

ML

① Linear Regression $\hat{y}_i = \beta_0 + \beta_1 x_i$

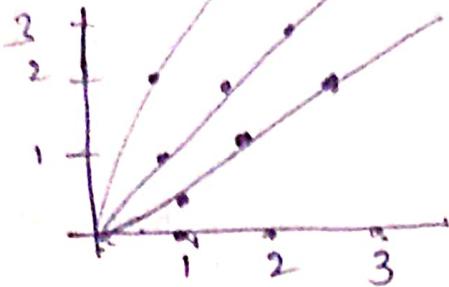
$$\hat{y}_i = \beta_0 + \beta_1 x_i$$

\downarrow

β_0

$$\hat{y}_i = \beta_1 x_i$$

$$\beta_1 = 1$$



$$\beta_1 = 1$$

$$\beta_1 = 0.5$$

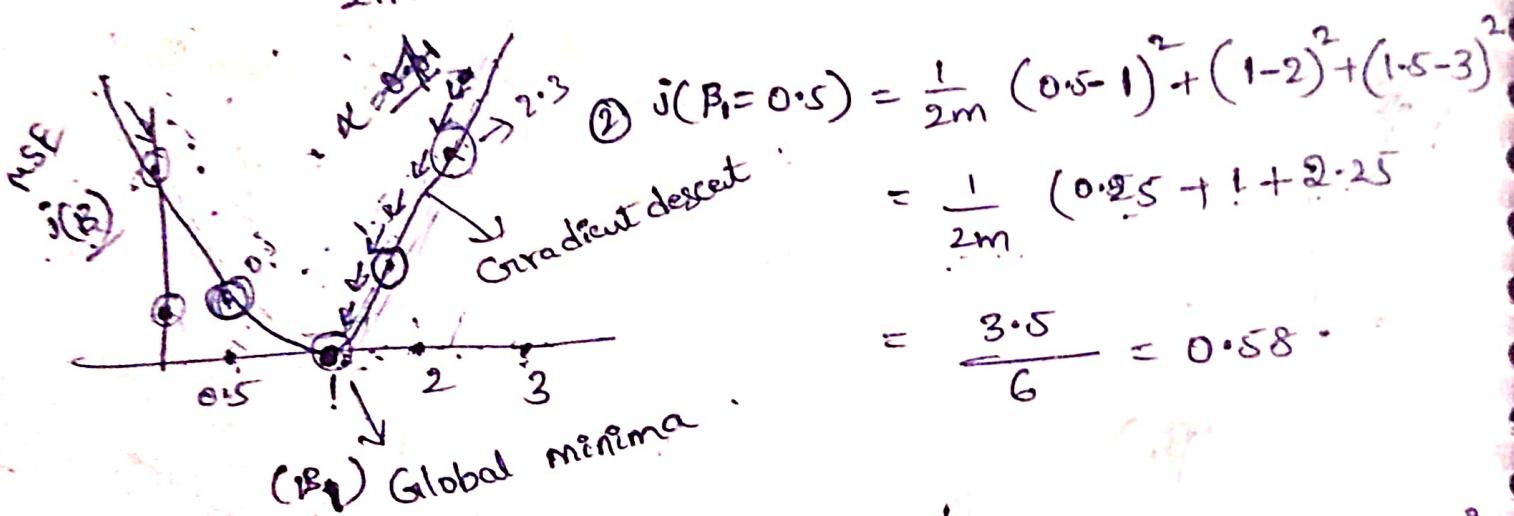
$$\beta_0 = 0$$

$$\hat{y}_i = \beta_0 + \beta_1 x_i$$

$$0.5$$

$$\textcircled{1} J(\beta_1) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2$$

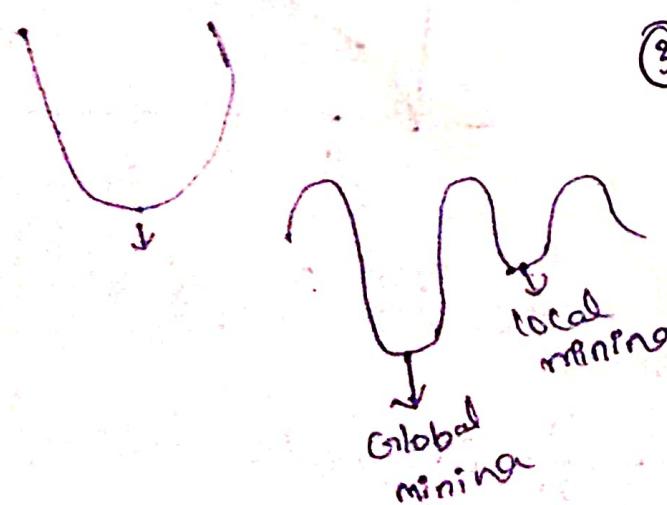
$$= \frac{1}{2m} (1-1)^2 + (2-2)^2 + (3-3)^2 = \frac{0}{2m} = 0$$



$$\textcircled{2} J(\beta_1 = 0.5) = \frac{1}{2m} (0.5-1)^2 + (1-2)^2 + (1.5-3)^2$$

$$= \frac{1}{2m} (0.25 + 1 + 4.25)$$

$$= \frac{3.5}{6} = 0.58$$



$$= \frac{1}{2m} 1 + 4 + 9$$

$$= \frac{14}{6} = \frac{7}{3} = 2.3$$

$$\theta_j = \theta_j - \alpha \cdot \cancel{\frac{\partial J}{\partial \theta_j}} \frac{dJ}{d\theta_j}$$

Gradient Descent for Linear Regression

$$\beta_0 = \beta_0 - \alpha \cdot \frac{\partial}{\partial \beta_0} \left(\frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 \right)$$

$$= \beta_0 - \alpha \cdot \frac{\partial}{\partial \beta_0} \left(\frac{1}{2m} \sum_{i=1}^m (\beta_0 + \beta_1 x_i - y_i)^2 \right)$$

$$= \beta_0 - \alpha \cdot \left(\frac{1}{2m} \sum_{i=1}^m (\beta_0 + \beta_1 x_i - y_i) \right)$$

$$= \beta_0 - \alpha \cdot \frac{1}{m} (\hat{y}_i - y_i)$$

$$\begin{aligned} & \frac{\partial (ax+b)}{\partial x} \\ &= 2(ax+b) \cdot \frac{\partial ax}{\partial x} \end{aligned}$$

$$\beta_1 = \beta_1 - \alpha \cdot \frac{\partial}{\partial \beta_1} \left(\frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 \right)$$

$$= \beta_1 - \alpha \cdot \frac{\partial}{\partial \beta_1} \left(\frac{1}{2m} \sum_{i=1}^m (\beta_0 + \beta_1 x_i - y_i)^2 \right)$$

$$= \beta_1 - \alpha \left(\frac{1}{2m} \sum_{i=1}^m (\beta_0 + \beta_1 x_i - y_i) \cdot \frac{\partial R_i(x)}{\partial \beta_1} \right)$$

$$= \beta_1 - \alpha \cdot \frac{1}{m} (\hat{y}_i - y_i) \cdot x_i$$

$$\boxed{\beta_j = \beta_j - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i) \cdot x_i^{(j)}} \quad \text{for all } j \in \{1, 2, 3, \dots, n\}$$

100
0.1

② Logistic Regression

Cost function

$$J(\beta) = -\frac{1}{m} \sum_{i=1}^m [y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)]$$

$y_i \in \{0, 1\}$

$$\hat{y}_i = h_{\beta}(x_i) = \frac{1}{1+e^{-\beta^T x_i}} \quad (\text{predicted probability})$$

Hence, m = no. of training observations

y_i = Actual class label (0 or 1)

\hat{y}_i = Predicted probability

β = model parameters.

likelihood function

$$P(y_i/\hat{y}_i) = (\hat{y}_i)^{y_i} (1-\hat{y}_i)^{1-y_i}$$

$$* y_i = 1 \rightarrow \hat{y}_i$$

$$* y_i = 0 \rightarrow 1 - \hat{y}_i$$

Log. likelihood Function

$$L(\theta) = \prod_{i=1}^m P(y_i/\hat{y}_i)$$

$$\because (\log(a^b) = b \log a)$$

$$(\log(ab) = \log a + \log b)$$

Taking log on both sides,

$$\Rightarrow \log L(\theta) = \log \left[\prod_{i=1}^n (\hat{y}_i)^{y_i} \cdot (1-\hat{y}_i)^{1-y_i} \right]$$

$$\Rightarrow \log L(\theta) = \sum_{i=1}^n \left[\log(\hat{y}_i)^{y_i} + \log(1-\hat{y}_i)^{1-y_i} \right]$$

$$= \sum_{i=1}^n \left[y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i) \right]$$

Negative log-likelihood,

$$J(\beta) = -\frac{1}{m} \sum_{i=1}^m \left[y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i) \right]$$



Gradient descent optimization



$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial J}{\partial \theta_j}$$

③ Naive Bayes → Binary classification



$$P(A \text{ and } B) = P(A) * P(B/A)$$

$$P(R) = 3/5$$

$$P(B \text{ and } A) = P(B) * P(A/B)$$

$$P(G) = 2/4$$

Equalling.

$$P(A) * P(B/A) = P(B) * P(A/B)$$

$$P(B/A) = \frac{P(B) * P(A/B)}{P(A)} \rightarrow \text{conditional probability.}$$

↙ Baye's Theorem.

Naive Bayes classifier.

$$x_1, x_2, x_3, \dots, x_n$$

Independent var

$$y$$

dependent var

$$\Rightarrow \begin{cases} y=1 & \text{Yes} \\ y=0 & \text{No} \end{cases}$$

$$P(y \underset{\text{yes}}{\overset{\text{no}}{\wedge}} / (x_1, x_2, \dots, x_n)) = \frac{P(y) * P((x_1, x_2, x_3, \dots, x_n) / y)}{P(x_1, x_2, \dots, x_n)}$$

$$= P(y) * P(x_1/y) * P(x_2/y) * P(x_3/y) * \dots * P(x_n/y)$$

$$\text{constant} \quad P(x_1) * P(x_2) * P(x_3) * \dots * P(x_n)$$

$$P(y=\text{no} / (x_1, x_2, \dots, x_n)) = \frac{P(0) * P(x_1/0) * P(x_2/0) * \dots * P(x_n/0)}{P(x_1) * P(x_2) * P(x_3) * \dots * P(x_n)}$$

$$\text{constant} \quad P(x_1) * P(x_2) * P(x_3) * \dots * P(x_n) \quad \# \text{fixed}$$

$$P(y=\text{Yes}/x_i) = 0.13$$

$$P(y=\text{No}/x_i) = 0.05$$



Need to do Normalization.

$\geq 0.5 \Rightarrow \text{class } 1$

$< 0.5 \Rightarrow \text{class } 0$

$$\left\{ \begin{array}{l} P(y=\text{Yes}/x_i) = \frac{0.13}{0.13+0.05} = 0.72 = 72\% \\ P(y=\text{No}/x_i) = 1 - 0.72 = 0.28 = 28\% \end{array} \right.$$

Normalized

$$P(y/x_i) = p(y) * \prod_{i=1}^n P(x_i/y) \quad y = \{c_k\} \text{ class labels}$$

$$= p(c_k) * \prod_{i=1}^n P(x_i/c_k)$$

$$P(c_k/x_i) = p(c_k) * \prod_{i=1}^n P(x_i/c_k)$$

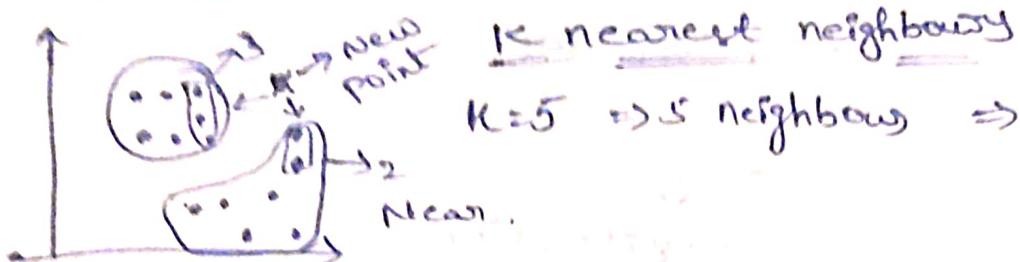
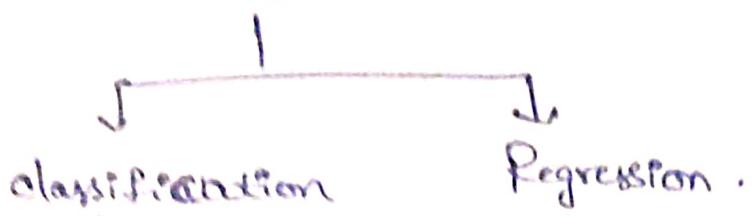
↓
pro of class
given input
features

↓
Probability
of class c_k

↑ feature
↓ pro x_i given by
class c_k

$$P(x_i/c_k) = \prod_{i=1}^n P(x_i/c_k) \Rightarrow \text{naive Assumption}$$

④ KNN Algorithm (Used for both)



- * New point goes to which class decision taken by distance.

