

Statistics

* What is Stats?

- It is a branch of maths which deals with collecting, organizing, interpreting & presenting data.
- Data summarizing / ~~explanation~~, hypothesis testing to make some decision.

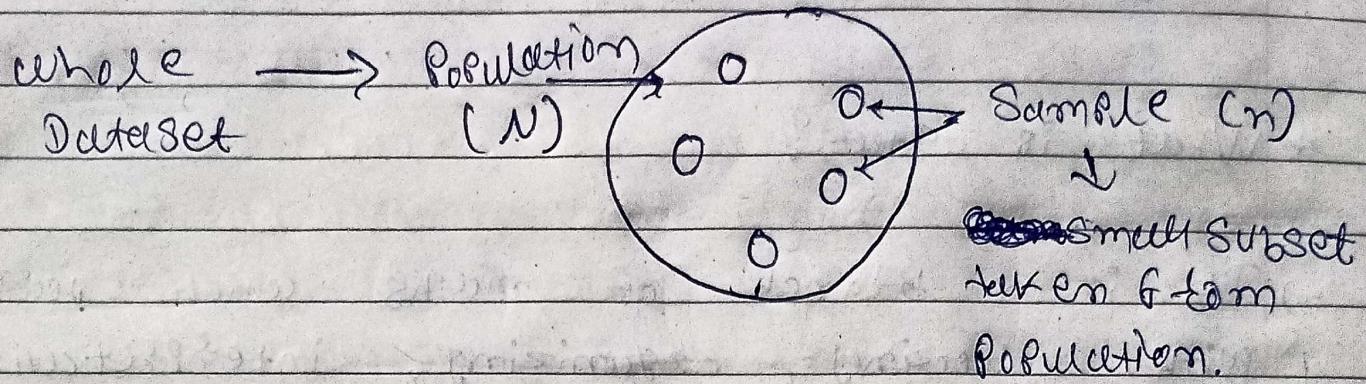
* What is Data?

- Facts or pieces of information that can be measured.
- DIKW (Data (Raw facts), Information, knowledge, wisdom (insights)).

* Types of Stats:

- 1). Descriptive → Organization & summarization.
- 2). Inferential. → Techniques using which we can conclude something from data or inference.

* Sample & Population :-



* Sampling Techniques :-

1. Simple Random Sampling :-

→ Every member of population have an equal chance of getting selected in your sample.

2. Stratified Sampling :-

→ Where we split the population into non-overlapping groups (strata).

→ For Gender → Male
→ Female

Age → 0 - 20
→ 21 - 40
→ 41 - 60

3. Systematic Sampling :-

→ Surveying on (N^{th}) individual

4. Convenience Sampling

* Variable :-

→ It is a property that can take any value.

→ Types of variable :-

1. Quantitative :- Measurable (numerical)

2. Qualitative :- Categorical value.

→ Types of quantitative variable :-

1). Discrete :- Whole no

2). Continuous :- Decimal value

* Variable Measurement Scale

1. Nominal :- Categorical - No order

2. Ordinal :- Ordered (order matters)

3. Interval :- [No zero/absolute Point], E.g. Temperature

4. Ratio:- Ex. 1 - Height / weight / Marks

* Frequency :-

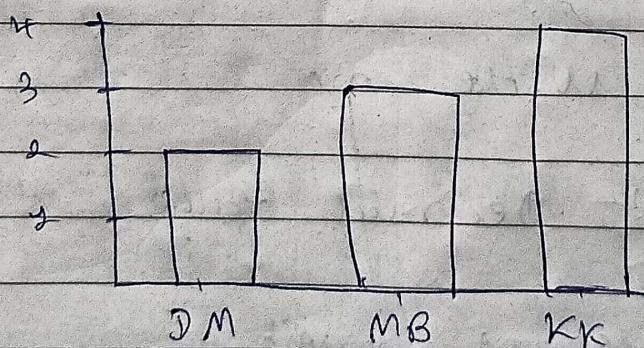
→ It means how many times value ~~repeats~~ occurs.

