

Special Topic → MACHINE LEARNING for NLP / Text Analytics

- ① Write
- ② Understand "MATHS"

AGENDA

① ABC ML

② Solve Algorithm (Maths and logic)

Text Analytics

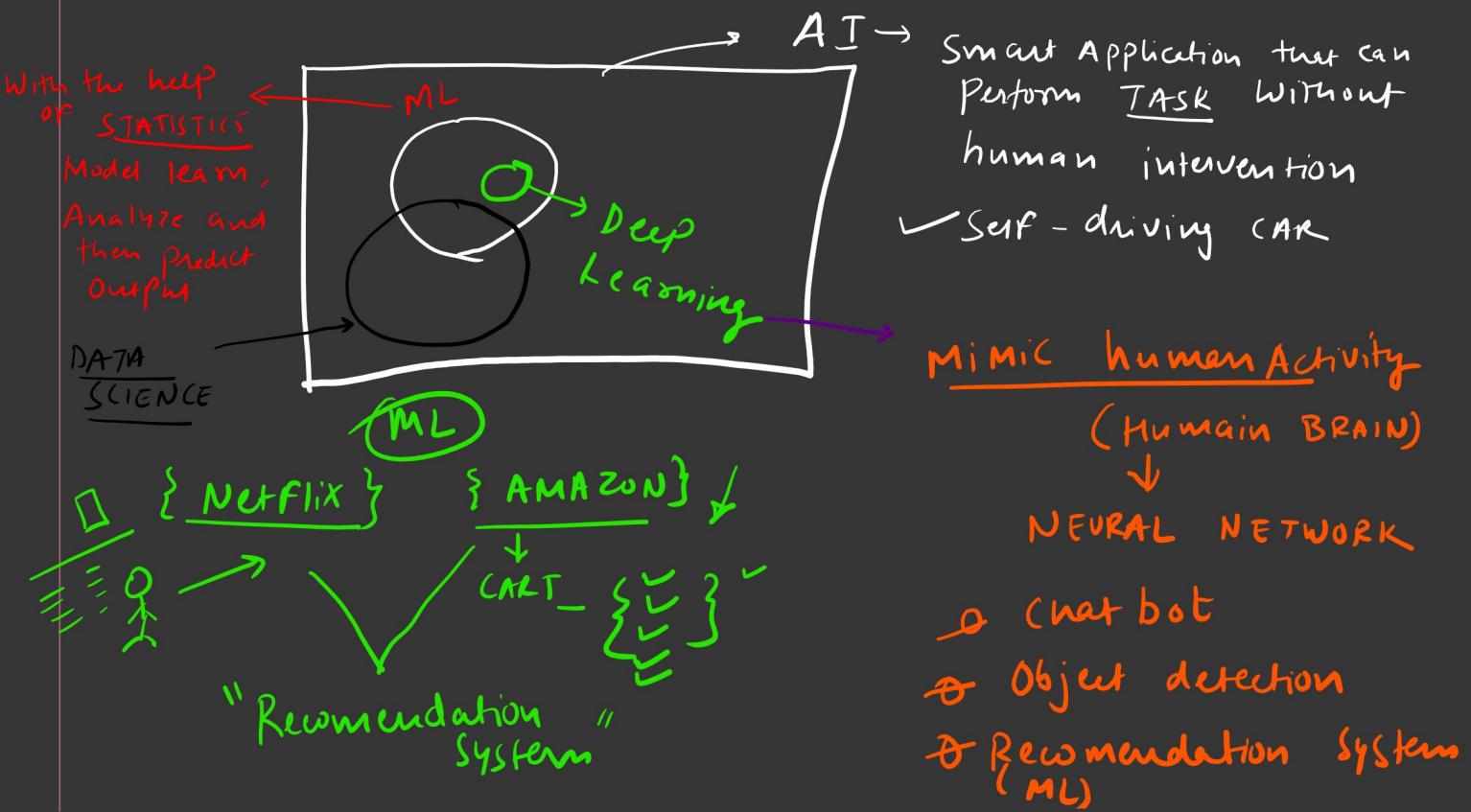
① Linear Regression

② Logistics Regression

③ Decision Tree

④ NAIVE BAYES

⑤ KNN



PRE - Requisite

1. Linear Algebra
2. Statistics
3. Probability
4. CALCULUS

STATISTICS - 101

Statistics is the science of Collecting, Organizing and

Analyzing the DATA

" Piece of information "