Euclidean distance

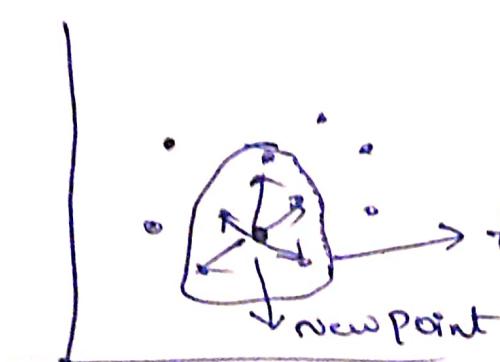
$$(x_1, y_1) \rightarrow (x_2, y_2)$$

$$\Rightarrow \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Manhattan distance

$$\rightarrow |(x_2 - x_1) + (y_2 - y_1)|$$

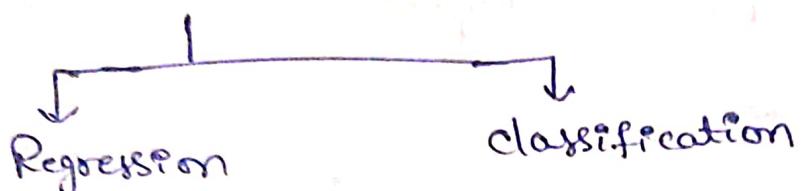
Regression



K Nearest work bad with

- * Outliers
- * Imbalanced dataset

⑤ Decision Tree (used for both)



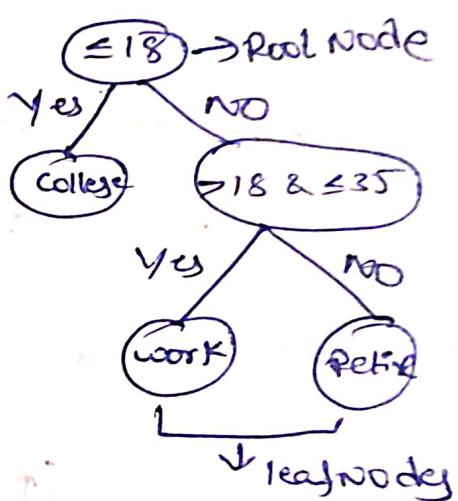
if (age ≤ 18):

 print ("college")

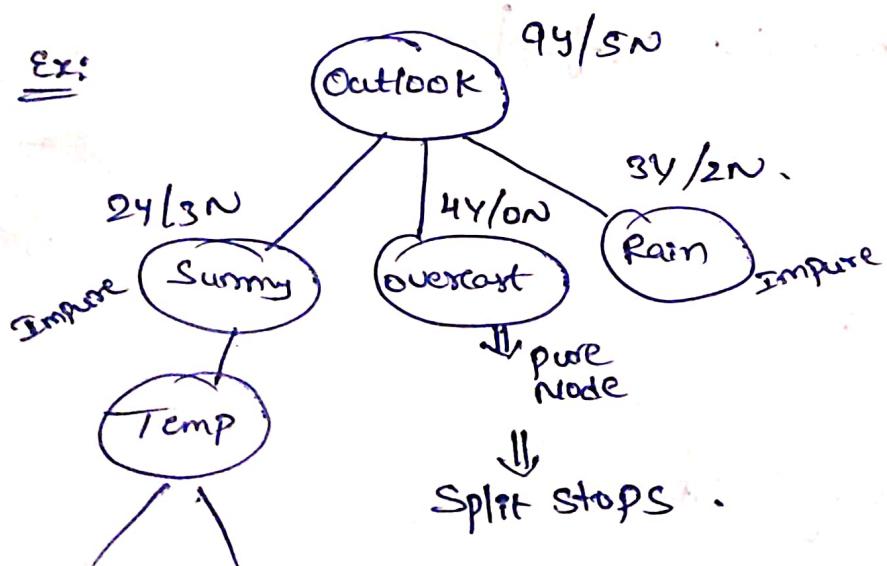
elif (age > 18 & age ≤ 35): \Rightarrow
 print ("work")

Else:
 print ("retire")

Decision Tree



Ex:



Pure split \rightarrow zero ND's

Impure split \rightarrow Yes

\Rightarrow until pure node i.e leaf node
comes this process repeats.

① How we Calculate Purity?

① Entropy

② Gini Coefficient (or) Gini Impurity

② How next node or root node features are selected.

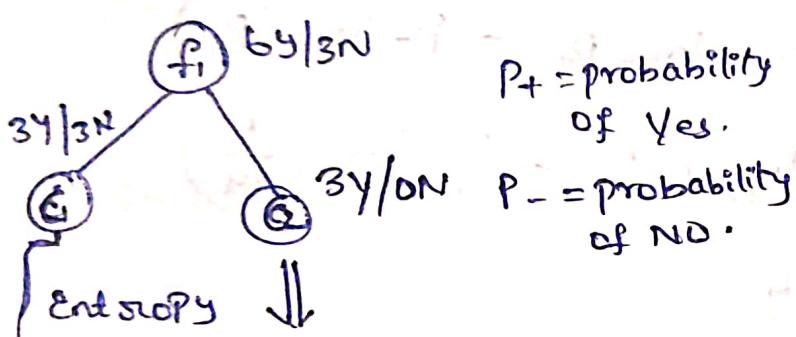
① Information Gain.

Entropy

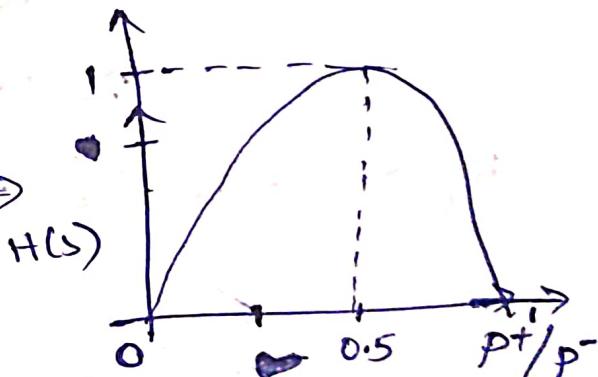
$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

Gini Impurity

$$G.I = 1 - \sum_{i=1}^n (P_i)^2$$



$$\begin{aligned} H(S) &= -\frac{3}{3} \log_2 \frac{3}{3} - \frac{0}{3} \log_2 \frac{0}{3} \\ &= -1 \log_2 1 = 0 \end{aligned}$$



$$\begin{aligned} H(S) &= -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} \\ &= 1 \end{aligned}$$

* Entropy is Between 0 & 1.

$$P^+ = 0.5$$

$$P^- = 0.5$$

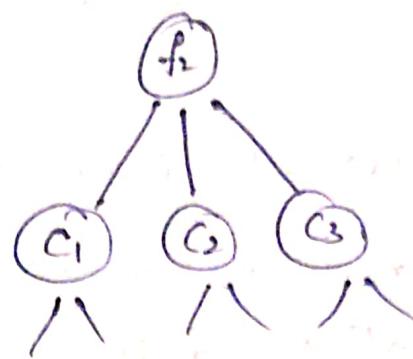
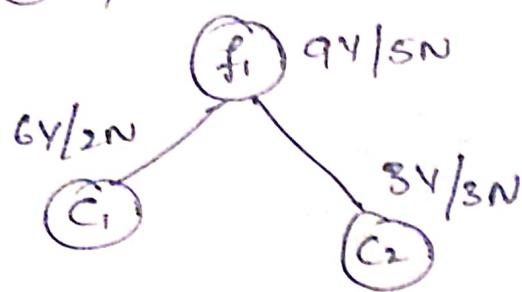
$$\downarrow$$
$$H(S) = 1.$$

Impure split.

$$H(S) = 0$$

Pure split.

② feature selection to Splitting



Information Gain

$$\text{Gain}(S, f_1) = H(S) - \sum_{v \in \text{eval}} \frac{|S_v|}{|S|} H(S_v)$$

$$\begin{aligned} \text{Entropy of root node} &= -P + \log_2 P_+ - P_- \log_2 (P_-) \\ &= -\frac{9}{14} \log_2 \left(\frac{9}{14}\right) - \frac{5}{14} \log_2 \left(\frac{5}{14}\right) \\ &\approx 0.94 \end{aligned}$$

$$H(S_v) \Rightarrow H(S_{vC_1}) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} \\ \approx 0.81$$

$$\Rightarrow H(S_{vC_2}) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{8} \log_2 \frac{3}{8} = 1.$$

$$\begin{aligned} \text{Gain}(S, f_1) &= 0.94 - \left[\frac{8}{14} * 0.81 + \frac{6}{14} * 1 \right] \\ &= 0.049. \end{aligned}$$

$\text{Gain}(S, f_2) = 0.051 \Rightarrow$ Select this.

$$\text{Gain}(S, f_2) > \text{Gain}(S, f_1)$$

✓

Gini Impurity

$$G \cdot I = 1 - \sum_{i=1}^3 (P_i)^2$$

$$= 1 - [(P_+)^2 + (P_-)^2]$$

$$= 1 - \left[\left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right] = 1 - \left[\frac{1}{2}\right] = 0.5$$

\Downarrow
but in Entropy = 1.

* Entropy takes more time to compute because of \log .

\Downarrow
So, when Decision Tree is more complex

\Downarrow

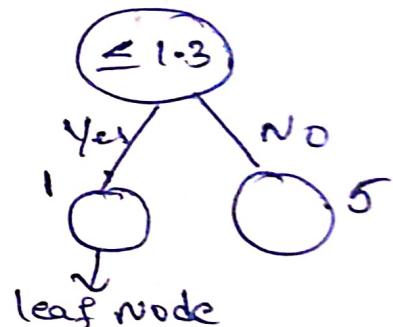
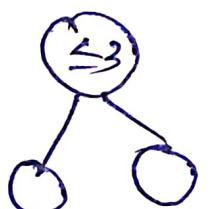
Use Gini Impurity.

Fast

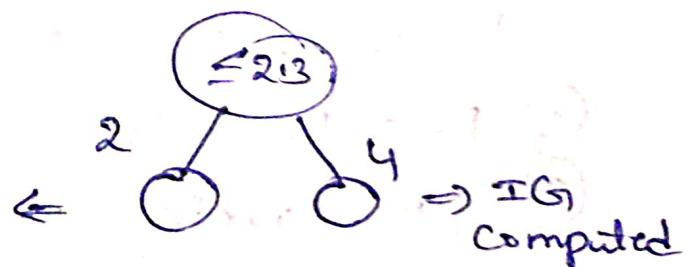
Gini $>$ Entropy.

Numerical Features

<u>f_1</u>	<u>%p</u>	$\Rightarrow f_1$ (sorted)
2.3		1.3
1.3		2.3
4		3
5		4
7		5
3		2



\Downarrow
Information Gain
computed ≤ 1.3



\Downarrow
After all this for which root node IG is high
that is Selected.

Decision Tree Regressor

$f_1, f_2 \text{ / } \%$

continuous

Mean
calculated

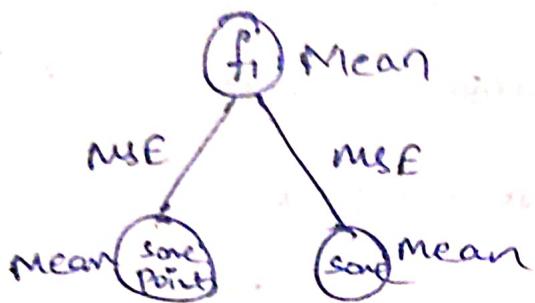


Mean



MSE (or) MAE

$$= \frac{1}{2m} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

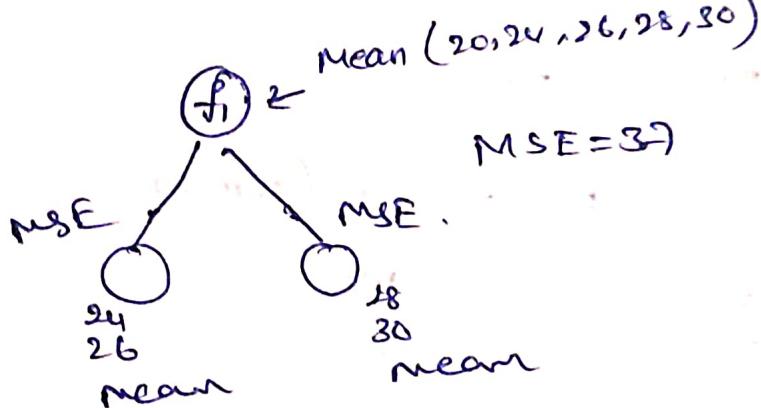


MSE reducing

going to reach leaf node.

Hyper Parameters

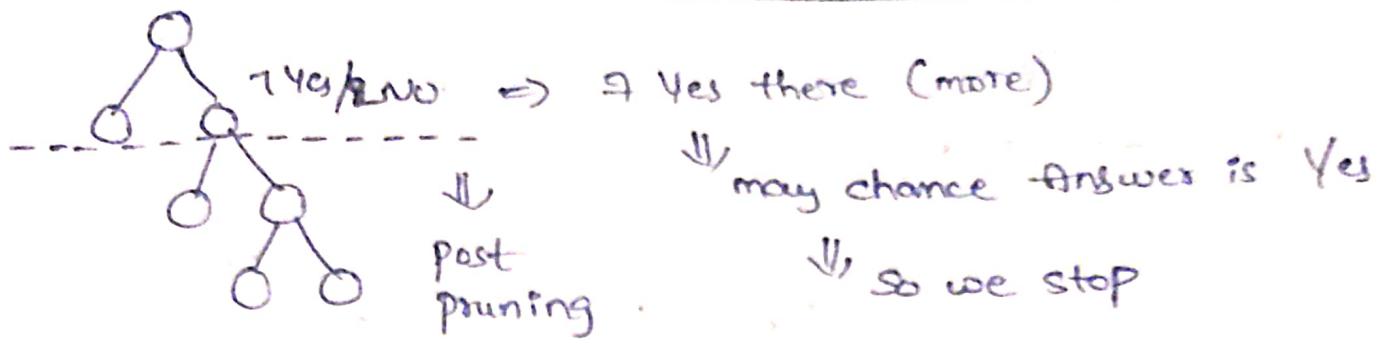
f_1	%
C_1	20
C_2	24
-	26
-	28
-	30



Hyperparameters

Decision Tree \Rightarrow Overfitting to avoid.

