## Statistics

Descriptive stats vs Inferential statistics.

1. Measures of central tendency
2. Measures of Variability
3. Skewness - Symmetry
4. Kurtosis -

2) Variance, Standard deviation.

1. Hypothesis testing
2. Confidence intervals
3. Analysis of Variance (ANOVA)
4. Z test
5. T test.

# Sampling techniques

- Simple Random Sampling.
- Stratified Sampling.
- cluster Sampling
- Systematic sampling.
- Convenience Sampling

## Types of Variables :-

1) Quantities - Measure of Numeric

  • Continous & ~~a~~ discrete

2) Qualitative - Measure of categories

## Descriptive Statistics :-

1, Mean, 2, Mode, 3, Median.

$\hookrightarrow \{2, 3, 2\} = $ Mean $= \frac{2+3+2}{3} = 2.8$.

Median

$\Rightarrow \{1, 2, 3, 4\}$

if

$= \frac{2+3}{2} = \frac{5}{2} = 2.5 \longrightarrow$ Median.

if #2 $\Rightarrow \{1, 2, (3), 4, 5\}$

$= \{3\} \rightarrow$ Median

## Mode → Categorical data

Mode - Most repeated element

$\{a, a, b, b, b, b, c\}$

$= \underline{b}$

→ whenever we have out lies
in data then we will go with
Median.

\* whenever i have Categorical
data we will use Mode.

- \* $\underline{Variance}$
  n

$V = \sum\limits_{i=1}^{n} (x_i - \bar{x})^2 \quad x = \{1, 2, 3, 4, 5\}$

$\downarrow$
Mean of Sample

$\bar{x} = \left(\dfrac{1+2+3+4+5}{5}\right)$

$V = \sum\limits_{i=1}^{n} \dfrac{(x_i - \bar{x})^2}{n-1}$

$\bar{x} = 3$

$= \dfrac{(1-3)^2 + (2-3)^2 + (3-0)^2 + (4-3)^2 + (5-3)^2}{5-1}$

$= \dfrac{4 + 1 + 0 + 1 + 4}{5-1} \qquad = \dfrac{10}{4} = 2.5$

$V = \underline{2.5}$

Standard deviation :

$$= \sqrt{V}$$

Skewness :-

Kurtosis :- Measures the peaked nus or flatnus of a distribution. It provides info about the presence of outliers or extreme values in dataset.

* Normal distribution :-

- when the data is symmetric to both sides then that distribution is called Normal distribution.

Standard Normal distribution :

if m=0 ; SD=1 then we can say that our data is Standard Normal distribution

## Day - 16_ python!

1) why numpy ?

2) why pandas

3)

\* Measure of central tendencies

  ↳ Mean, Median, Mode.

+

## Percentile

- At what particular percentile my value is present at ?

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

Ex 1- $n = \{1, 2, 5, \overset{\downarrow x}{6}, 10\}$

$$= \frac{4}{0} \qquad\qquad = \frac{3}{5} \times 100 \qquad = \frac{6}{10} \times 100 = 0.6$$

$$\text{percentile} = 0.6\%.$$

- what if he gives percentile and we haveto know the element.

  Ex:- 75% ?

$$\{2, 5, 5, 7, 8, 9, 10\}$$

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

$$= \frac{\overset{3}{75}}{100} \times (\overset{2}{8}) = \boxed{6}$$

$$A_1$$

<u>6 is index of given element</u>

$$\frac{75}{100} \times \overset{3}{\cancel{6}} = \frac{9}{2}$$
$$\frac{9}{2} = 4.5$$
$$= 5$$

**\* 5 - Number Summary :**

$\{1 \; \underset{25\%}{5} \; \underset{\text{Med}}{7} \; 8 \; 9\}$
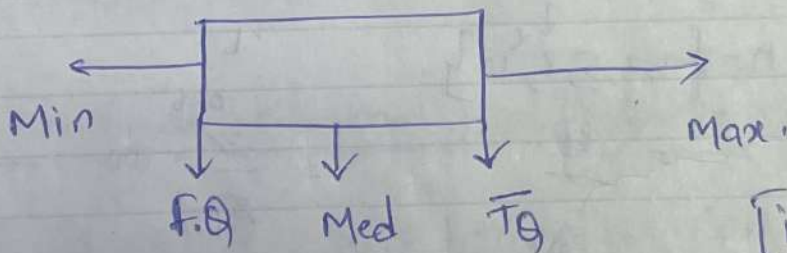
1) - Mine - (1) minimum Value

2) first quartile → <u>25%</u> - ② => <u>5</u>

3) Median - 7

4) third quartile - 8.

5) Max - 9.

we will do our Box plot with this.