STATISTICS

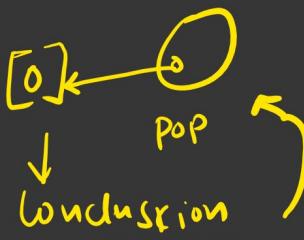
Descriptive

ORGANIZING +
SUMMARIZING
DATA

- ① M - central Tendency
MEAN - MEDIAN - MODE

Influential Stats

Measured DATA + conclusion

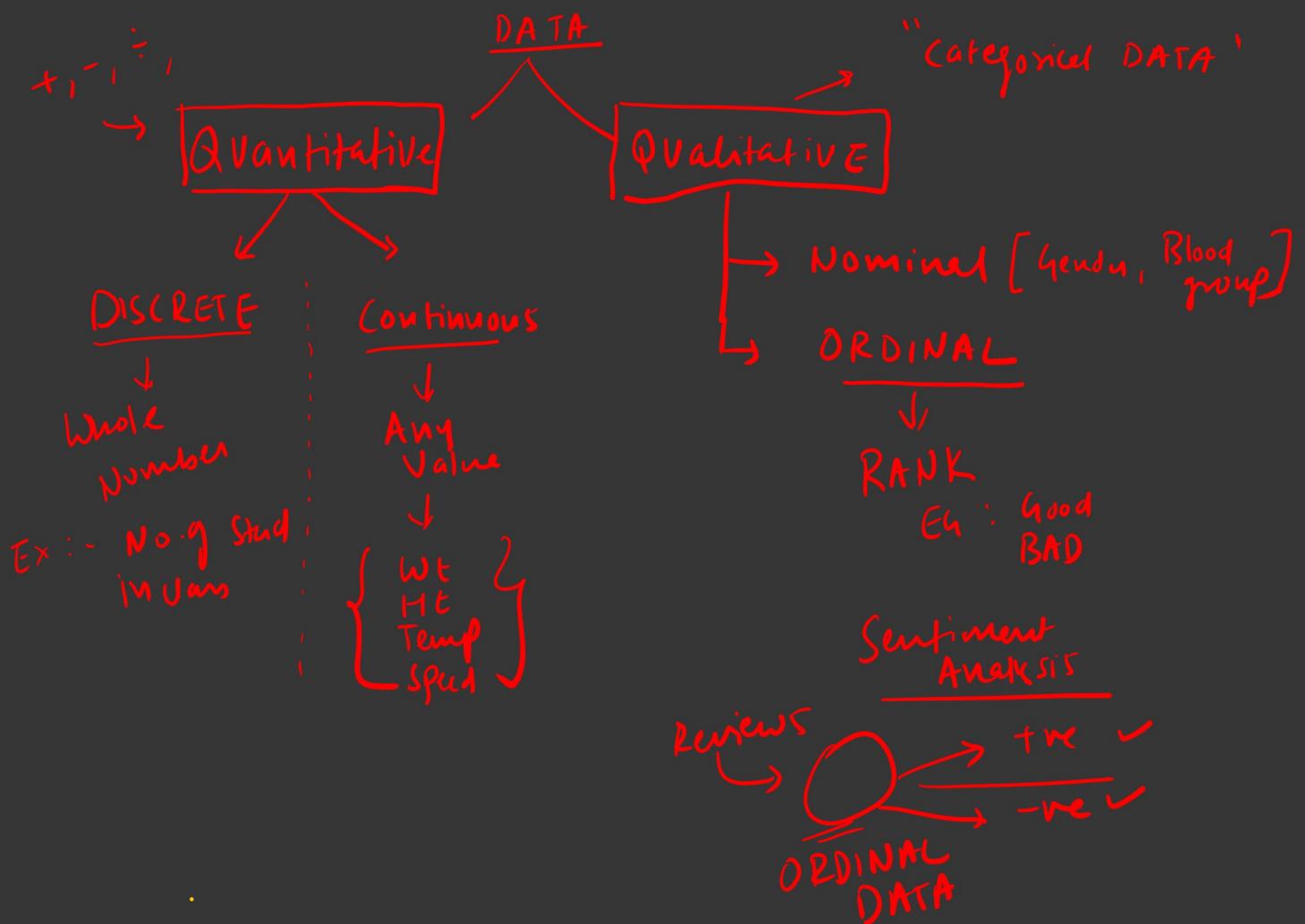


- ⊖ Z-test
- ⊖ T-Test
- ⊖ Hypothesis Testing

- ② MEASURE OF DISPERSION
 - VARIANCE
 - Std dev

Summarizing DATA

- ✓ HISTOGRAMS
- ✓ BAR CHART
- ✓ PIE CHART



PPR Reg SCALE OF MEASUREMENT OF DATA

1. NORMAL Scale DATA [ORDER does not Matter]
2. ORDINAL Scale DATA [RANK + ORDER matters]
3. Interval Scale DATA [$0 \rightarrow$ does not have zero starting value]
4. Ratio Scale DATA [ORDER / RANK matters & Ratio]

B-Tech → M-Tech → Phd → Postdoc

(4th) (3rd) (2nd) (1st)

Pre-Reg from STATISTICS PERSPECTIVE

① MEAN, MEDIAN, MODE

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

MEDIAN $x_i = \{4, 5, 2, 3, 1\}$ * SORTING IS MUST

$$\text{SORT} = \{1, 2, 3, 4, 5\}$$

Median value

MODE = {Frequency Maximum}

$$x = \{1, 1, 1, 2, 3, 4\} \quad \boxed{\text{Mode} = 1}$$

VARIANCE

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

x_i = Data point

μ = Population Mean

N = Size of Population

STANDARD DEVIATION

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

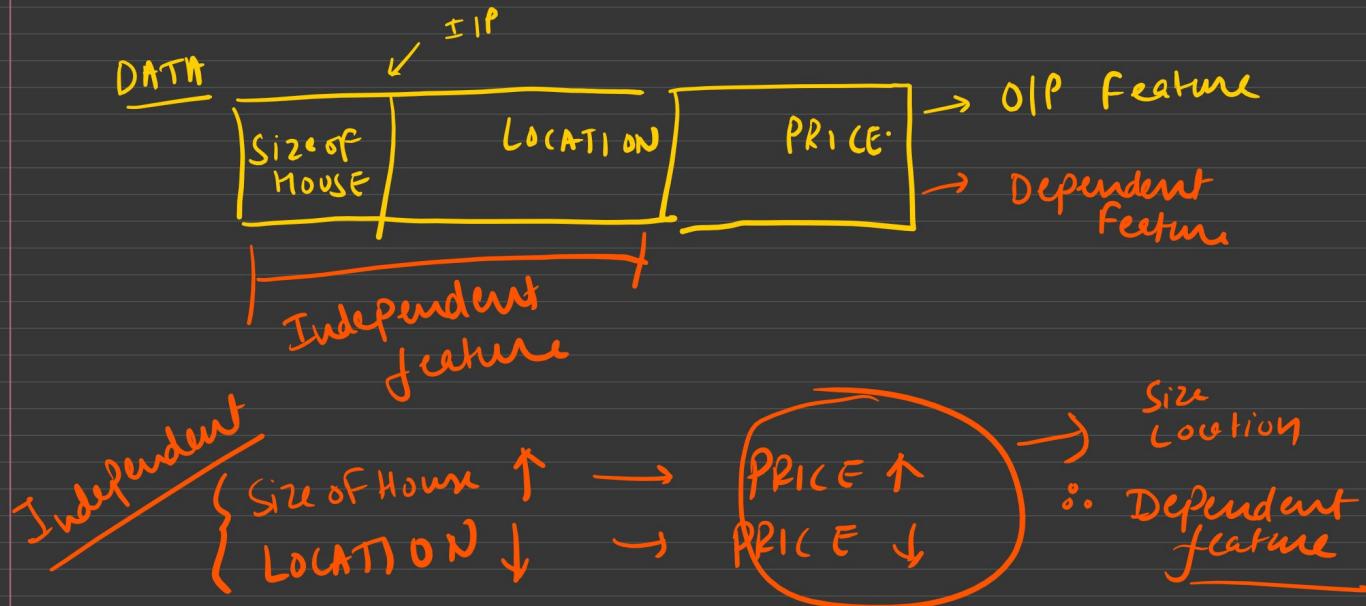
* Probability \rightarrow Random Variable