- ① post pruning
- ② pre pruning

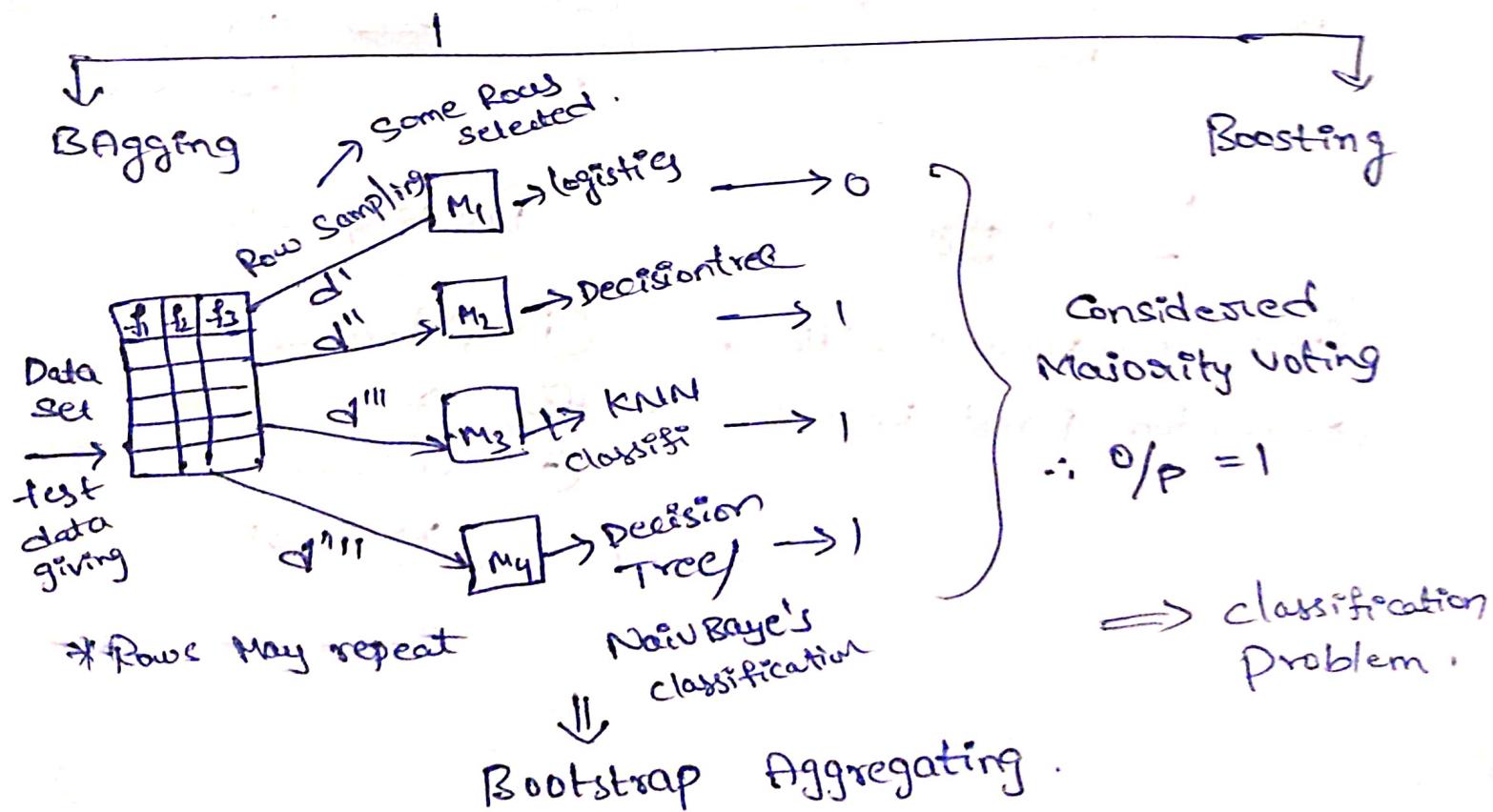


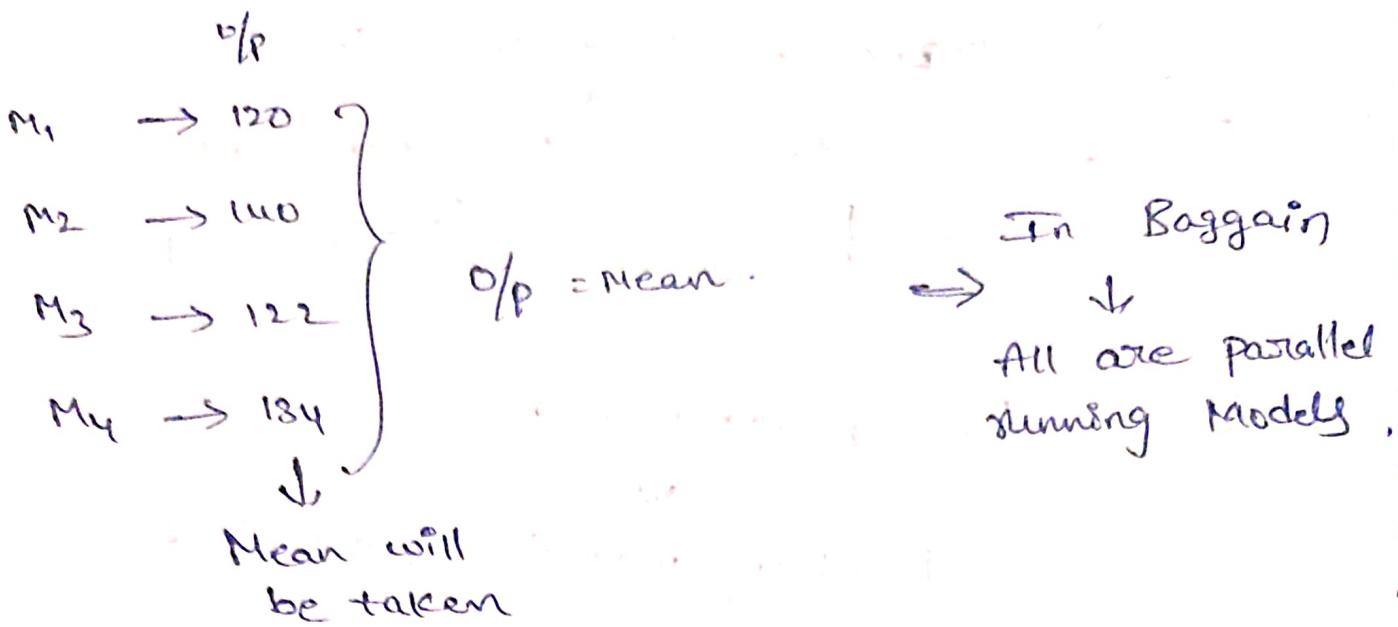
Prepruning done by GridsearchCV

\Downarrow
Based on max-depth, max-leaf

Ensemble Techniques

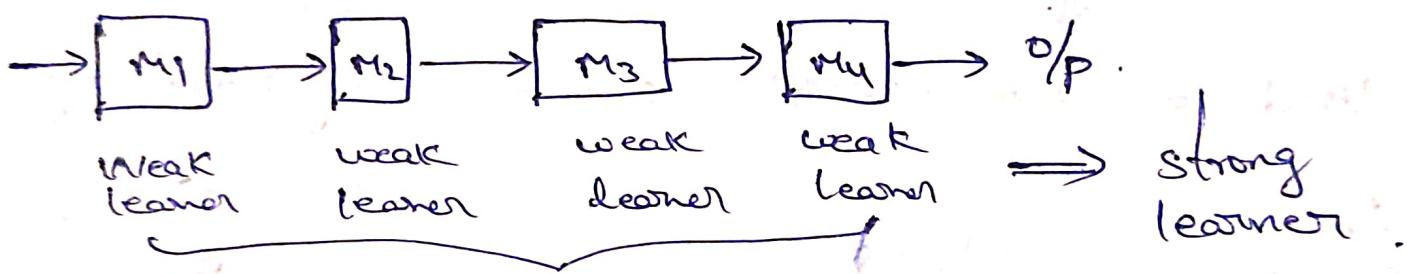
\Downarrow
Used to solve Multipl' a problem using Multiple Algorithms.





Boosting

⇒ Sequential Models Combined to give strong learner.



Bagging Algo

- ① Random Forest classifier
- ② Random Forest Regressor

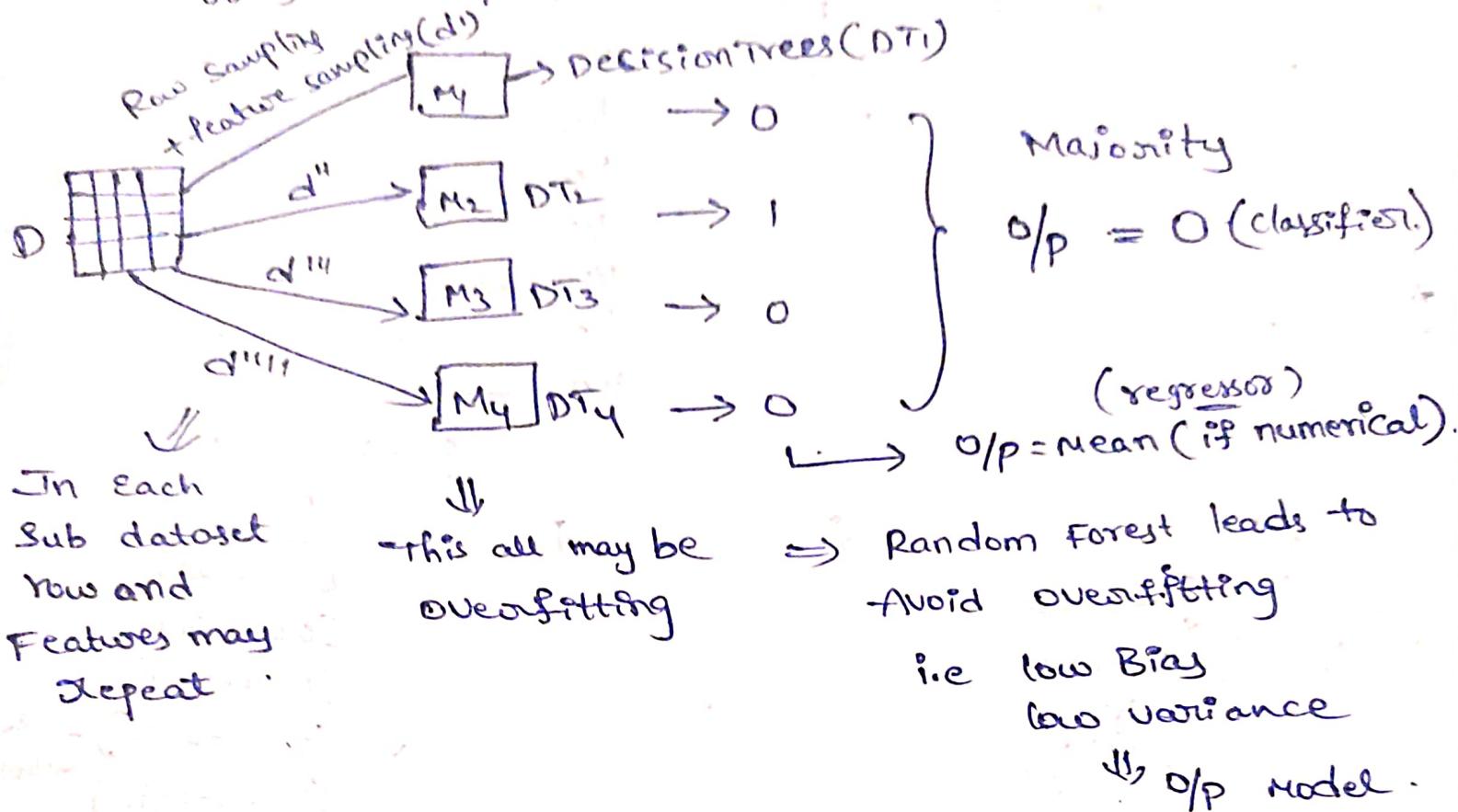
Boosting Algo

- ① AdaBoost
- ② Gradient Boost
- ③ Extreme gradient

xg Boost

⑥ Random Forest classifier & Regressor

→ Bagging Technique. \Rightarrow All Models are Decision Trees



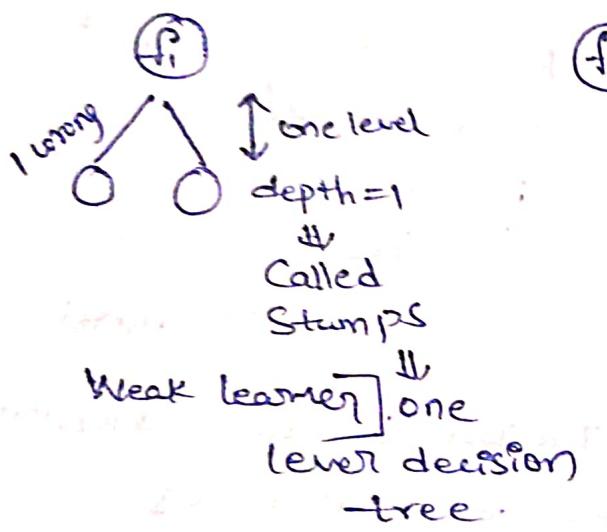
- * Normalization Not required in Decision Tree or random forest
- * Standardization needed in KNN for distances.
- * Random Forest not impacted by Outliers.
- * KNN impacted by Outliers.

