline{\underline{2.8}} = \frac{32}{10} = \underline{\underline{3.2}}$$

Central Tendency → refers to the measure used to determine the centre of the distribution of data.

- Mean
- Median
- Mode

$$\{1, 1, 2, 2, 3, 3, 4, 5, 5, 6, \boxed{100}\}$$

outliers

$$\frac{1}{2} - \frac{6+7}{102} = \frac{12}{26} = \frac{6+7}{23} = \frac{13}{23} = \frac{11}{2} = 5.5$$

11 Date: / 15-5

### Median

{ 1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100 } odd number

$$\frac{11}{2} = 5.5$$

① Sort the Numbers

In Case Odd Number =  $\frac{(N+1)}{2}$   $\text{th element}$  =  $\frac{11}{2} = 5$ th element

In Case Even Number =  $\frac{(N)}{2}$ th element +  $\frac{(N+1)}{2}$ th element =  $\frac{3}{2}$

Summary (Summary)

$$\text{Median} = \frac{3+4}{2} = 3.5$$

(Median) : { Median works well with outliers }

Mode = { 1, 2, 2, 3, 4, 5, 6, 6, 6, 7, 8, 100, 100, 100, 100 }

= { Most Frequent Element }

Petals

2.0

3.0

4.0

Mode works in categorical data

Type of flowers      petal length      petal width

Ran

Lily

Sunflower

Missing Value  $\rightarrow$  Most frequent occurring element

~~1, 2, 3, 5~~  
~~80%~~

Date: / /

## Measures of Dispersion $\rightarrow$ Spread

a) Variance

b) Standard Deviation

Population Variance

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

Sample Variance

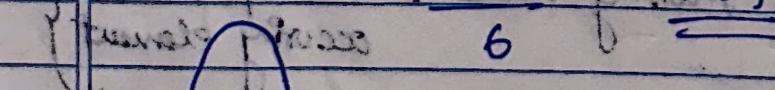
$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

k	$x_i$	$x_i - \bar{x}$	$(x_i - \bar{x})^2$
1	2.83	-1.83	3.34
2	2.83	-0.83	0.6889
3	2.83	0.17	0.03
4	2.83	1.17	1.37
5	2.83	2.17	4.71
M	2.83		10.84

Average from  $\underline{10.84} = 11.8$

Sample size

6



Low Variance

More Variance

Data is more dispersed

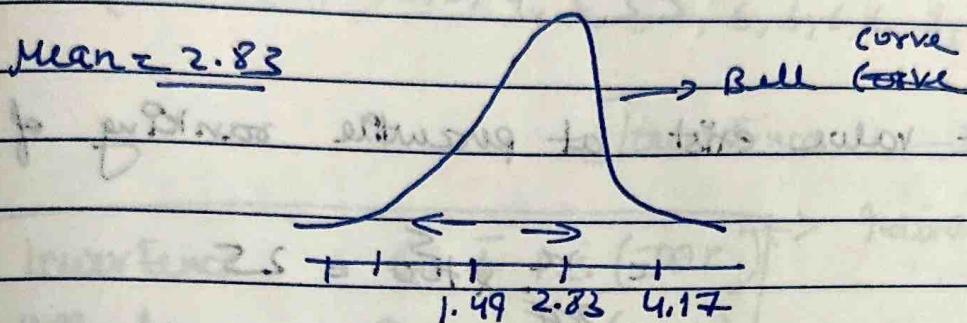
Variance means more spread

$$2.2 = \frac{14}{4} \rightarrow \frac{14}{4} = 1+2+3+4 \\ \text{Date: } \frac{10}{4} = \underline{\underline{2.5}}$$

SD

$$\sigma = \sqrt{\text{Variance}} = \sqrt{1.81} = 1.345$$

$$\text{Mean} = 2.83$$



$$\frac{2.83}{4.17} + 1.34$$

→ Variance tells how the data spread

→ SD tells the range of the data

## Percentiles & Quartiles [Find Outliers]

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

Ig

Dataset : 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 9, 9, 10, 11, 11, 12

Q What is the percentile ranking of 10?