COVARIANCE - CORRELATION

X	Y
2	3
4	5
6	7
8	9

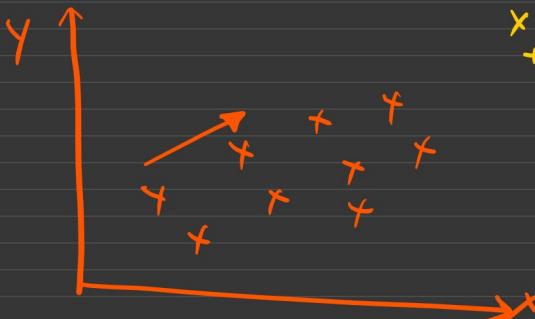
Relationship b/w X and Y





Scenario #1

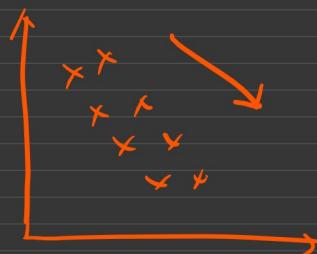
$$\begin{array}{ll} X \uparrow & Y \downarrow \\ X \downarrow & Y \uparrow \end{array}$$



X, Y has
positive Covariance

Scenario #2

$$\begin{array}{ll} X \downarrow & Y \uparrow \\ X \uparrow & Y \downarrow \end{array}$$



X, Y has Negative Covariance

COVARIANCE (X, Y)

$$\Rightarrow COV(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1} \quad \text{+ Relationship between DATA POINTS}$$

$$VARIANCE(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \quad \text{+ } COV(X, Y)$$

$$VAR(X, X) = \frac{\sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})}{n-1} \quad \text{+ "Spread of DATA")}$$

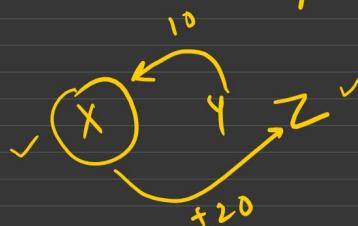
X	Y
2	3
4	5
6	7

$\bar{x} = 4$ $\bar{y} = 5$

$$\begin{aligned}\text{Cov}(x, y) &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1} \\ &= \frac{(2-4)(3-5) + (4-4)(5-5)}{3-1} \\ &= \frac{4+0+4}{2} = 4\end{aligned}$$

$$\boxed{\text{Cov}(x, y) = 4}$$

+ve Related

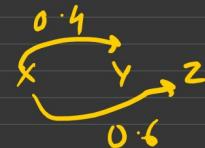


{ covariance does not have any specific limit }

$[-1 \text{ to } 1]$

* Pearson Correlation Coefficient

\downarrow
 $[-1 \text{ to } 1]$



$$r_{x,y} = \frac{\text{Cov}(x,y)}{\sigma_x \cdot \sigma_y} = [-1 \text{ to } 1]$$

*

size
of house

Rooms

location

\wedge \wedge \times

-ve correlated +ve (0-related)



Pearson Coefficient

Price

Text ANALYTICS or Pearson Correlation Coefficient

for your knowledge

$$\boxed{\textcircled{3} \text{ SPEARMAN RANK CORRELATION}}$$

$$r_s = \frac{\text{Cov}(R(x), R(y))}{\sigma_{R(x)}, \sigma_{R(y)}} \quad \begin{array}{l} R(x) = \text{Rank of } x \\ R(y) = \text{Rank of } y \end{array}$$

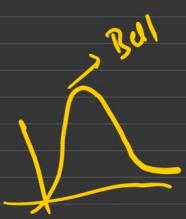
⇒ HISTOGRAMS

↳ Stemmers
↳ Box Plot



⇒ PROBABILITY DISTRIBUTION

1. Normal / Gaussian



2. Bernoulli distribution

PMF = Prob mass function

3 Uniform distribution

PDF = Prob density function

4 log Normal distribution

5 Poisson distribution

6 Power law distribution

7 Binomial distribution



① Percentiles and Quantiles

Percentile : 1, 2, 3, 4, 5, 6

$$\% \text{ of No. that are odd} = \frac{3}{6} = \frac{1}{2} = 50\%$$

PERCENTILE : Value below which a certain $\% \text{ of}$ observation OR DATA POINT LIES

Percentile

Quartiles

$Q_1 \rightarrow 25\text{ percentile}$

$Q_2 \rightarrow 50\text{ percentile}$

$Q_3 \rightarrow 75\text{ percentile}$

→ Summarize the DATASET (Boxplot)