② Boosting

i) AdaBoost

	f_1	f_2	f_3	f_4	O/P
1	-	-	-	-	Y ₁
2	-	-	-	-	NO
3	-	-	-	-	-
4	X	-	-	-	-
5	-	-	-	-	-
6	-	-	-	-	-
7	-	-	-	-	-
8	-	-	-	-	-
9	-	-	-	-	-
10	-	-	-	-	-

Define weight	Updated	Normalized weight
$\frac{1}{7} \rightarrow 0.05$		0.07
$\frac{1}{7} \rightarrow 0.05$		0.07
$\frac{1}{7} \rightarrow 0.05$		0.07
$\frac{1}{7} \rightarrow 0.349$		0.532
$\frac{1}{7} \rightarrow 0.05$		0.07
$\frac{1}{7} \rightarrow 0.05$		0.07
$\frac{1}{7} \rightarrow 0.05$		0.07
	<u>0.649</u>	
		↓



f_3

Buckets

[0 - 0.07]

[0.07 - 0.14]

[0.14 - 0.21]

✓ [0.21 - 0.749] → diff 0.349
[0.749 - 0.751]

=

Correct Records

③ New sample = weight of e^{-P_s}

$$= \frac{1}{7} \times e^{-0.895} = 0.05$$

Incorrect records = weight * e^{P_s}

$$= \frac{1}{7} * e^{0.895} = 0.349$$

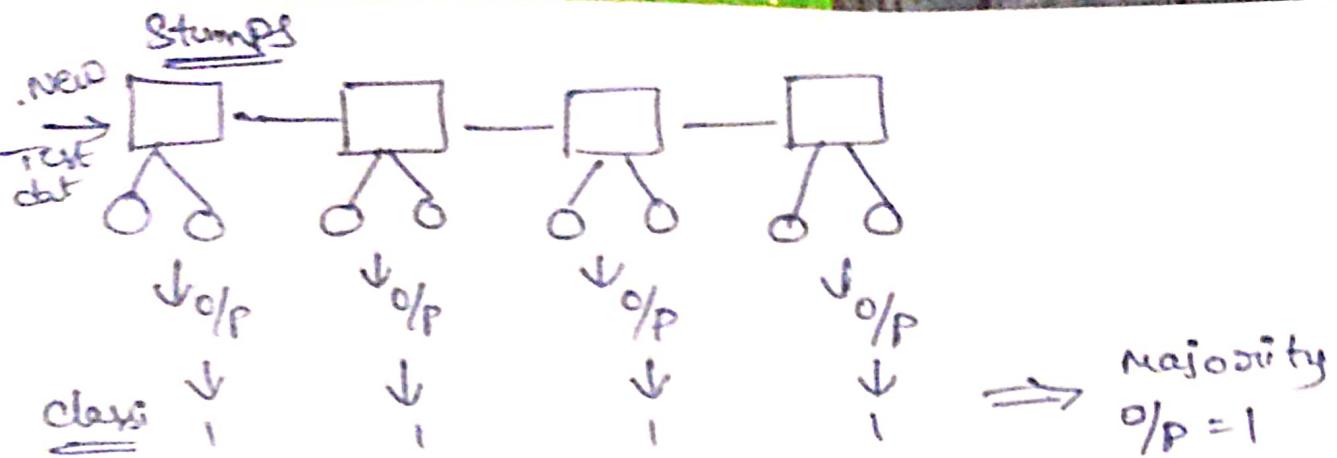
① Total Error = $\frac{1}{7}$

② Performance of stump

$$= \frac{1}{2} \log_e \left(\frac{1-TE}{TE} \right)$$

$$= \frac{1}{2} \log_e \left(\frac{1-\frac{1}{7}}{\frac{1}{7}} \right)$$

$$= 0.895$$



Regression \Rightarrow o/p = mean.

Black box models vs white box Models

- Linear Regression \rightarrow white box (able to see all β_i changing)
 - Random Forest \rightarrow Black box (impossible to see all decision trees & working)
 - Decision Tree \rightarrow white box (Able to know how data splitting)
 - ANN \rightarrow Black box (nodes & weights not able to see).

Probabilities

① Classical probability $P(E) =$

$$\frac{\text{No. of favorable outcome}}{\text{Total no. of outcomes}}$$

Ex: $P(H) = \frac{1}{2}$
 \downarrow
 Head

② Empirical:

$$P(E) = \frac{\text{No. of times event occurred}}{\text{Total No. of trials}}$$

Ex: If you roll a die 100 times get 4 on 18 times

$$P(4) = \frac{18}{100} = 0.18$$

③ Subjective

→ Based on personal belief, intuition or expert judgment,
 not actual data.

Ex: A doctor believes there is 70% chance a
 patient will recover.

④ Axiomatic

Based on mathematician Andrey Set of axioms & rules.

Axioms

1. $0 \leq P(E) \leq 1$

2. $P(S) = 1$, where S is the sample space
 \downarrow
 possible outcomes set

3. For any two mutually exclusive events A and B .

Not occurred at same time

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

Bayesian statistics

① prior probability - $P(A)$

* The probability of an event before observing new data

$$P(\text{Disease}) = 0.01$$

② likelihood - $P(B|A)$

* data under the assumption. (probability based on condition)

③ posterior - $P(A|B)$

* updated probability of the hypothesis after seeing new data.

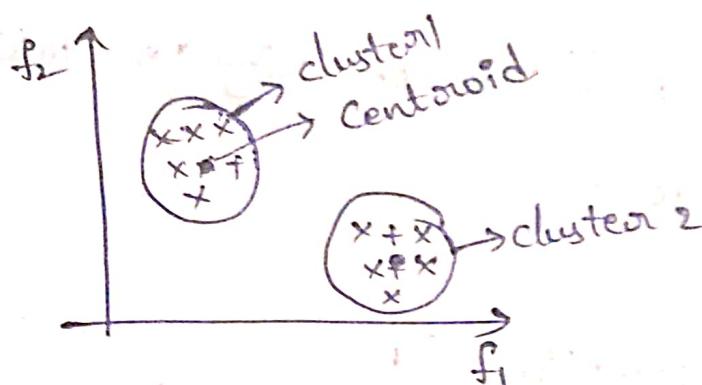
$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

④ Marginal probability. - $P(B)$

The total probability of observing the evidence, considering all possible causes.

Unsupervised Algorithms :- Don't have specific output

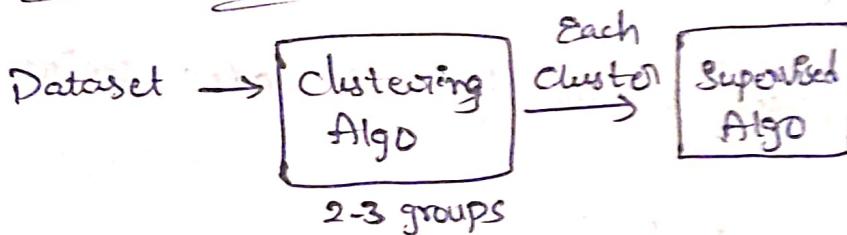
⑧ K-Means clustering \Rightarrow clustering



↓
Similar kind of data
groups -

$\Leftrightarrow K=2$

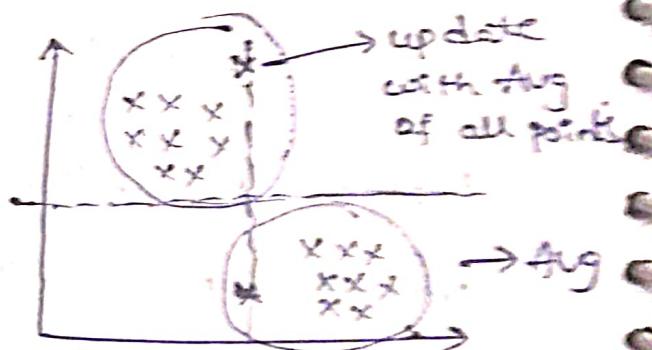
Custom Ensemble technique :-



K Mean



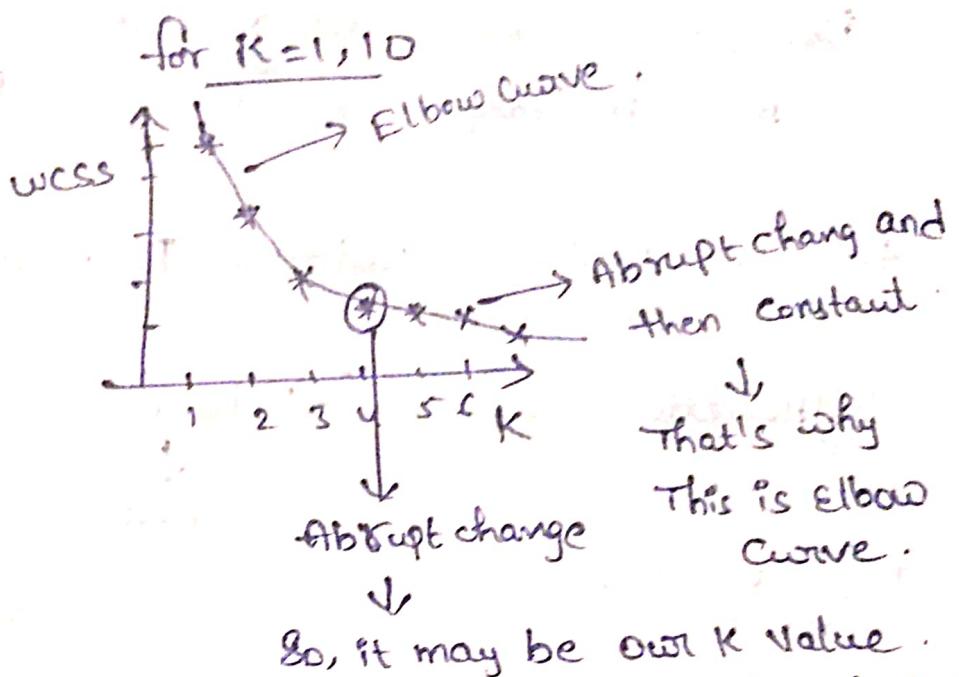
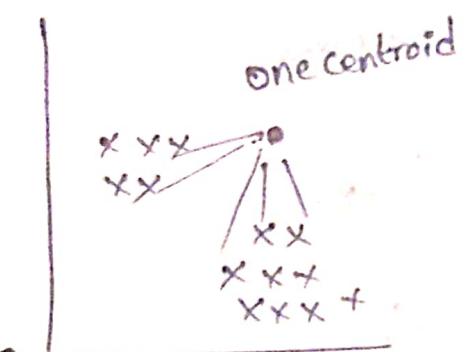
$K = \text{Centroids}$



- ① We try different K values \Rightarrow suitable $K=2$
- ② Initialize K no. of Centroids.
- ③ Groups Centroids nearest datapoints as one cluster using Euclidean distance formula.
- ④ Compute the Avg to update Centroids.
* It's better to initialize Centroids very far to get Centroid Exactly in center. But it may be chance of getting More clusters. So, KMeans++ (Ensuring Centroids are far).

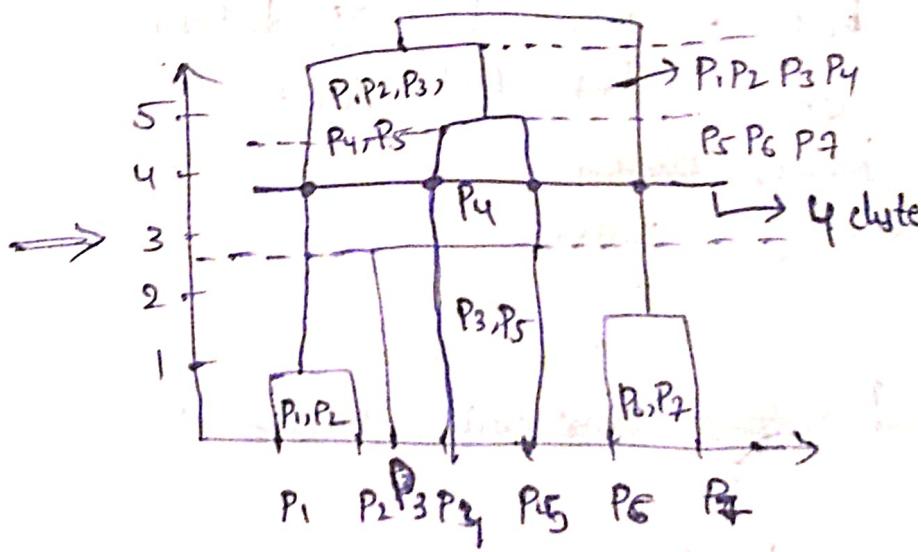
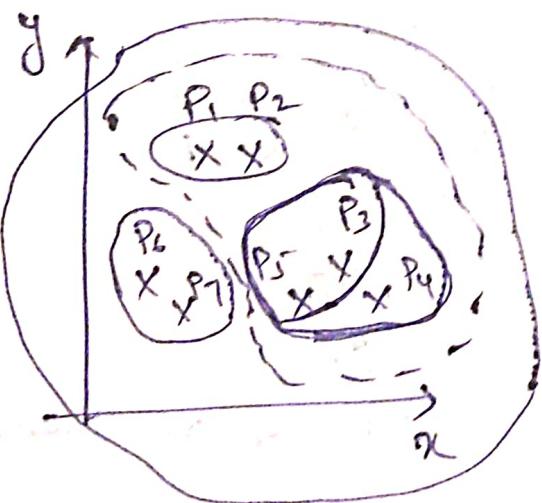
How we decide K Value?