Percentile Rank of  $x = \frac{\# \text{ of values below } x}{n} \times 100$

$n = \text{sample size}$

$$= \frac{16}{20} \times 100 = \underline{\underline{80\%}}$$

80% of the entire distribution is less than 10

$$\frac{5+6}{2} = \frac{11}{2} = \underline{\underline{5.5}}$$

Date: / /

$$= \frac{17}{20} \times \frac{5}{100} = \underline{\underline{85\%}}$$

Q

What value exist at percentile ranking of 25%?

$$= \frac{n}{20} \times \frac{5}{100} = 25$$

$$5n = 25$$

$$n = 5$$

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

$$= \frac{25}{100} \times 21 = 5.25 \rightarrow \begin{matrix} \text{Index} \\ \text{Position} \end{matrix}$$

$$\frac{75}{100} \times 21 = 15.75 \rightarrow \underline{\underline{16}}$$

## Five Number Summary

- ① Minimum
- ② Q<sub>1</sub> (Quartile)
- ③ Median
- ④ Q<sub>3</sub> (Quartile 3)
- ⑤ Max

Date: / /

## Removing the Outliers

Q: 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 9, (27)

↑ outlier

[Lower Fence  $\longleftrightarrow$  Higher Fence]

$$\text{Lower Fence} = Q_1 - 1.5 \text{ (IQR)} \rightarrow \text{Interquartile Range}$$

$$\text{Upper fence} = Q_3 + 1.5 \text{ (IQR)}$$

$$\text{IQR} = Q_3 - Q_1$$

$$Q_1 = (25\%)$$

$$Q_3 = (75\%)$$

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

$$\frac{25}{100} \times 20 = 5 \rightarrow \text{Index} \\ = \underline{\underline{3}}$$

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

$$\therefore 3 - 6 = \boxed{-3}$$

$$\frac{75}{100} \times 20 = 15 \rightarrow \text{Index} \\ = \underline{\underline{7}}$$

$$\text{Higher Fence} = Q_3 + 1.5 \text{ (IQR)}$$

$$= \boxed{13}$$

$$[-3 \longleftrightarrow 13]$$

anything lower than 3

anything greater than 13

$$\frac{18}{2} = 9 \leftarrow S \quad \text{Date: } 1 / 1$$

## Remaining Data

1, 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, ~~2~~

$$\text{Min} = 1$$

$$Q_1 = 3$$

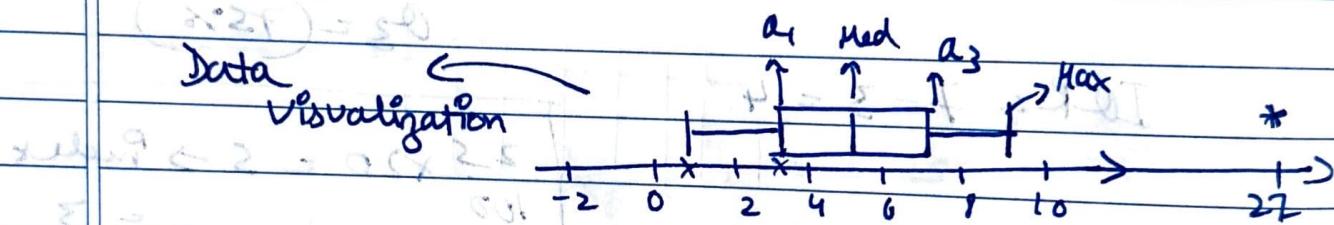
$$\text{Median} = 5$$

$$Q_3 = 7$$

$$\text{Max} = 9$$

5 Number Series

Box Plot



Data Visualization

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

Basal correction

Degree of freedom

Why Sample Variance is  $n-1$ ?

Population ( $N$ )

$$\mu = \frac{\sum_{i=1}^N x_i}{N}$$

Sample ( $n$ )

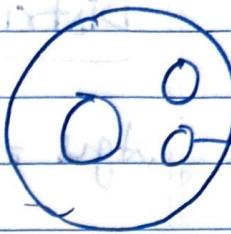
$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

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

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

Date: / /

Ags



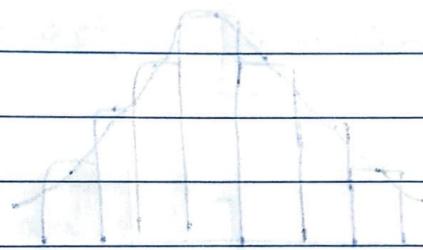
$$\bar{x} \rightarrow \mu$$



all are in even side of mean

free spirograph

standard deviation of all variables



## ① Distribution

Probability density function

↳ Normal (gaussian)