1. Minimum = 1

2. $Q_1 = 2$

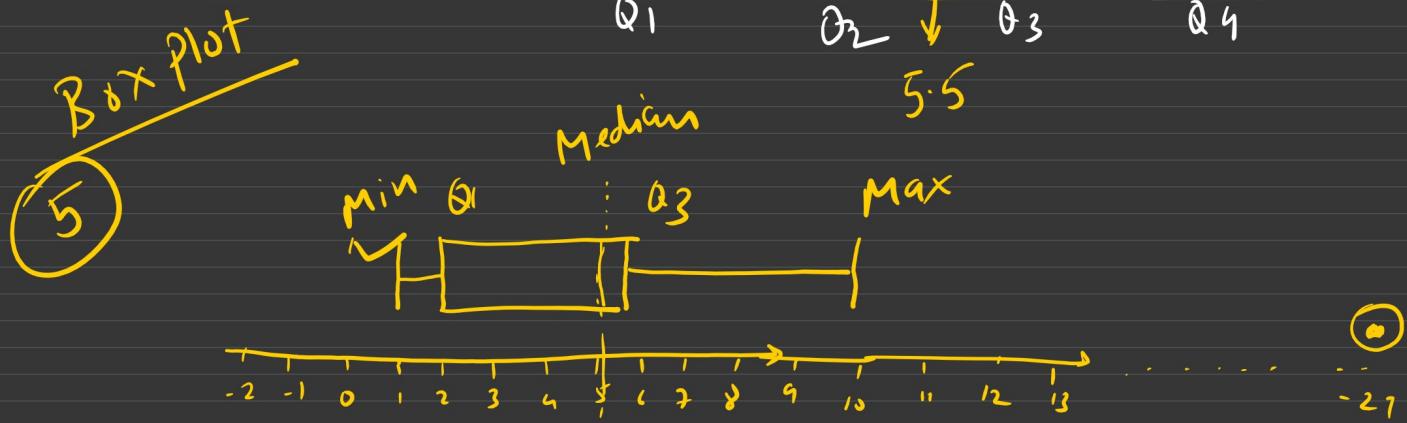
3. MEDIAN = 5.5 $X = \{1, 2, 2, 2, 3, 3, 4, 5, 5, 6, 6, 6, 6, 7, 8, 9, 10, 29\}$

4. $Q_3 = 6$

5. MAX = 10

1. Outlier

$$X' = \left\{ \left[\frac{1, 2, 2, 2}{Q_1} \right], \left[3, 3, 4, \underset{5}{\boxed{5}}, \underset{5}{\boxed{6}}, \underset{6, 6}{\boxed{6, 6}} \right], \left[\underset{Q_3}{7, 8, 9, 10} \right] \right\}$$



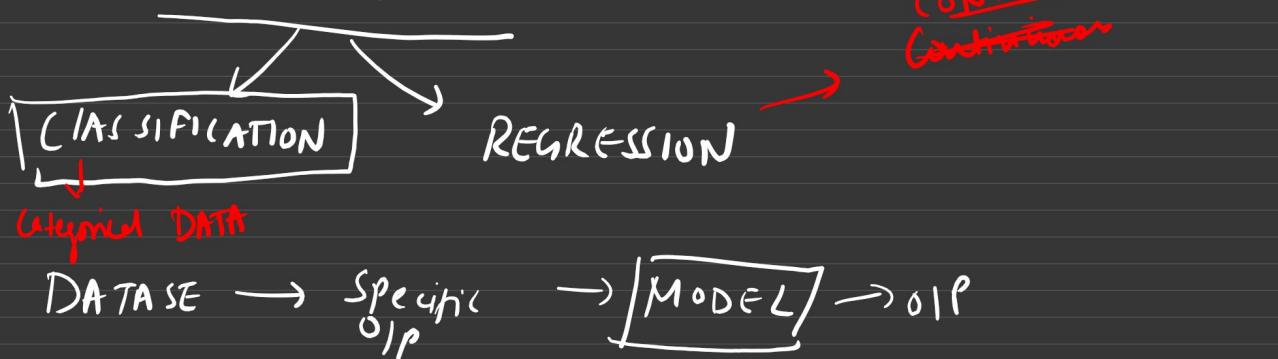
MACHINE LEARNING

Types of ML

1. SUPERVISED ML
2. UNSUPERVISED ML
3. SEMI - SUPERVISED ML

4 Reinforcement ML

① SUPERVISED ML



* CLASSIFICATION ML

PLAY HOURS	STUDY HOURS	PASS / FAIL
8	2	PASS
7	3	FAIL
3	6	FAIL

⇒ Categorical Feature
 PASS → FAIL

Regression

Size of House	No. of Rooms	Price
1M	1.2M	\$1.6M
1.2M	1.4M	\$1.8M
1.4M	1.6M	\$2.0M