↓
Elbow Method.



⑨ Hierarchical clustering

↳ step by step process. → Evaluate or Identify nearest points & Combine



* How to find how many groups there?

⇒ Need to find the longest vertical line that has no horizontal line passes through it.

↓
Dendogram.

* Max time is taken by hierarchical clustering than K-Means clustering.

Dataset is small \Rightarrow go with hierarchical
" is large \Rightarrow use " K-Means.

How we validate clustering Problem (Evaluation Metric?)

Silhouette Score (-1 to 1)

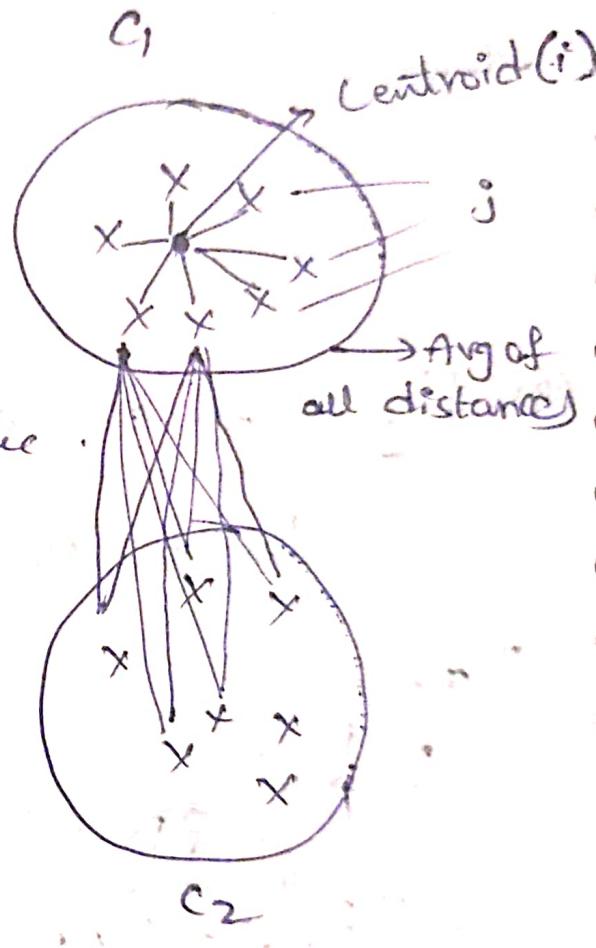
↓
Best

For K-Means

$$① a(i) = \frac{1}{|C_i|-1} \sum_{j \in C_i, i \neq j} d(i, j) \Rightarrow \text{Avg distance}$$

$$b(i) = \min_{j \neq i} \frac{1}{|C_j|} \sum_{j \in C_j} d(i, j)$$

↓
Nearest cluster
among All clusters



For good clustering Model

$$b(i) \gg a(i)$$

↓
This means distance b/w clusters are more
That is clearly grouped as diff clusters.

* $a(i) \gg b(i) \times$ bad.

Silhouette clustering formula

8

$$S(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$

$b(i) > a(i)$
 \uparrow
 $a(i) > b(i)$
 \downarrow

DBScan clustering

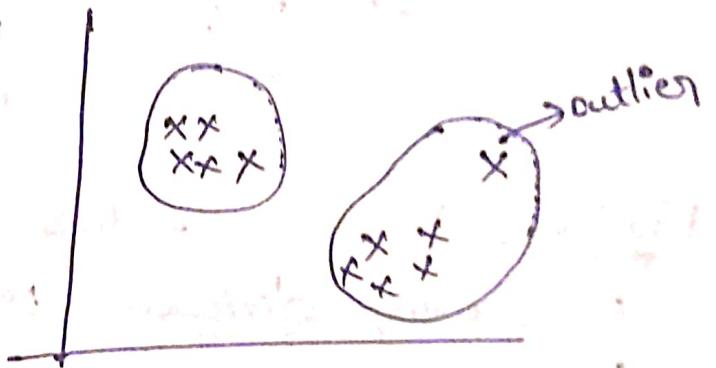
Density-Based Spatial clustering of Applications with Noise

⇒ 1) Min points

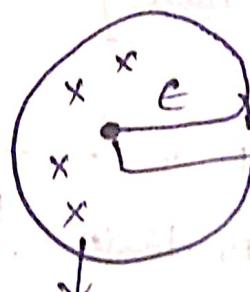
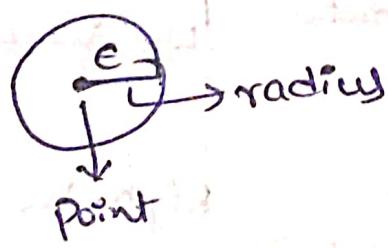
2) Core points

3) Border points

4) Noise points

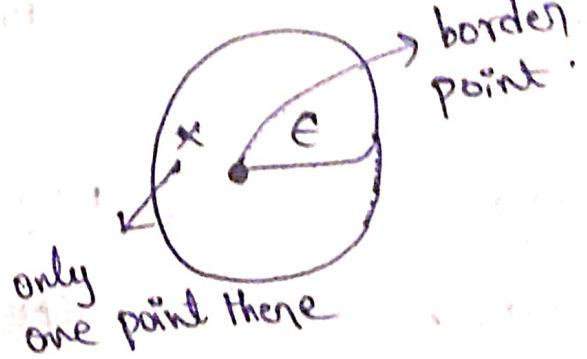


Min pts = 4 ⇒ Hyperparameter.

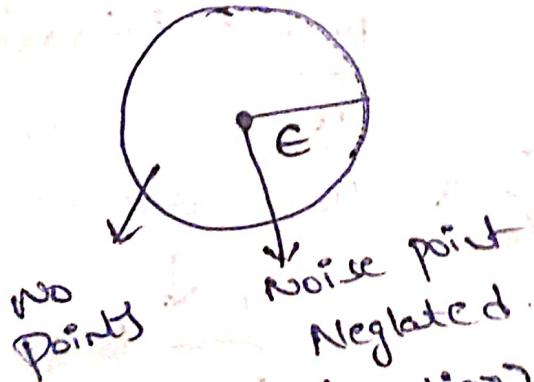


Min points = 4

Core point
because min point = 4



only one point there



No Points

Neglected
(outlier)

never taken into group

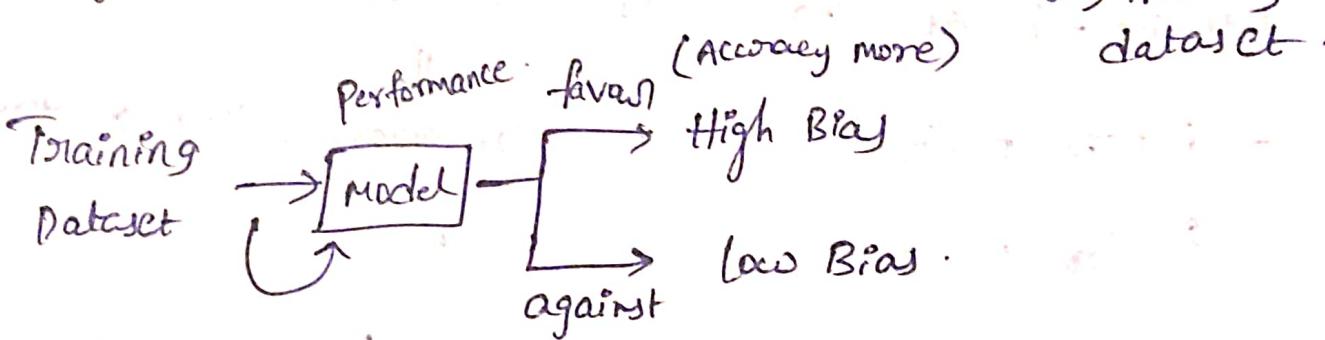
Depth of Bias and Variance

Training Dataset = 90% } \Rightarrow overfitting

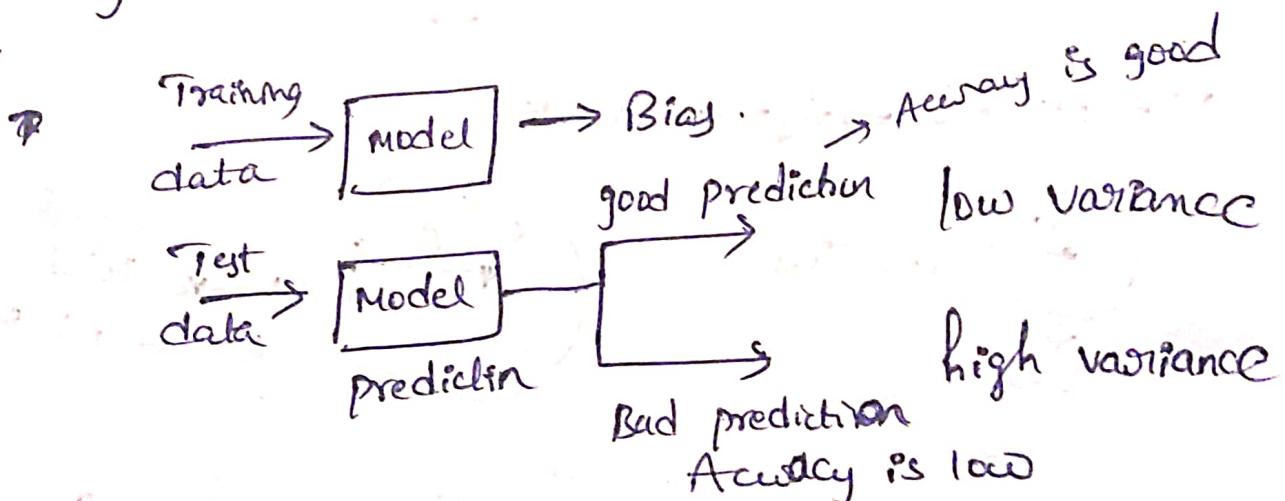
Test dataset = 10%

↓
low Bias
high Variance

Bias :- It is a phenomenon that skews the result of an algorithm in favour or against an idea.



Variance: Variance refers to the changes in the model when using different portions of the training or test data.



Model 1

Train acc = 90%

Test acc = 75%

↓
low Bias

high Variance

Model 2

Train acc = 60%

Test acc = 55%

↓
High Bias

high Variance

Model 3 (Generalize model)

Train acc = 90%

Test acc = 92%

↓
Low Bias

Low Variance

⑪ Xgboost classifier \rightarrow extreme Gradient Boosting

Salary	credit	Approved	Final	Classification
$<=50$	B	0	Regression	
$<=50$	G1	1		
$<=50$	G1	1		
>50	B	0		
>50	G1	1		
>50	Normal	1		
<50	N	0		

Base model \rightarrow probability = 0.5 (for classification)

Weak learner

Residual

$$0 - 0.5 = -0.5$$

$$0.5$$

$$0.5$$

$$-0.5$$

$$0.5$$

$$0.5$$

$$-0.5$$

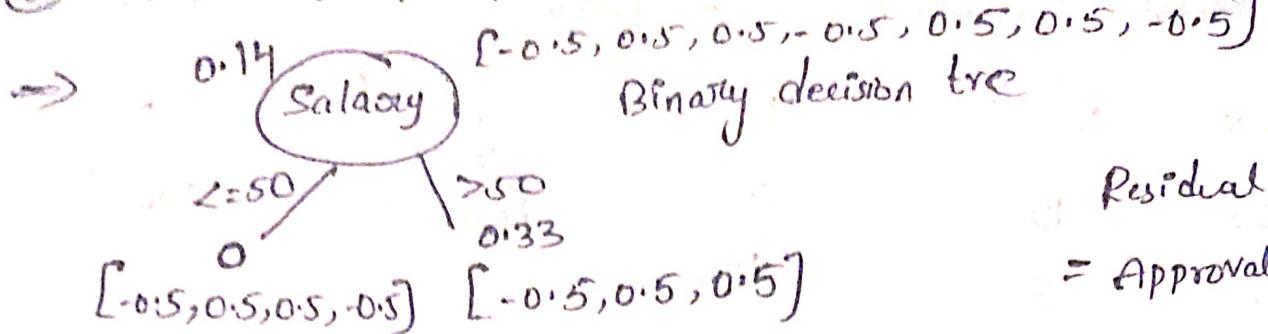
Creating decision tree after base model.

① Create a Binary Decision Tree using the features.