↳ Standard Normal

↳ Z score

↳ Log Normal Distribution

↳ Bernoulli Distribution

↳ Binomial Distribution

discrete probability

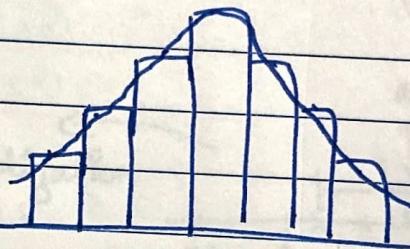
Date: / /

## Distribution

$$\text{Ages} = \{ 24, 26, 27, 28, 30, 32, \dots \}$$

↳ we need to focus how we see the data  
in visualize way

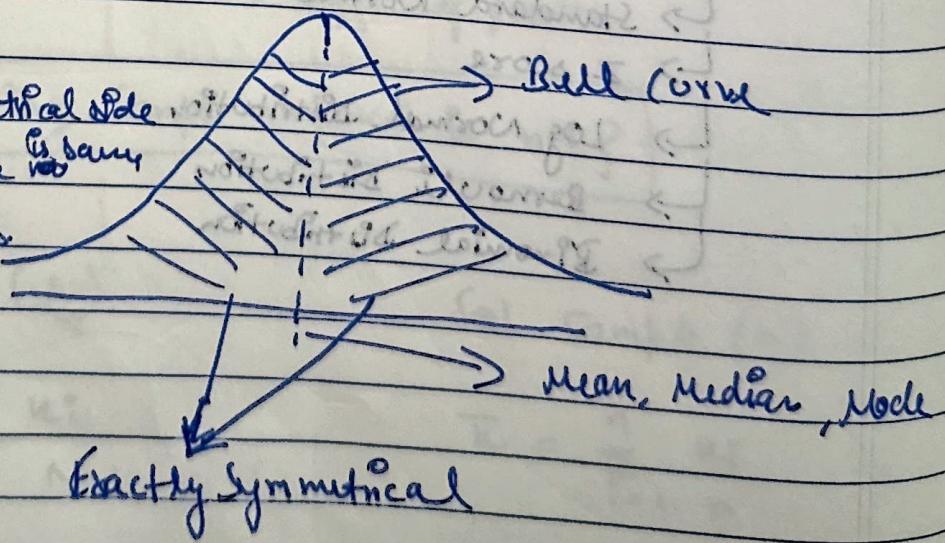
↳ Our objective is to visualize the data.



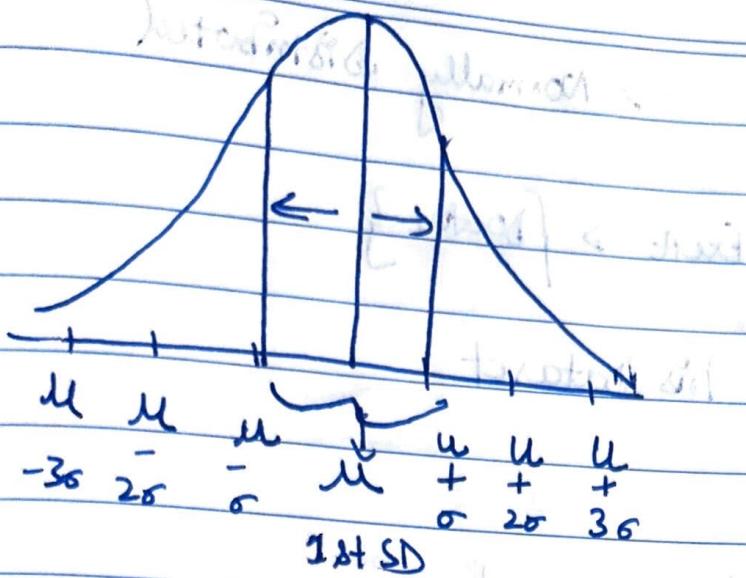
## Gaussian / Normal Distribution

### Property

- Same symmetrical side,
- Amount of data <sup>is same</sup> in both parts.