Continuous

② UNSUPERVISED ML

DATA → Cluster or Similar group

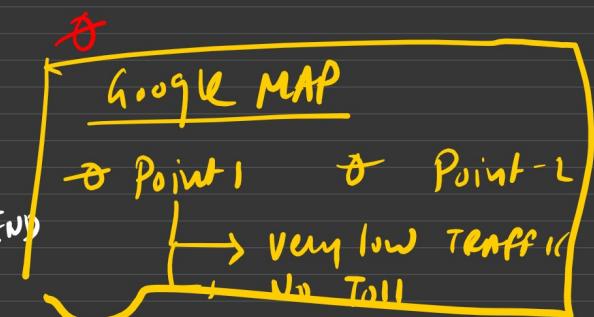
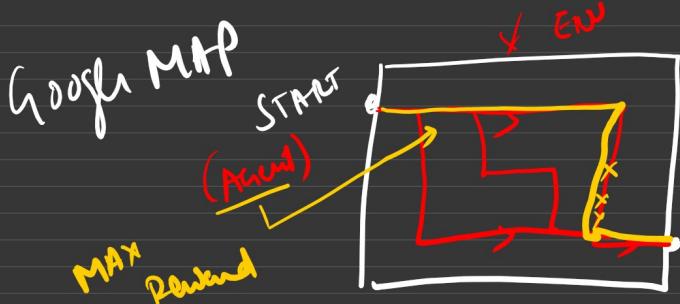
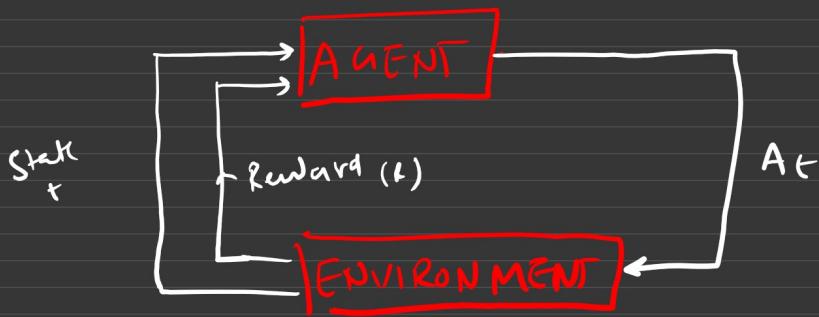


③ Semi Supervised

Supervised + Unsupervised

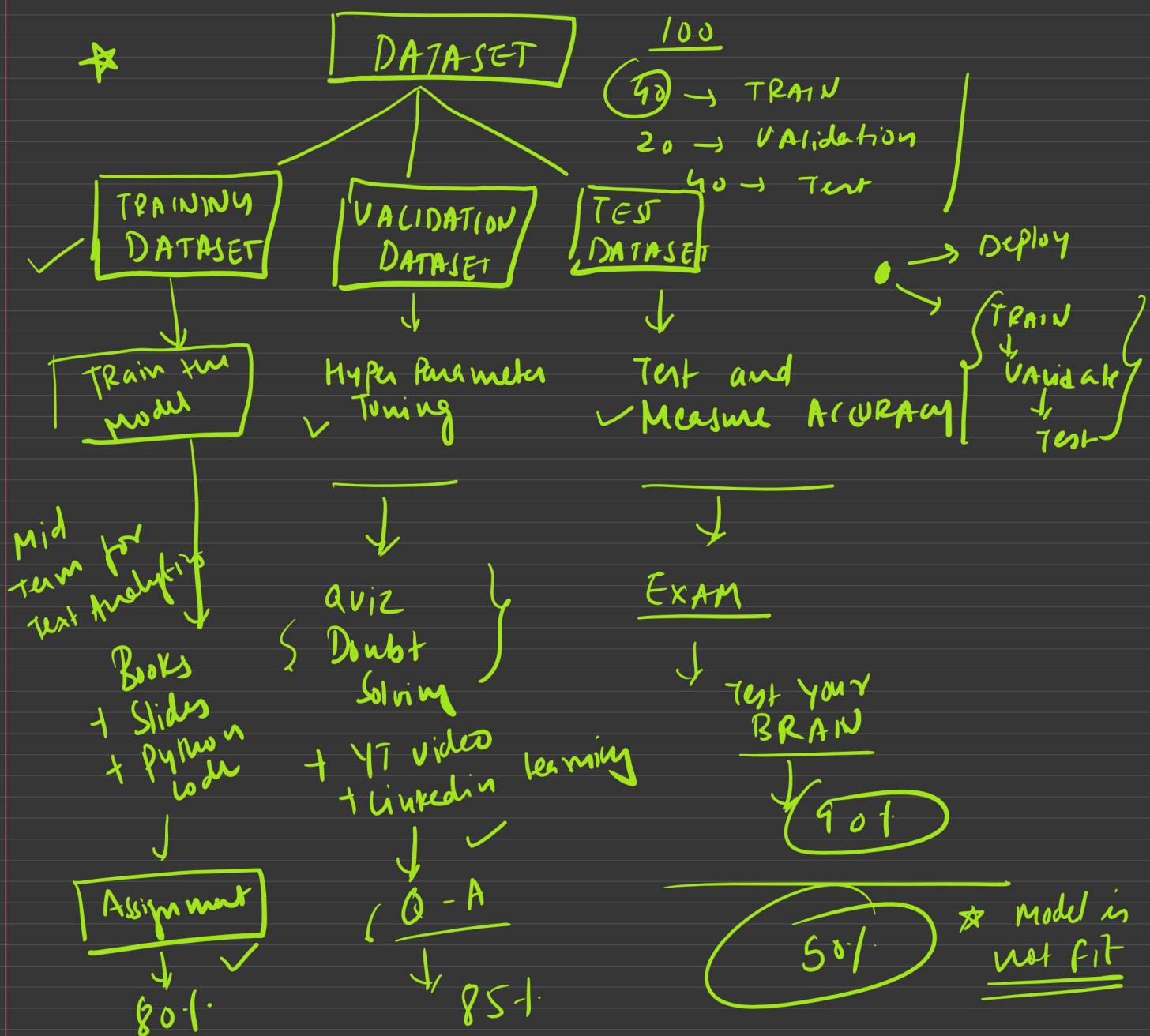
↓ ↓
* Classification : / "Clustering"
Regression

④ Reinforcement learning



Reinforcement learning

Area of ML concerned with how intelligent Agents ought to take Actions in an Environment in Order to Maximize the notion of Cumulative Rewards.