② calculate Similarity weight

$$= \frac{\sum (\text{Residual})}{\sum (pr(1-pr) + \lambda)}$$

③ Calculate Information Gain



Residual

= Approval-probability

⇒ Calculate similarity weight



$$= \frac{(-0.5 + 0.5 + 0.5 - 0.5)^2}{0.5(1-0.5) + 0.5(1-0.5)} \\ + 0.5(1-0.5) + 0.5(1-0.5)$$

$$= 0$$

$x=0$



$$= \frac{(-0.5 + 0.5 + 0.5)^2}{0.5(1-0.5) + 0.5(1-0.5)} \\ + 0.5(1-0.5)$$

$$= \frac{0.25}{0.75} = 0.33$$

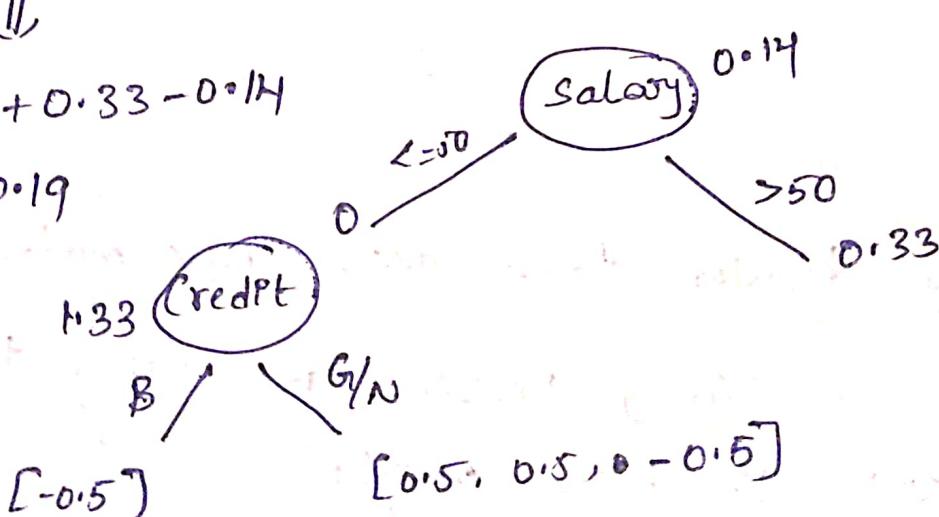
Root node similarity weight = $\frac{0.25}{1.75} = \frac{1}{7} = 0.14$

⇒ Information gain.



$$= 0 + 0.33 - 0.14$$

$$= 0.19$$



$$Smw = \frac{0.25}{0.5(1-0.5)} = \frac{0.25}{0.25} = 1$$

$$\frac{0.25}{0.75} = 0.33$$

$$\Rightarrow 1 + 0.33 \\ - 0 \\ = 1.33$$

In Case of Base model.

$$\log\left(\frac{P}{1-P}\right) = \log\left(\frac{0.5}{0.5}\right) = 0$$

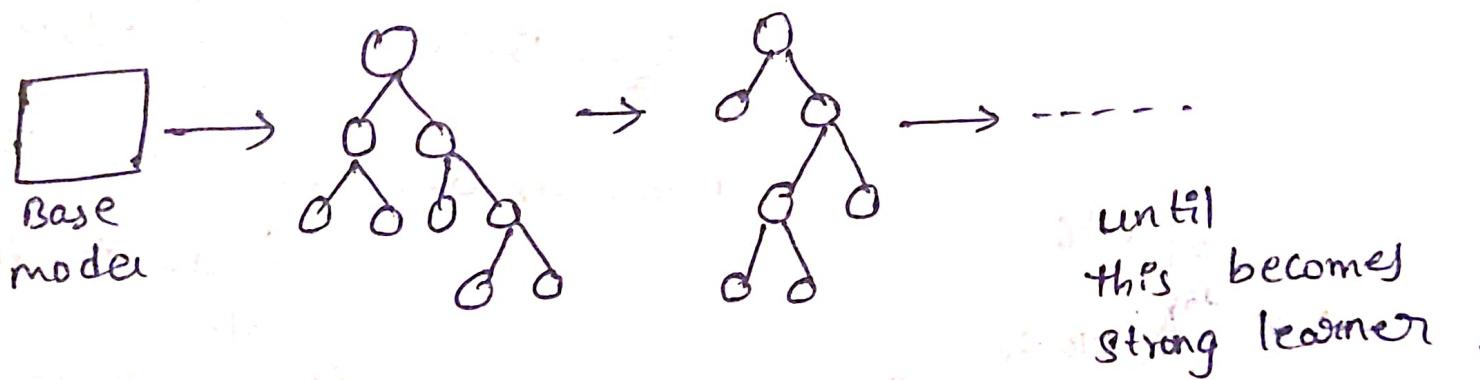
$\sigma [0 + \alpha(1)]$ ^{similarity}
weight
 \downarrow learning rate (0 to 1)

$\alpha \Rightarrow$ fixed by
Cross Validation

activation function.

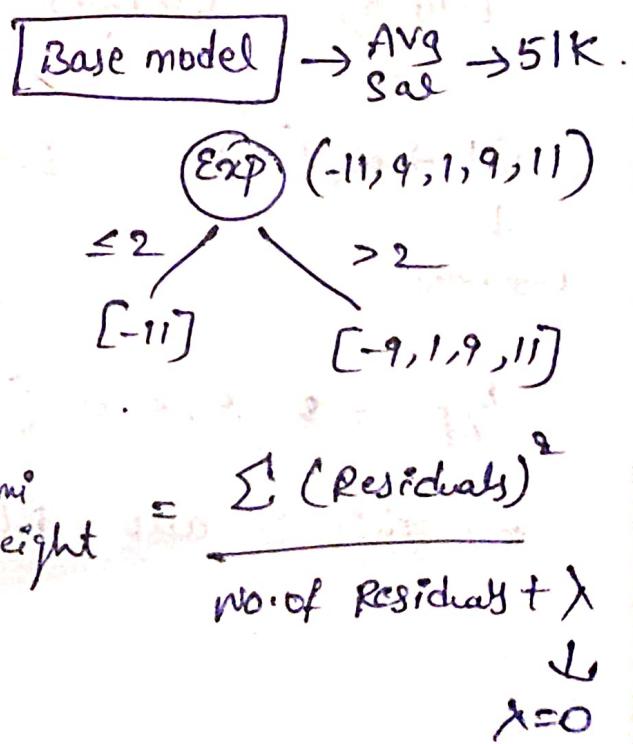
$$\sigma [0 + \alpha_1(DT_1) + \alpha_2(DT_2) + \alpha_3(DT_3) + \alpha_4(DT_4) + \dots + \alpha_n(DT_n)]$$

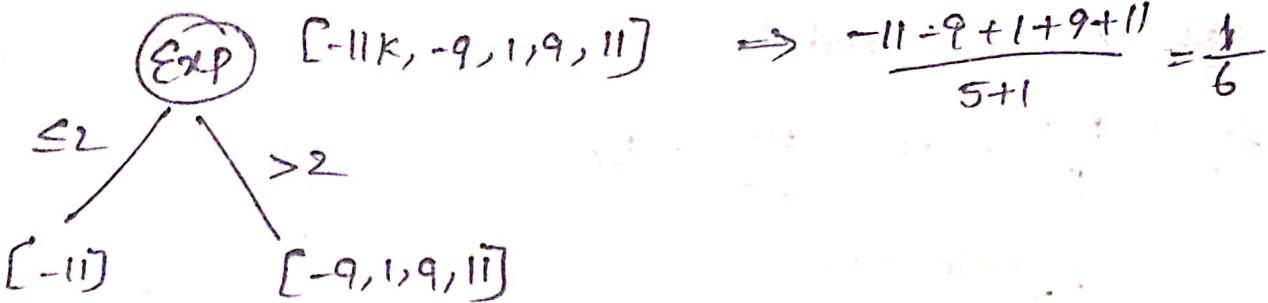
* xgBoost is a Black Box model (prone to overfitting)



xgboost Regressor (Bagging)

Exp	Graph	sal	Resid
2	Yes	40K	-11K
2.5	Yes	40K	-10
3	No	52K	1
4	No	60K	9
4.5	Yes	62K	11





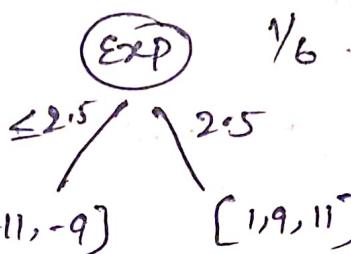
$$\frac{|2|}{1+1} = 6.05 \quad \frac{(-9+1+9+11)^2}{4+1} = \frac{144}{5} = 28.5 \quad (\because x=1)$$

Informant

$$\text{Gain} = 6.05 + 28.5 - \frac{1}{6}$$

~~$= -0.33$~~

$= 89.13$



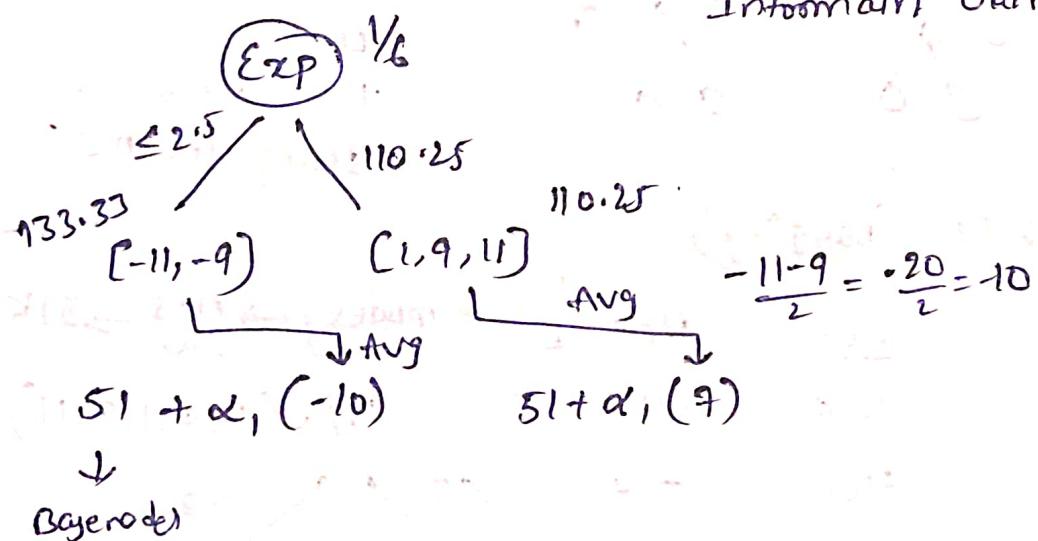
$$\text{Sinfo} = \frac{(-11-9)^2}{2+1} = \frac{400}{3} = 133.33$$

$$\text{Sum} = \frac{(1+9+11)^2}{3+1} = \frac{21^2}{4} = 110.25$$

$$\text{Informant Gain} = (1.33 + 110.25 - \frac{1}{6})$$

greater than previous

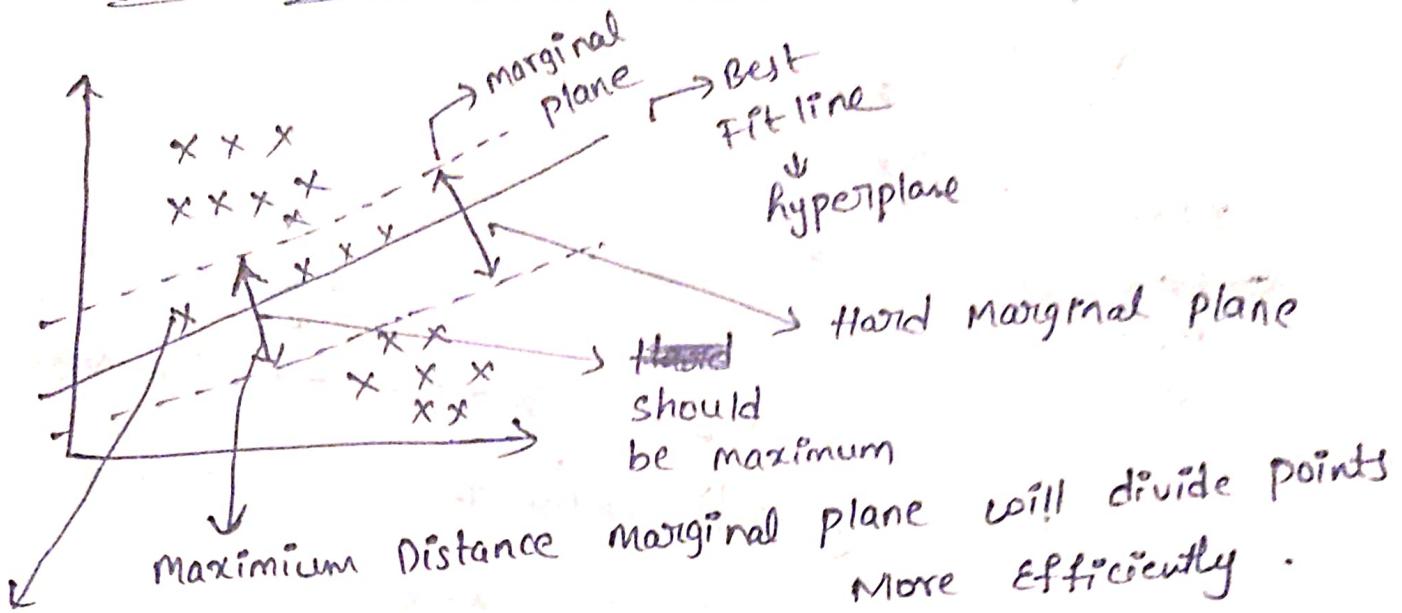
↓
so we choose
this split.