Date: / /



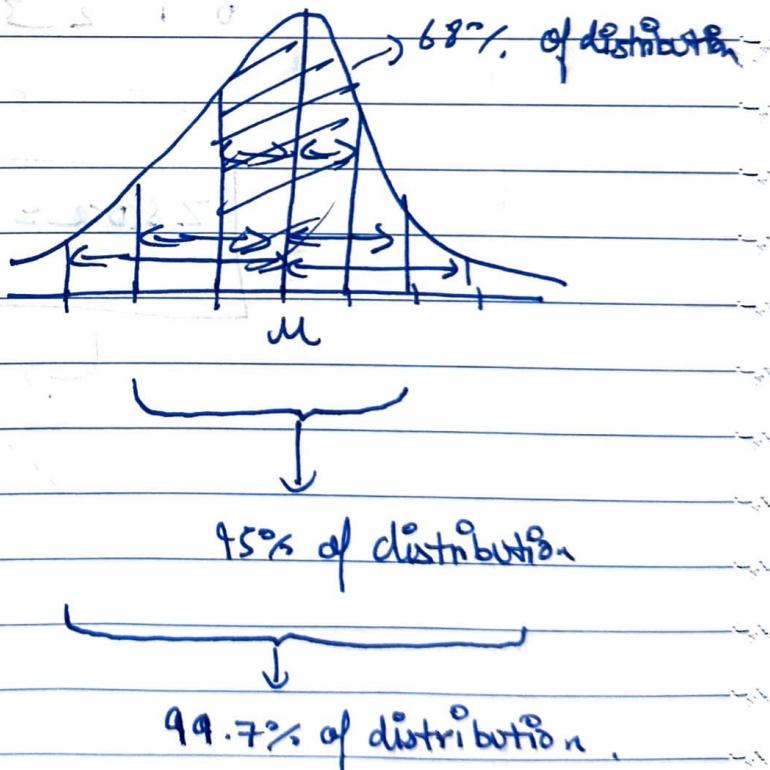
### Empirical Formula

68

~~68 - 95 - 99.7~~

~~68 - 95 - 99.7% Rule~~

DataSet [100 Data points]



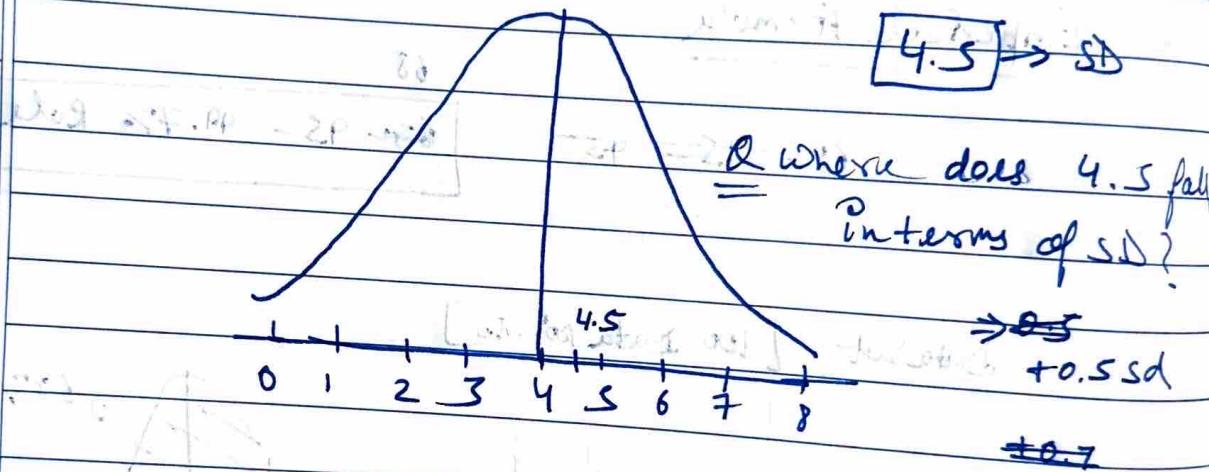
Date: / /

E.g. Height  $\rightarrow$  Normally Distributed

↓  
Domain Extent  $\rightarrow \{ \text{Doctor} \}$

Weight, 9ns dataset

E.g.  $\mu = 4$ ,  $\sigma = 1$



$$Z\text{ score} = \frac{x_i - \mu}{\sigma}$$

will basically help to find out how much standard deviation a value is away from its mean.