→ Ex 1 -

2. [DM, MB, KK, KK, KK, MB, DM, KK, MB]

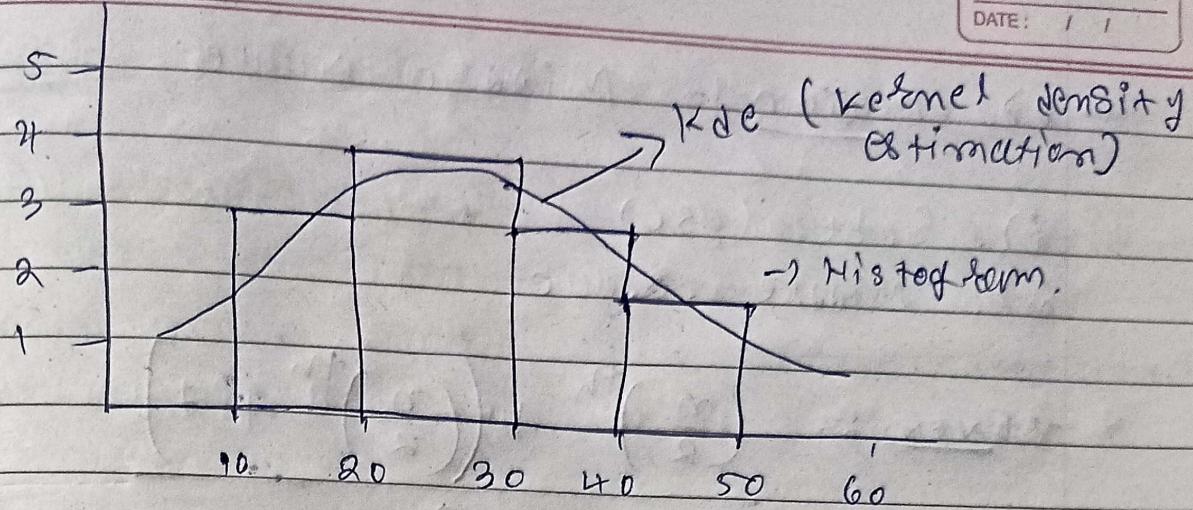
<u>Chocolates</u>	<u>Frequency</u>	<u>Cumulative freq.</u>
DM	2	2
MB	3	5
KK	4	9

Bar Graph



2. Age → [11, 12, 15, 25, 26, 28, 24, 33, 37, 39, 42, 44]

Age Range	Frequency
10 - 20	3
20 - 30	4
30 - 40	3
40 - 50	2



* Descriptive Statistics :-

→ Descriptive stats use methods for summarizing & understanding the key characteristics of data.

1. Measure of Central Tendency:-

→ Measure of CT describes the "center" of typical value of a dataset. They provide a summary of the data with a single value that represent the dataset as a whole.

2. Mean :- Avg → refers to the measure used to determine the center of distribution.

$\mu \rightarrow \text{Pop avg}$

$\bar{x} \rightarrow \text{sample avg}$

$$\mu = \frac{\sum xi}{N}$$

$$\bar{x} = \frac{\sum xi}{n}$$

2) Median :- Middle value

i). Sort (asc)

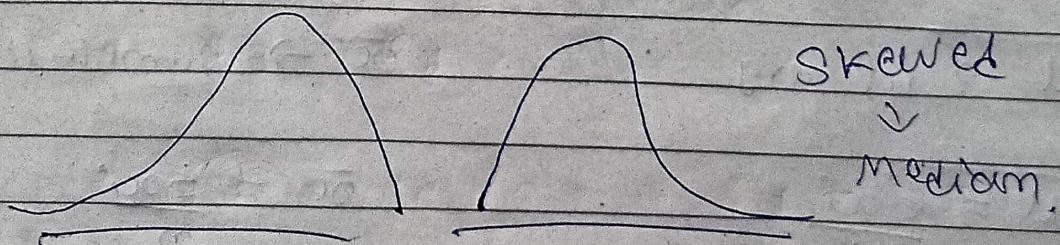
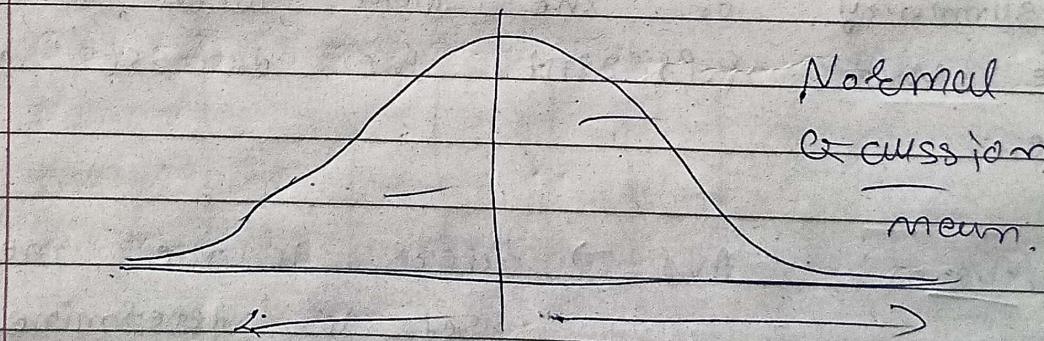
ii). Middle value