$$D/P = 51 + \alpha_1(DT_1) + \alpha_2(DT_2) + \alpha_3(DT_3) + \dots + \alpha_n(DT_n)$$

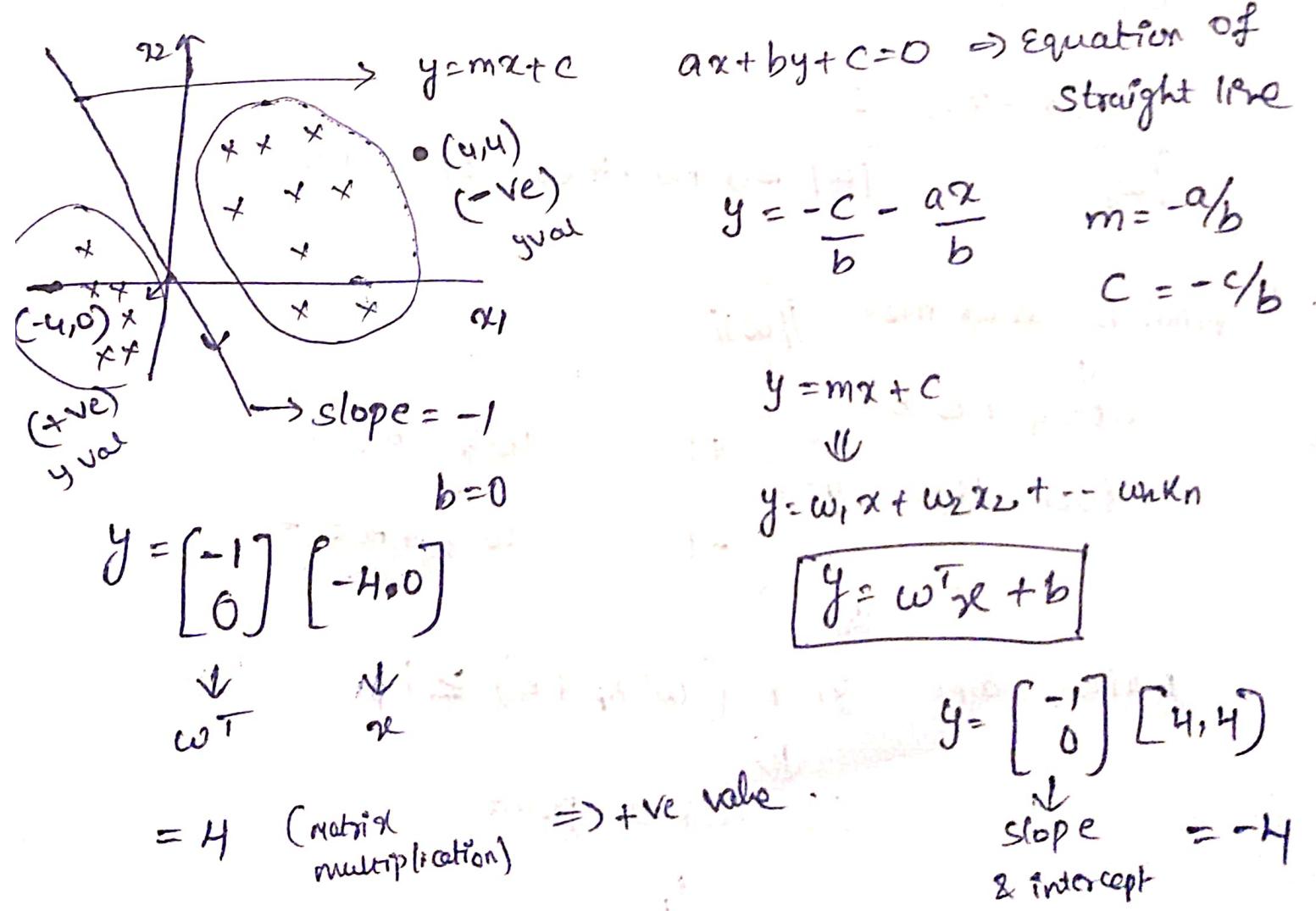
⇒ This is also Black box model (Can't visualize all of them).

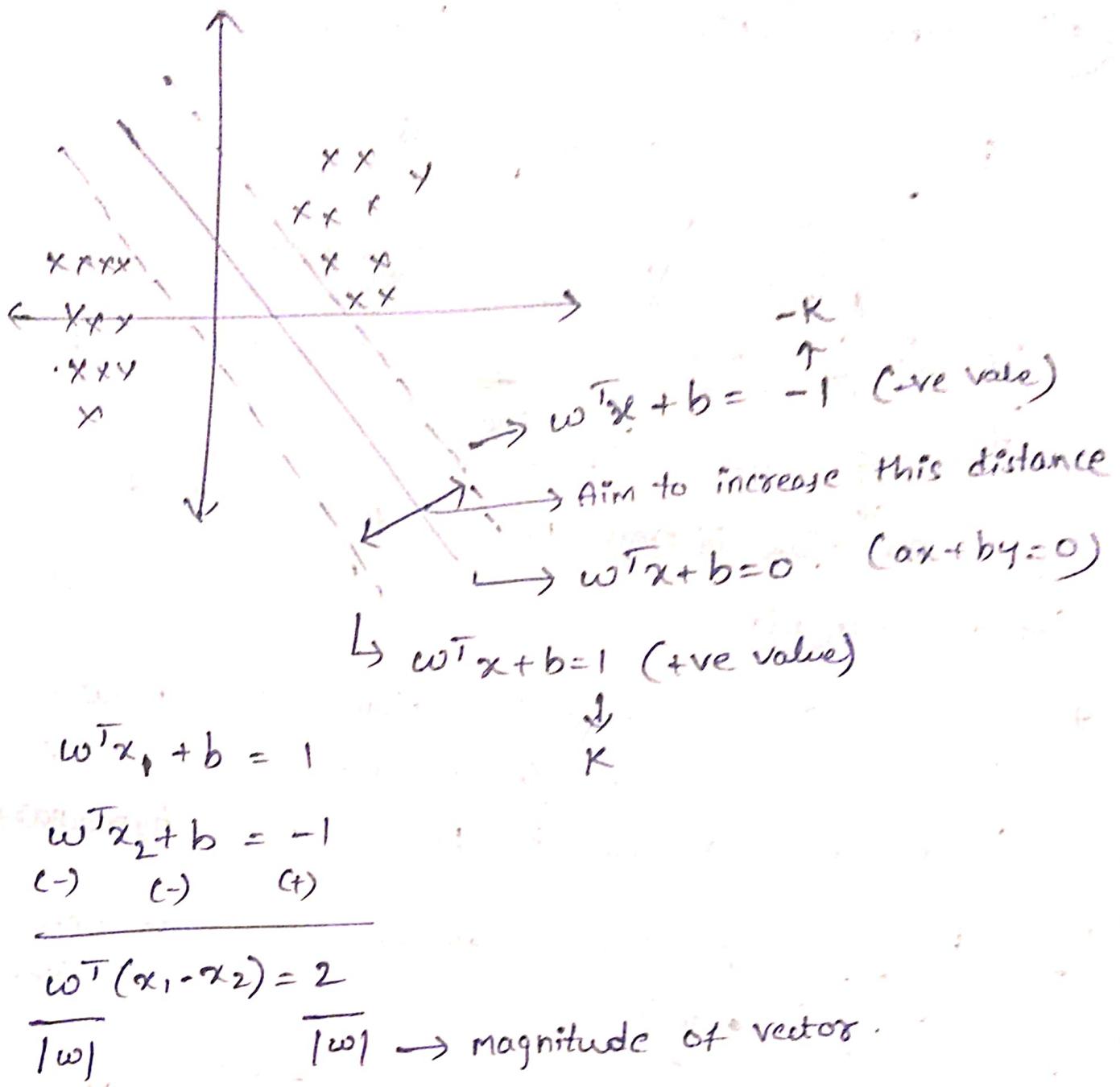
(12) SVM (Support Vector Machine)



Soft marginal plane:

* Aim is creating marginal plane with maximum distance Even though there are some errors we consider it in solving by proving some hyperparameters





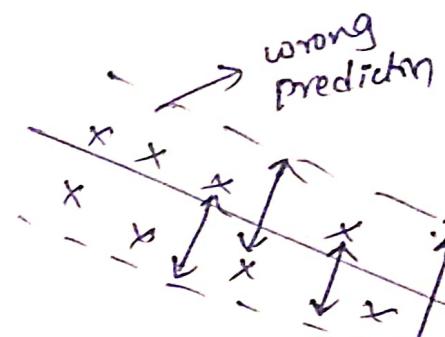
Aim is maximize $\frac{2}{\|w\|}$
 using (w, b)
 such that $y_i \begin{cases} +1 & w^T x_i + b \geq 1 \\ -1 & w^T x_i + b \leq -1 \end{cases}$
 Major aim $y_i * (w^T x_i + b) \geq 1$
 for correct point

$$\text{maximize}_{(\omega, b)} \frac{2}{\|\omega\|} \iff \min_{(\omega, b)} \frac{\|\omega\|}{2}$$

$$\min_{(\omega, b)} \frac{\|\omega\|}{2} + C \sum_{i=1}^n \xi_i$$

↓

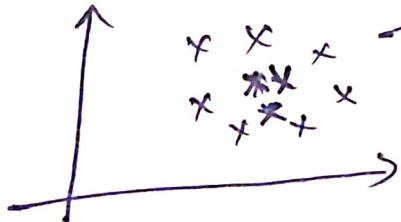
How many errors we can have \Rightarrow summation of the distance of the wrong point



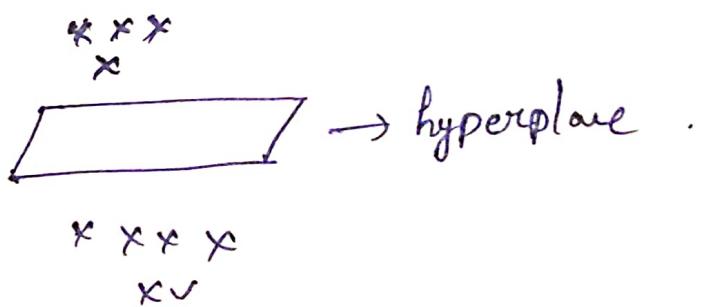
SVM

$$C \sum_{i=1}^n \xi_i \Rightarrow \text{only will change}$$

SVM Kernel



Cannot do with straight line
so we convert it int 3 dimension

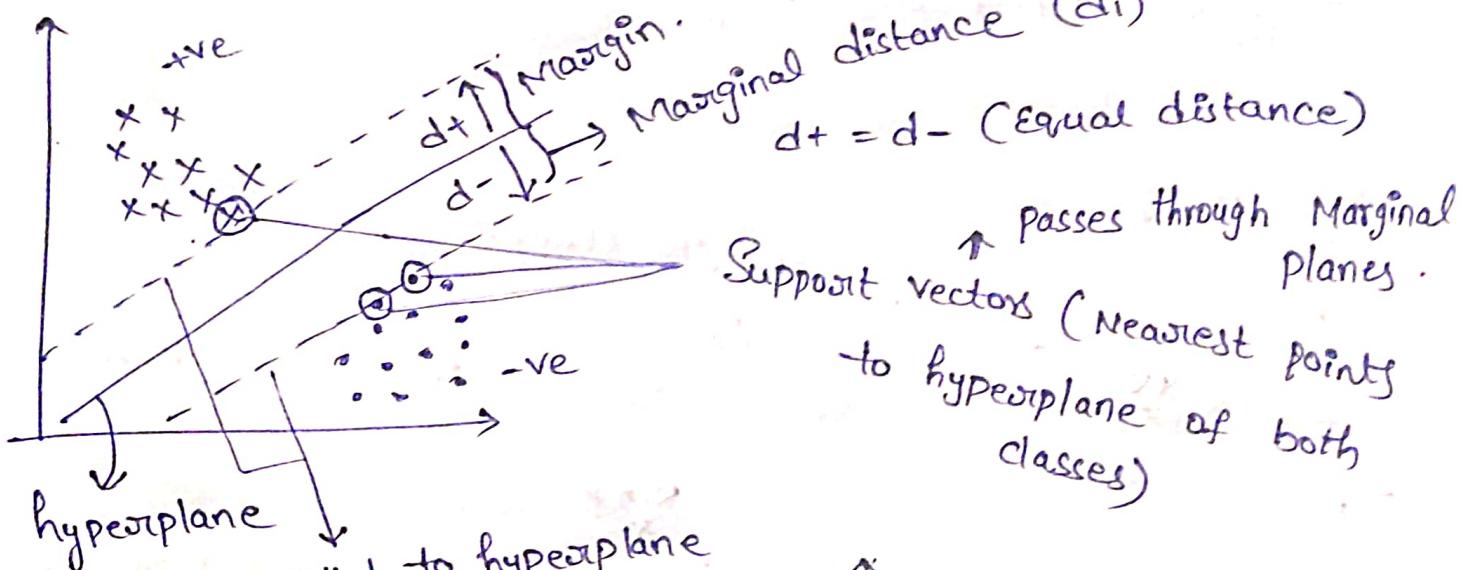


Support Vector Machines

→ It is used to solve both classification & regression problems.