Min          Max.

f.Q   Med   TQ        [iqr = Inter Quartile Range]

Lower fence => $Q_1 - 1.5 * iqr$

Upper fence => $Q_3 + 1.5 * iqr$ [∵ $iqr = Q_3 - Q_1$]



Lf   $Q_1$  Me  $Q_3$   UF.
         dian

Standard Normal Distribution.

$$Z = \frac{x_i - \mu \,(Mean)}{\sqrt{\sigma}\,(S\text{-}D)}$$

N·D $\longrightarrow$ S·N·D

Day - 17 - python

Questions :- How to choose Datastructures
  $\rightarrow$ Lambda function
  $\rightarrow$ oops, class, object, Inheritance, polymorphism
Encapsulation
  $\rightarrow$ why numpy!
  $\longrightarrow$ why panda!

\* probability.

Emperical Rule

68      95      99.7
                95%.
68% of data    lies in    99.7data
lies in 1st S·D  data Second  lies in 3rd
     ±ve "      S·D        S·D.

Z Score

# Probability — Measure of likelihood of an event

Ex:- Roll of Dice $\{1, 2, 3, 4, 5, 6\}$

$$P(4) = \frac{1}{6}.$$
$$P(1) = \frac{1}{6}.$$

Toss of a Coin $\{H, T\}$

$$P(H) = \frac{1}{2} \; ; \; P(T) = \frac{1}{2}$$

## Additive Rule $P(A \text{ or } B)$

$$\{H, T\} \xRightarrow{\text{As we know}} P(H) \to \frac{1}{2}$$

$$P(T) = \frac{1}{2}$$

what is $P(H \text{ or } T) \Rightarrow P(H) + P(T)$

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

$$P(H \text{ or } T) = 1$$

## Mutually Exclusive Events :

### In a Dice.

→ Is it possible to get 1 at the same time when we get 6?

No

→ Answer to this question is wrong ✗

for Example we can take Coin
- So what if we want to get H & T at same
time ?

No

- Answer would be wrong ✗ again

## Non-Mutually Exclusive Events

for Example

To Deck of cards [52]

$$P(Q) \, \& \, P(♡) \quad i.e; \quad P(Q \, \& \, ♡) \, ?$$

is it possible to get both at same time ?
→ Answer to this question is yes ✓ we can
this type of events are called Non Mutual Exclus
events.

for mutual Exclusive.

what is $P(Q) = \dfrac{4}{52}$

and $P(\heartsuit) = \dfrac{13}{52}$

- $P(Q \text{ or } \heartsuit) = \dfrac{4}{52} + \dfrac{13}{52}$

$$= \dfrac{17}{52}$$

and $P(Q \text{ and } \heartsuit) = \dfrac{1}{52}$

## Non-Mutual Exclusive.

$$P(Q \text{ or } \heartsuit) = P(Q) + P(M) - P(Q \text{ and } \heartsuit)$$

$$= \dfrac{4}{52} + \dfrac{13}{52} - \dfrac{1}{52} = \dfrac{16}{52}$$

# Conditional probability :-

## Dependent event :-

## Permutation f combination

$$nPr = \frac{n!}{(n-r)!}$$

$= \{a, b, c, d, e, f\}$

$n = 6; \quad r = 3$

(creating a pair of 3 elements)

$\bullet \quad nPr = \dfrac{6 \times 5 \times 4 \times 3 \times 2 \times 1}{3 \times 2 \times 1}$

$= 120$

So we can have 120 permutations from

Ex :- → abc
       cba
       deb
       a c b
         |
         |
         |
         |

Combination :- The main difference b/w permutation & Combination is there is no chance of repetion like what we see in permutations.

$$= \{a, b, c, d, e, f\}$$

$$\therefore n_{c_r} = \frac{n!}{r!\,(n-r)!}$$

$$n = 6\,; \quad r = 3 \quad = \quad \frac{6 \times 5 \times 4 \times 3 \times 2 \times 1}{3 \times 2 \times 1\,(3 \times 2 \times 1)}$$
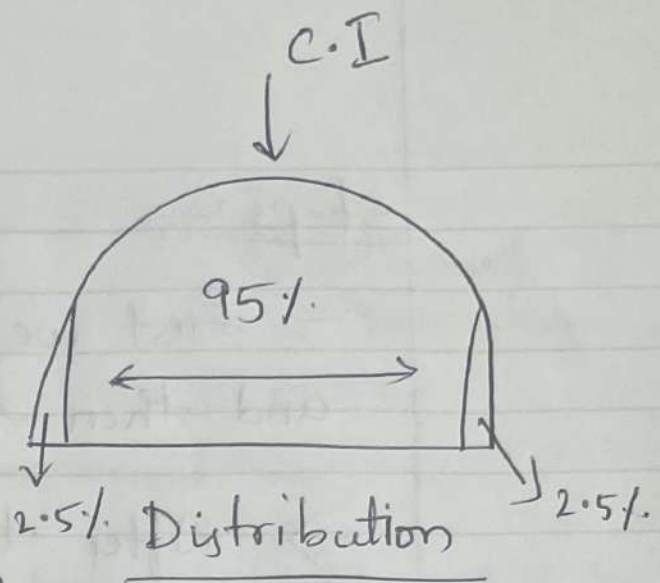
$$= 5 \times 4$$

$$= 20.$$

So we can get 20 Combinations

*** Inferential Statistics :-

- Drawing Conclusions or making predictions about larger dataset based on Sample data.