→ Even :-
$$\frac{\left(\frac{n}{2}\right)^{th} + \left(\left(\frac{n}{2}\right)^{th} + 1\right)}{2}$$

→ Odd :-
$$\left(\frac{n+1}{2}\right)^{th}$$

→ Outliers :- a datapoint(s) who doesn't follow the pattern at end of the data / or are extreme points.

3) Mode :- Most frequency value



2. Measure of Dispersion :- Spread

→ MOD describes the spread of variable of dataset. They indicates how much individual data points differs from the central value.

3). Variance:-

→ Variance measure how far each data pt in the dataset is from the mean.

→ High variance → more spread

→ Low " → less "

pop

sum

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

$$s^2 = \frac{\sum (\bar{x}_i - \bar{x})^2}{n-1}$$

Bessel's
correction.

2). Standard Deviation:-

→ Square of variance.

pop

sum

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

$$s = \sqrt{\frac{\sum (\bar{x}_i - \bar{x})^2}{n-1}}$$

→ It gives you a measure of spread that is in the same units as the original data, making it easier to interpret them variance.

* What is Percentile and Per centile?

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

$$\text{Per centile} = \frac{\text{no. of values below } x}{n}$$

$$\text{Q}_1 \rightarrow \text{per centile} = \left(\frac{\text{percentile}}{100} \times n \right) + 1$$

⇒ Five Number Summary :-

1. Minimum
2. $Q_1 \rightarrow 25^{\text{th}}$ percentile
3. Median
4. $Q_3 \rightarrow 75^{\text{th}}$
5. Maximum

IQR = Inter quartile Range

$$= Q_3 - Q_1$$

$$\begin{aligned} \text{Lower fence / min} &= Q_1 - 1.5 \text{IQR} \\ \text{Upper fence / max} &= Q_3 + 1.5 \text{IQR} \end{aligned}$$

* Distribution :-

→ The possible values a variable can take & how frequently they ~~can~~ occur.

$y \rightarrow$ the actual outcome of an event.

$y \rightarrow$ one of the possible outcomes.

* Probability :-

Sample Space → Total possible outcomes.

1. Mutually exclusive

2. Collectively exhaustive

⇒ Event :-

⇒ Subset of Sample Space.

⇒ Axioms :-

1). non-negative. or $P(A) \geq 0$

2). Normalisation :- $P(\Omega) = 1$

3). Additivity :- if $A \cap B = \emptyset$

$$P(A \cup B) = P(A) + P(B)$$

if $A \cap B \neq \emptyset$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

⇒ Multiplication Rule

→ To check 2 or more events can occur together.

→ Two versions :-

1). Independent → 2 outcome independent

2). Dependent → 2 outcome depend

1). Independent event -

$$P(A \cap B) = P(A) \times P(B)$$

2) Dependent event -

$$P(A \cap B) = P(A) \times P(B/A) \rightarrow \text{conditional Prob.}$$

Prob B whom A has already occurred

$$P(B/A) = \frac{\text{Prob of what we want to know}}{\text{Prob of what we already know}}$$

⇒ Probability Distribution function (crossed out) :-

→ It gives the probability that a random variable takes a specific value or lies within a certain range.

i) Discrete \rightarrow Probability Mass fn (PMF)

ii) Continuous \rightarrow Probability Density fn (PDF).

* Types of Distribution :-

~~Discrete & Continuous~~

1. Discrete :- limited no. of outcomes

2. Continuous - infinite no. of outcomes

$$\rightarrow X \sim N(\mu, \sigma^2)$$

Variable type of Dis.