⇒ DATASET

⇒ Model Performance

BOOK → TRAIN → Model is TRAINED

Accuracy → 95%

EXAM → TEST → Model is Tested

Accuracy → 50%

OVER
FITTING

★ { low BIAS
HIGH VARIANCE }

TRAIN → ACCURACY 55%

TEST → ACCURACY 50%

UNDERFIT Model

{ high BIAS
High VARIANCE }

Overfit → Retrain the Model
With New Algorithm / Method

Underfit → Retrain the Model
With More DATA

GOAL :-

GENERALIZED Model

[TRAIN → Acc ↑↑ (85%)
TEST → Acc ↑↑ (85%)] ⇒ Low BIAS
Low VARIANCE

FEATURE EXTRACTION

⇒ F.E is process of selecting and extracting the **Most important** feature from raw data

ML App → 1000 feature



All Features are Not
Imp

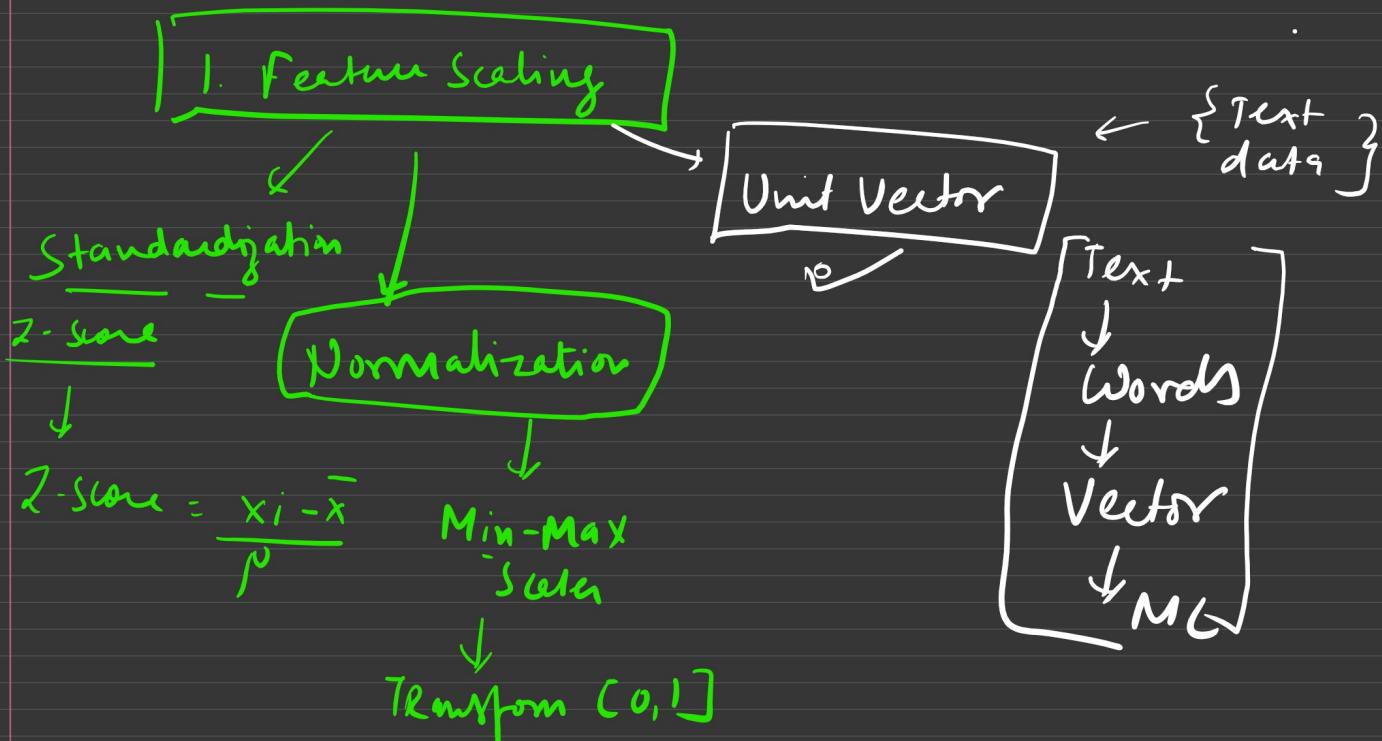


Select "Features"



TRAIN the Model

1. Feature Scaling
2. Feature Selection
3. P.C.A → PRINCIPAL Component Analysis



UNIT VECTOR

→ conversion to Unit Vector → Vector
 $\sqrt{3^2 + 4^2}$ → Magnitude of vector
 \rightarrow Magnitude of vector
 \rightarrow Vector coordinates for Unit Vector
 \rightarrow Magnitude of $UV = 1$

$$\vec{x} = (3, 4)$$



Pythagoras theorem

$$|\vec{x}| = \sqrt{3^2 + 4^2} = \sqrt{9 + 16} = \underline{\underline{5}}$$

$$\hat{p} = \left(\frac{3}{\sqrt{25}}, \frac{4}{\sqrt{25}} \right) = \left(\frac{3}{5}, \frac{4}{5} \right)$$