$$\frac{1-4}{1} = -3 \quad -\frac{3+0}{1} = -3$$

Date: / /

$$\mu = 4$$

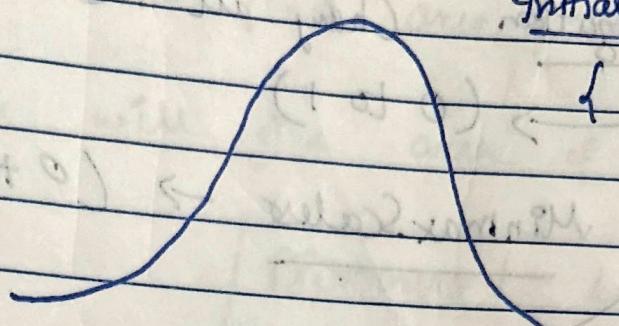
$$\sigma = 1$$

Initially distribution

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

Apply

$$Z \text{ Score} = \frac{x_i - \mu}{\sigma}$$



$$1, 2, 3, 4, 5, 6, 7$$

$$-3\sigma, -2\sigma, -1\sigma, 0, +1\sigma, +2\sigma, +3\sigma$$

$$\{-3, -2, -1, 0, 1, 2, 3\}$$

$$\{1, 2, 3, 4, 5, 6, 7\} \rightarrow \text{Normal distribution}$$

$$Z \text{ score}$$

$$\{-3, -2, -1, 0, 1, 2, 3\} \rightarrow \text{Standard Normal distribution}$$

$$(\mu=0, \sigma=1)$$

satisfying

### Practical Application

Dataset (years)

Age

24

25

26

27

(Rs)

Salary

40k

80k

60k

70k

(kg)

Weight

70

80

55

45

To convert the dataset into this →

$$\mu=0, \sigma=1$$

Standardization

Z score

Apply

Z score

& convert into

SND

-1.02

0:25 s.

~~0.01~~ ~~0.05~~

0.02

0.26

0.2 Date 0.96

Normalization Change all the values in 0 to 1

$\rightarrow$  to 1)

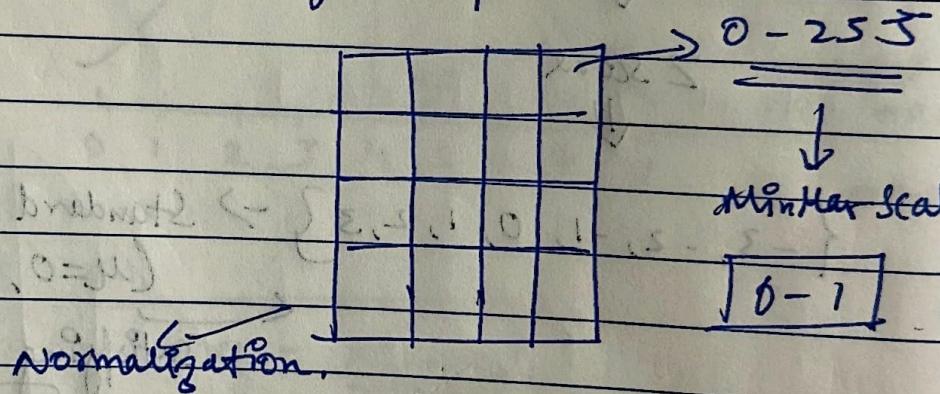
$\leftarrow$  (0 to 1)

Minmax Scales  $\rightarrow$  (0 to 1)