Characteristics.

i) Discrete Dist.

i) Bernoulli dist.

out $\rightarrow 2^{-0}$ \rightarrow Equiprobable

ii) Binomial :- Bernoulli \rightarrow extension.

iii) Poisson dist

~~Continuous Dists~~

ii) Continuous :-

Normal (Gaussian) :-

ii) Chi-Square dist

3) Uniform dist [discrete].-

→ $U(a, b)$ → Range of values in a set.

→ Variable X follows uniform dist with range a to b .

$$X \sim U(a, b)$$

→ All outcomes → equi probable
 $p(a) = p(b)$

2) Bernoulli dist

→ Bern(p) → Prob of success

$$X \sim \text{Bern}(p)$$

→ Event:-

1). 1 + 0's

2). 2 possible outcomes ($1/p$)

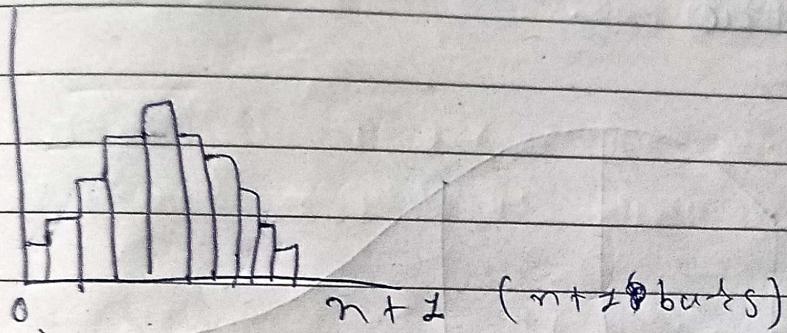
$$\rightarrow E(X) = 1 \times p + (0 \times (1-p))p$$

3) Binomial Dist.

→ ~~sample~~ sequence of Bernoulli events
Identical

$$X \sim B(n, p)$$

↳ Prob of success in each trial



$$\rightarrow P(Y) = {}^n C_Y \cdot p^Y \cdot (1-p)^{n-Y}$$

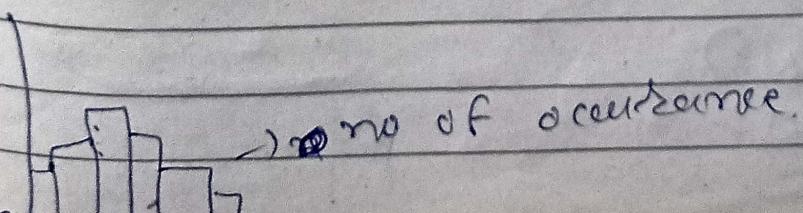
$$\rightarrow E(X) = n \cdot p$$

4) Poisson Dist.

$p_0(\lambda)$, probability

$$P(Y) = \frac{\lambda^y e^{-\lambda}}{y!} \quad e = \text{euler's}$$

$$E(Y) = \lambda$$



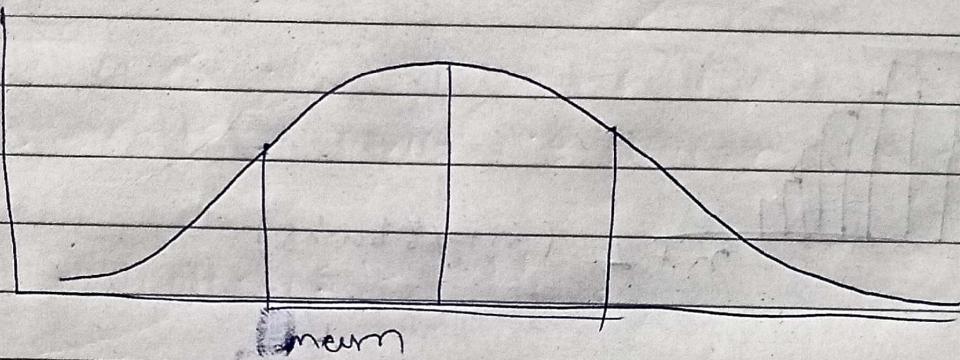
5) Normal dist (Gaussian) :-

$$X \sim N(\mu, \sigma^2)$$

~~Not bell shaped~~

ND \rightarrow Bell shape majority defer \rightarrow mean.