Calculate Magnitude of Unit Vector = 1

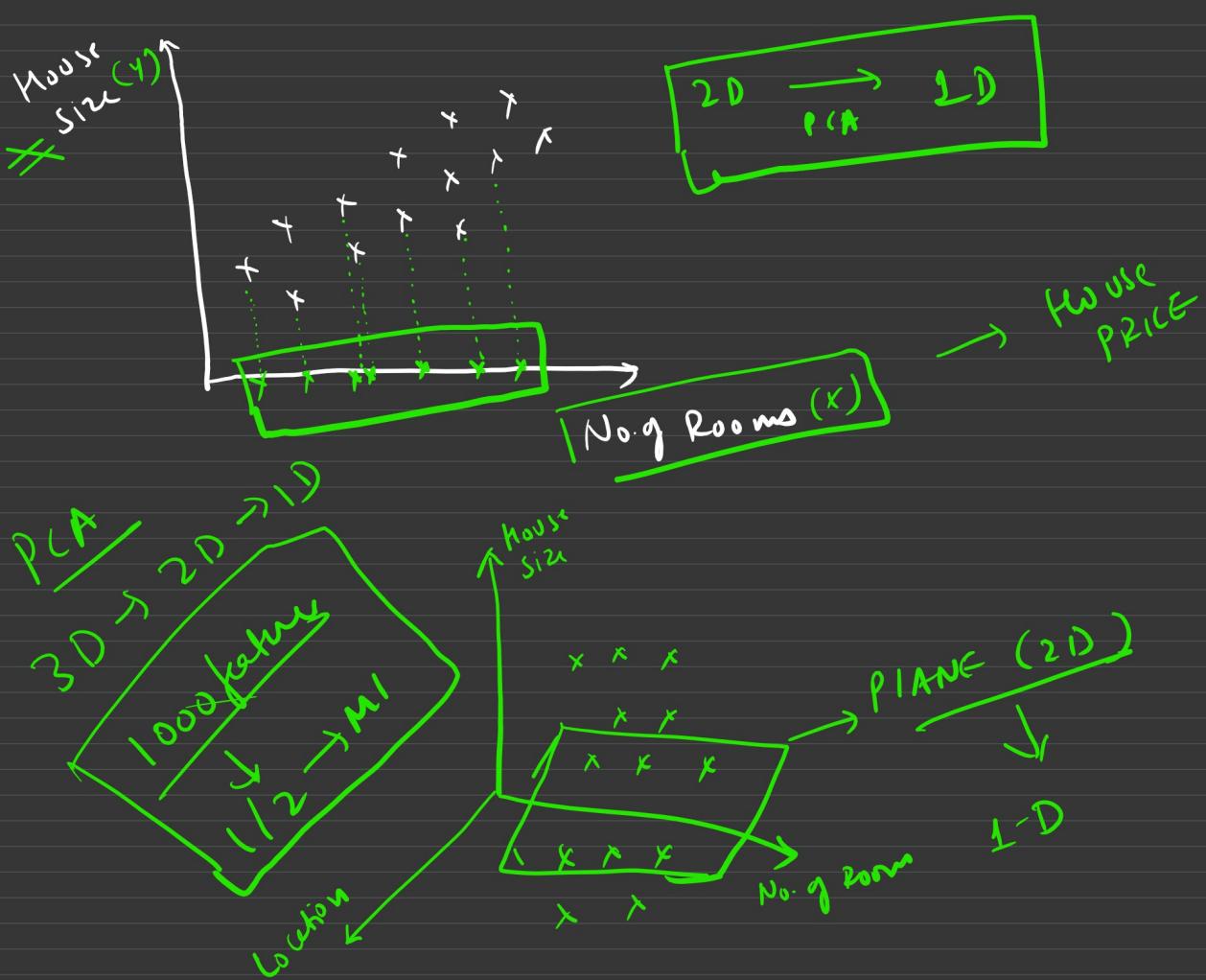
$$|\hat{p}| = \sqrt{\left(\frac{3}{5}\right)^2 + \left(\frac{4}{5}\right)^2} = \sqrt{\frac{9+16}{25}} = \sqrt{\frac{25}{25}} = 1$$

PCA → PRINCIPLE COMPONENT ANALYSIS

Dataset → 100 features → 10 features → ML Algorithm

Disadvantage

~~Loss of DATA~~



⇒ loss of data

PCA → Hypothesis testing

DATA Scientist

Q → COVID-19 vaccine discovery

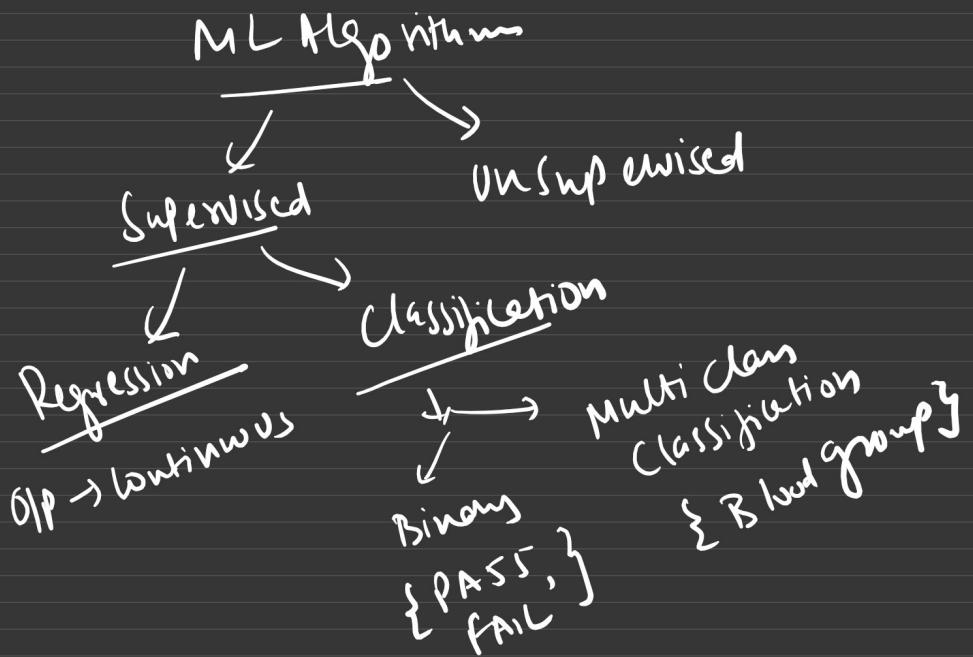
Age	Nationality	RNB	DNA	Blood group	WT	SURVIVAL RATE	PCAT
1	2	3	4	5	6		