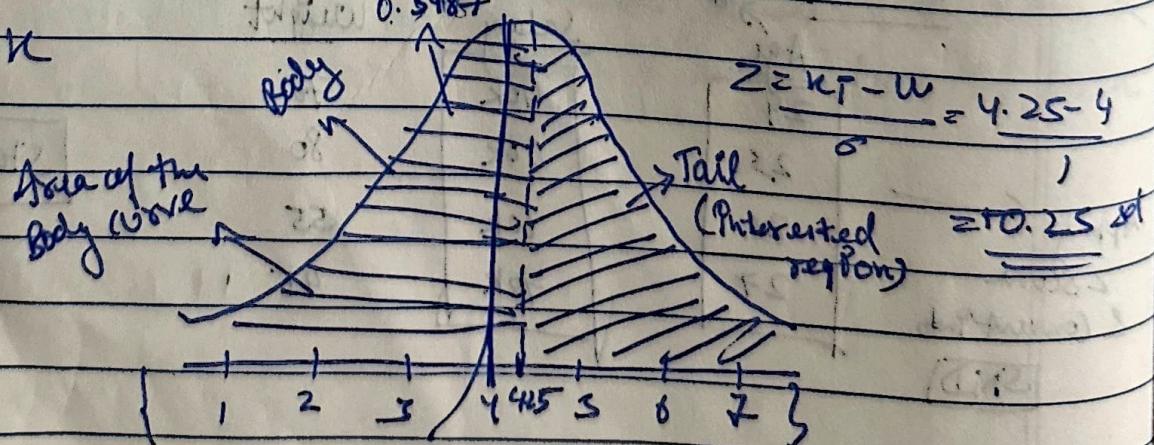
It is a process where you can define the lower bound & upper bound & you can convert the data between them

E.g.

CNN  $\rightarrow$  Image Classification



## Stats Interview Question



Q What percentage of scores fall above 4.25?

$$\frac{200}{\sigma \text{ or } 25} = \underline{\underline{25}} \quad -1.02 \quad -1.03 \quad \frac{4}{7} \times 100 = \frac{400}{7} = 57 \quad \frac{3}{7} \times 100 = \frac{300}{7} = 42.8 \quad \text{Date: } / /$$

→ Z scores helps to find area of the body curve

→ Z tables will give area of the body curve

$$\text{Right Area} = 1 - \text{Left Area}$$

$$= 1 - 0.5987$$

$$\downarrow$$

40%

In India, ~~Axes~~

$$\mu = 100$$

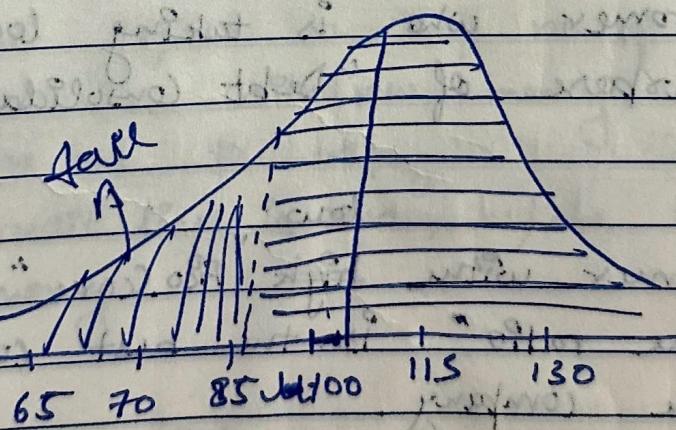
$$\sigma = 15$$

$$x_i^o = 85$$

$$Z = \frac{x_i^o - \mu}{\sigma}$$

$$Z = \frac{x_i^o - \mu}{\sigma}$$

$$Z = \frac{85 - 100}{15} = -1$$



$$\text{Left} = 1 - \text{Right Area}$$

Date: / /

## Practical

### Google Colab

```
import seaborn as sns  
import numpy as np  
import matplotlib.pyplot as plt  
import statistics
```

# mean, median, mode

```
df = sns.load_dataset('tips')
```

```
df.head()
```

```
np.mean(df['total_bill'])
```

```
np.median(df['total_bill'])
```

```
statistics.mode(df['total_bill'])
```

```
sns.boxplot(df['total_bill'])
```

```
sns.histplot(df['total_bill'], kde=True)
```

↳ probability density function.

```
sns.countplot(df['species'])
```

↳ count the data grouping by categories

nb. percentile (df['<sup>new</sup>mpg'][25, 75])