## Hypothesis testing :-

→ Confidence Interval . (C·I)

→ Significance value

→ Null hypothesis ⇒ $H_0$

→ Alternate Hypothesis ⇒ $H_A$

→ P-Value.

$$C.I = 95\%$$

$$= 1 - \frac{CI}{100}$$

Significance Value ⇒ 1-0.95

$$S.V (\alpha) \text{ value} = 0.05$$



2.5% Distribution

Next checking the P value

if probability of touching at 2.5% is 0.01

and the $\alpha$ value = 0.05

we will check $(P < \alpha)$ in our case.

$$0.01 < 0.05$$

it satisfies the condition.

then we will reject our Null hypothesis

## Steps :-

→ First we create our Null Hypothesis and then Alternate Hypothesis.

→ After that we will define our Significance Value. and p value (after Conducting Some tests)
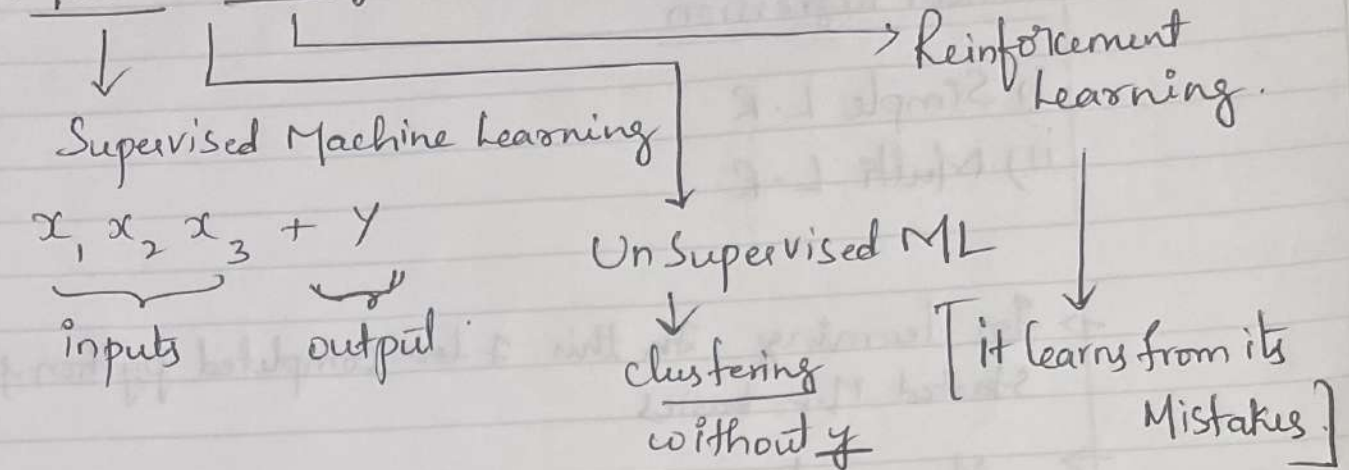
→ And Based on that p value we will decide which one to choose.

## Different kind of tests!

1) One Sample T-test
2) One Sample Z-test
3) One Sample proportion test.
4) Two Sample T-test
5) Two Sample Z-test
6) Two Sample proportion test.
7) paired-T test
8) Anova-test
9) Chi-Square test.
10) 1 tail & 2 tail test.

ML
Deeplearning
& NLP

Day-18 :-

## Machine Learning

Supervised Machine Learning → Reinforcement Learning.

$x_1, x_2, x_3 + y$

inputs    output

Un Supervised ML
↓
clustering
without $y$

[it learns from its Mistakes]

## Supervised M·L :-

→ If my output data is Continous then it is Regression problem.
   Variable

→ If my o/p Variable is discrete then it is ~~Continous~~ Classification problem.

— we have two different Models i) Parametric models
                                ii) Non parametric Mode

i) P·M's

   a) Linear Regression ·→ R
   b) Logistic Regression ·→ C

ii) Non P·M's

   a) Decesion tree (DT)
   b) Random forest (RF)
   c) Support Vector Machine (SVM)
   d) Ada Boost, Extra Variant boost