⇒ hypothesis

Blood group (A^+) → less chance of survival in covid

True ↘ False

#1 LINEAR REGRESSION



SIMPLE LINEAR REGRESSION

DATA SET

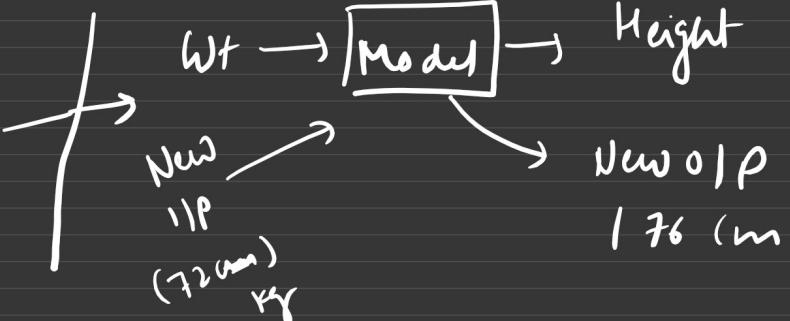
Independent Dependent "O/P"

Weight Height

Indep	Depend
Wt	Height
74	170
80	180
75	175
-	-
-	-
72	-

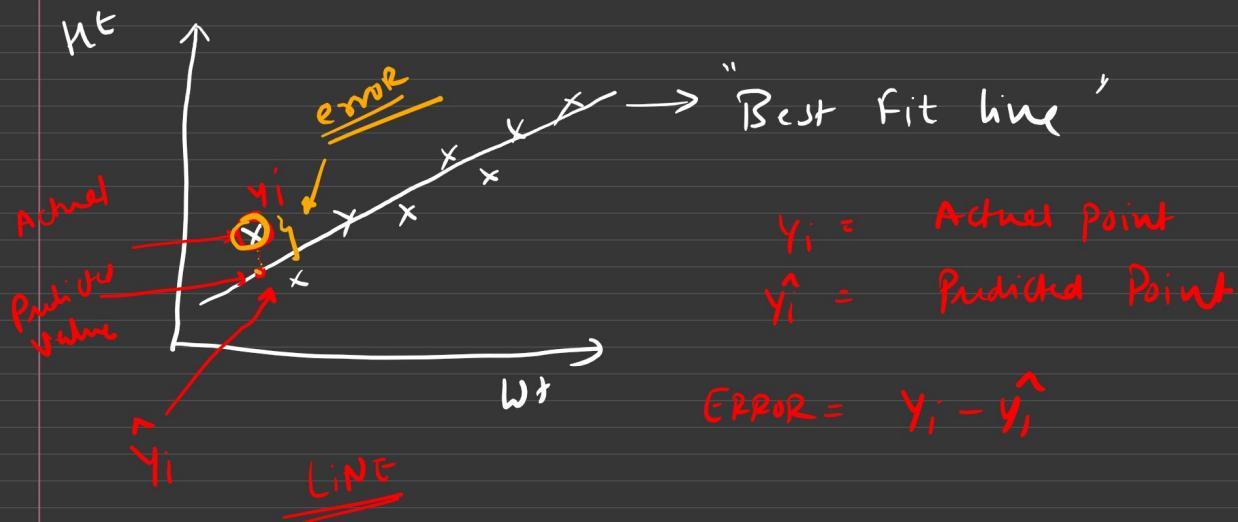
wt (given) → ht (predict)

⇒ Prediction Model

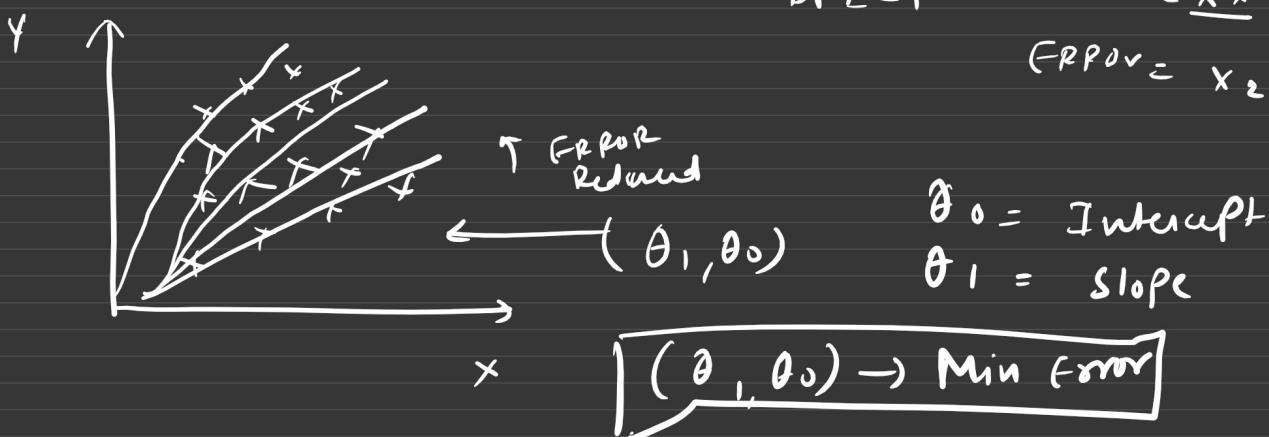


Ø Simple linear Regression

Goal → Find Best Fit line, in such a way
that the sum of error is Minimum



Goal → Minimize Error



* Manned (θ_0, θ_1) → Reduce ERROR

Optimizer -