## #outliers

dataset = [your own dataset], ~~new~~

outliers = []

def detect\_outliers(data):

threshold = 3

mean = np.mean(data)

std = np.std(data)

for i in data:

z-score = (i - mean) / std

if np.abs(z-score) > threshold:

outliers.append(z-score)

return outliers

$$IQR = Q_3 - Q_1$$

Date: / /

#  $IQR = \text{np.percentile}(df['A'], [75, 25])$  will work in

→ sort the data

↓

calculate  $Q_1$  &  $Q_3$

↓

$$IQR = Q_3 - Q_1$$

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

$$\text{Upper fence} = Q_3 + 1.5(IQR)$$

(I = outlier)

dataset = sorted(dataset)

$$Q_1, Q_3 = \text{np.percentile}(dataset, [25, 75])$$

# Find the lower fence & higher fence

$$IQR = Q_3 - Q_1$$

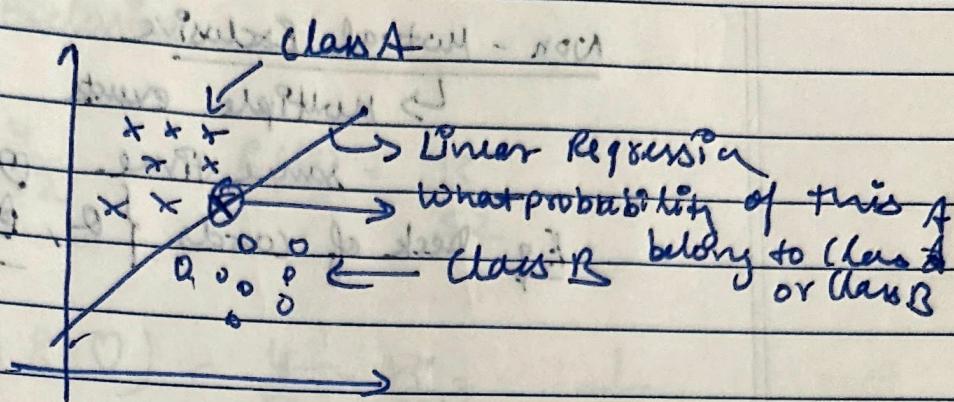
$$\text{lower-fence} = Q_1 - (1.5 * IQR)$$

$$\text{higher-fence} = Q_3 + (1.5 * IQR)$$

(outlier)  $\rightarrow$  outlier

Date: / /

## Probability



It is a measure of the likelihood of an Event

E.g.: Flipping a die {1, 2, 3, 4, 5, 6}

$$Pr(6) = \frac{1}{6} = \frac{\# \text{ of ways an event can occur}}{\# \text{ of possible outcomes}}$$

Toss a coin {H, T}

$$(3)^1 + (4)^1 = (2 \times 3) = 6$$

$$Pr(H) = \frac{1}{2}$$

$$\frac{1}{2} = \frac{1}{2} + \frac{1}{2} =$$

### Addition Rule ('or')

#### Mutually Exclusive Events

Two events are mutually exclusive if they ~~are~~ cannot occur at the same time  
e.g. rolling a dice {1, 2, 3, 4, 5, 6}

and not overlapping from a probability  
of probability adding up to one  
true

Non - Mutual Exclusive

↳ Multiple events can occur at the same time.

E.g. Deck of cards f.e., {A, B}

Q. If I toss a coin, what is the probability of the coin landing on heads or tails?

A. Non Mutual Exclusive = 1 = (1) 1/2

P(H or Tails)

$$P(A \text{ or } B) = P(A) + P(B)$$

$$= \frac{1}{2} + \frac{1}{2} = \underline{\underline{1}}$$

Q. Roll a dice

$$P(1 \text{ or } 3 \text{ or } 6) = \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{3}{6} = \frac{1}{2}$$

Q. Picking a card randomly from the deck.  
What is the probability of choosing a card that is queen or a ~~king~~ <sup>ace</sup> heart?

Date: / /

Any

### Non-Mutually Exclusive

$$P(A) = \frac{4}{52} \quad P(B) = \frac{13}{52}$$

$$P(A \text{ and } B) = \frac{4}{52} \times \frac{13}{52} = \frac{1}{52}$$

### Addition Rule for Non-Mutually Exclusive Events

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

$$P(A \text{ or } B) = \frac{4}{52} + \frac{13}{52} - \frac{1}{52}$$

$$= \frac{16}{52} \quad \text{Total number of cards with 4 or 13}$$

and not both

(2)

### Multiplication Rule

{Independent Events}

e.g. Rolling a die {1, 2, 3, 4, 5, 6}

1, 1, 2, 2, 3, ... → one event <sup>is</sup> independent to another event.

{Dependent Events}