Graph \rightarrow symmetric



$$E(X) = \mu$$

$$\text{mean} = \text{median} = \text{mode}$$

$$68\% \rightarrow \mu \pm \sigma$$

$$95\% \rightarrow \mu \pm 2\sigma$$

$$99.74\% \rightarrow \mu \pm 3\sigma$$

outliers are extremely low/high

\rightarrow Standardization :-

$$z = \frac{\text{score} - \mu}{\sigma}$$

Standard ND

$$\mu = 0 \quad \sigma = 1$$

$N \rightarrow S N D$

($n > 30$)

+ \rightarrow dist

$y \sim t(3)$ \rightarrow Degree of freedom

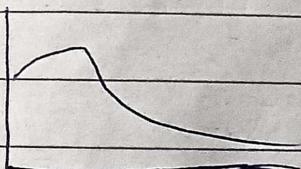
g) Chi - square dist:

\rightarrow asymmetric

\rightarrow skewed

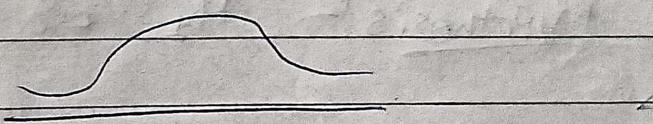
$$y \sim t(k), x \sim \chi^2(k)$$

$$y^2 \sim \chi^2(k), \sqrt{x} \sim t(k)$$

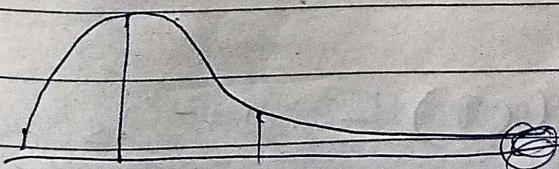


$E(X) = k$

\rightarrow Skewness

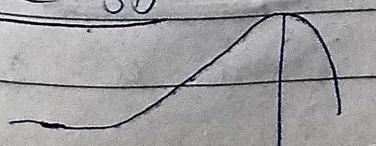


i. +ve SD:



mode < median < mean.

ii. -ve SD:



mean < median < mode

* Inferential Statistics:-

⇒ Hypothesis testing :-

- A hypothesis test is a method of statistical inference used to decide whether the data at hand is sufficient to support a particular hypothesis.
- It allows us to make probabilistic statements about population parameters.

1) NULL Hypothesis — H_0

- It assumes that ~~there~~ ^{two} there is no relationship b/w two variables [simply . sayd nothing new is happening]

2) Alternate Hypothesis — H_1 or H_a

- It states that there is significant relationship b/w two variables.
- Confidence Interval (CI) — It is a range of values within which we expect a particular population parameter to fall.
- It's a way to capture ~~uncertainty~~ uncertainty around an estimate obtained from a sample of data.

→ Confidence Level → It refers to the degree of certainty that the statistical estimate will contain the true population parameter.

90% → Narrower CI → ~~more~~ precise

95%.

99% → wider CI → less precise.

⇒ Loss of Significance) - (α)

→ It is threshold for deciding when to reject the H_0 . It defines the probability of making a Type I error, which is the error of rejecting H_0 when it is actually true.

* Rejection Region Approach:-

1. H_0, H_a State
2. $\alpha \rightarrow$ significance level $\alpha = 0.05$.
3. Assumption / value noted down.
4. Decide test
5. Value
6. Test conduct.
7. Reject / accept decision
8. Results.

→ Types of Test

- Type - I) - Rejecting H_0 when it is actually correct.
- Type - II) - Accepting H_0 when H_0 is actually incorrect.

* Z-test - $Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$

1. Random Sampling (Normalised)

2. σ

3. $n \geq 30$

* T-test - $df = n - 1$

→ One-sample t -test.

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$$

$$df = \left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} \right)^{-1} \left(\frac{s_1^2}{n_1} \right)^{-1} + \left(\frac{s_2^2}{n_2} \right)^{-1}$$

→ Independent - 2 sample t -test.

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

→ Paired t -test.

$$t = \frac{\bar{d}}{s_d/\sqrt{n}}, \quad \bar{d} \rightarrow \text{mean difference}$$

$$s_d \rightarrow \text{SD}$$

* Chi-Square Test :-

$$\chi^2 = \sum \frac{(O-E)^2}{E}, \quad O = \text{observe freq} \\ E = \text{Expected val}$$

E = Row total \times col total
Grand total.

$$df = (r-1) \times (c-1)$$