- 1) Support Vectors
- 2) Hyperplanes
- 3) Marginal Distance
- 4) Linear Separable points
- 5) Non-linear Separable points

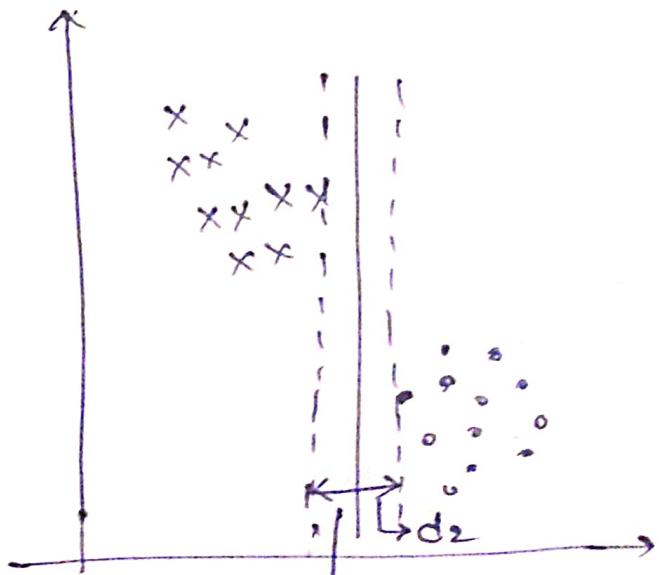
Geometrical Intuition



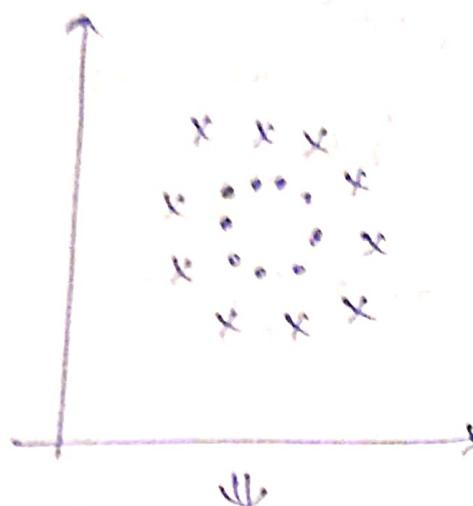
2 lines parallel to hyperplane which passes to a point near to hyperplane.

* Aim is to Maximize Marginal distance. That becomes best hyperplane.

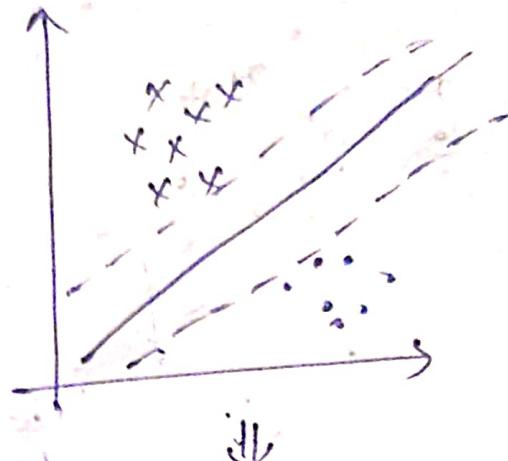
$$d_1 >> d_2$$



* In 2D we say hyperplane which separate the classes maximally. But in 2D hyperplane is straight line.



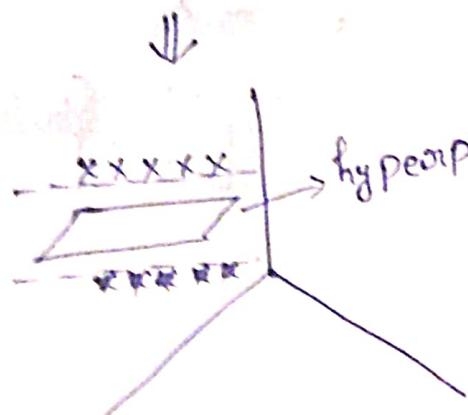
* Not linearly separable because we can't divide this 2 classes by a straight line



* linearly separable because we can divide 2 classes by a straight line

We can solve this type of problem Using SVM Kernel.

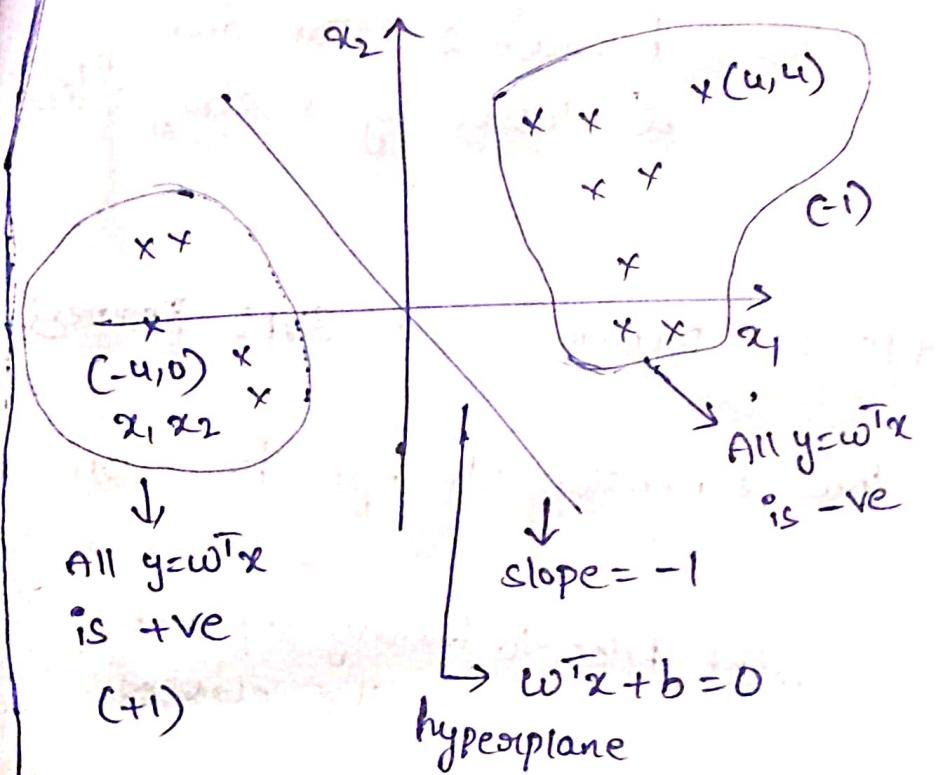
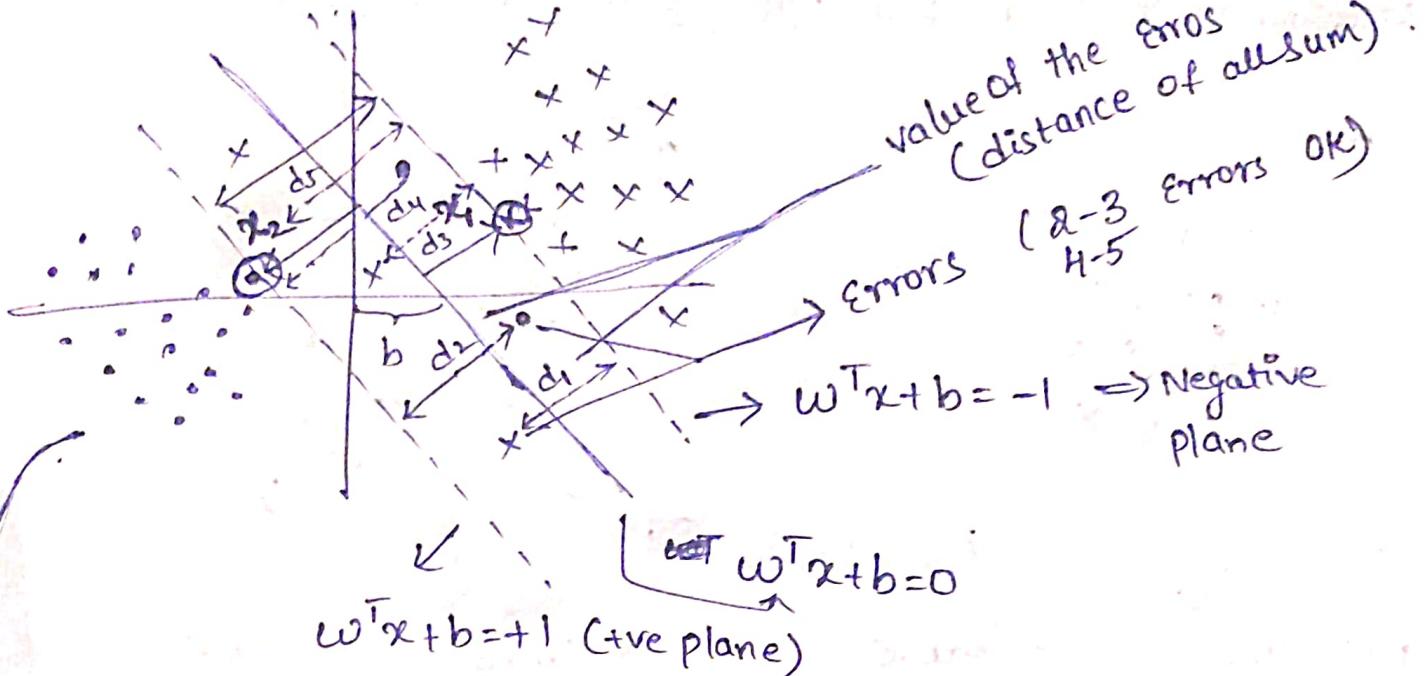
SVM Kernel Converts low Dimension to high dimension.



hyperplane which helps to linearly separate the classes.

* Aim is to identify a hyperplane that have maximum marginal distance which separates classes.

SVM Math Intuition



$$m = -1$$

$$\cos b = 0$$

$$y = w^T x + 0$$

$$= \begin{bmatrix} -1 \\ 0 \end{bmatrix} \begin{bmatrix} -4 & 0 \end{bmatrix} = 4$$

+ve value

$$y = w^T x$$

$$= \begin{bmatrix} -1 \\ 0 \end{bmatrix} \begin{bmatrix} 4 & 4 \end{bmatrix} = -4$$

-Ne

* $x_2 - x_1 = \text{Marginal distance}$.

$$w^T x_1 + b = -1$$

$$w^T x_2 + b = +1$$

$$\frac{w^T x_2 + b - w^T x_1 - b}{2}$$

$$w^T(x_2 - x_1) = 2$$

$$\Rightarrow \frac{w^T}{\|w\|} (x_2 - x_1) = \frac{2}{\|w\|}$$

To remove w^T

* optimization function is $\frac{1}{2} \frac{||w||^2}{||w||} \rightarrow$ we need to Maximize this.

update (\hat{w}, b) max $\frac{1}{2} \frac{||w||^2}{||w||}$

such that
$$y_i \begin{cases} +1 & w^T x_i + b \geq 1 \\ -1 & w^T x_i + b \leq -1 \end{cases}$$

$$y_i * w^T x_i + b \geq 1$$

if it is < 1 it is Misclassification.

* In real world scenario lot of overlapping of points will be there. we can't simply classify data.

$$(\hat{w}, \hat{b}) = \min \frac{1}{2} ||w|| + C \sum_{i=1}^n \xi_i$$

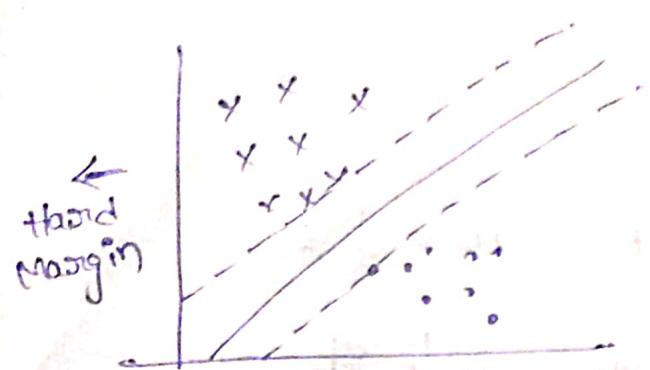
Hinge loss
value of the error

How many
Errors modeled
Can Consider

Regularization.
(we get by hyperparameter
tuning).

* Hard SVM \Rightarrow strictly divides points into different classes.

SVM Kernels Indepth Intuition



① soft Margin (Have Errors of $C=1, 2, 3 \dots$)

② Hard Margin (NO errors)



① Polynomial Kernel

② RBF Kernel

③ Sigmoid Kernel

Not linearly separable

SVM Kernels

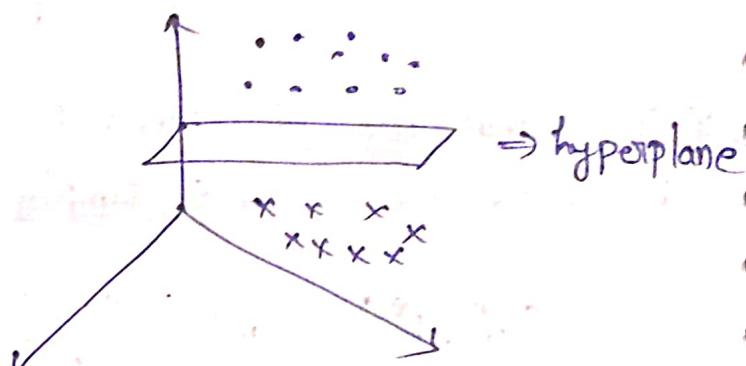
Converts low-to-high dimension

Here 2D to 3D

SVM Kernels does this by some Mathematical Formulas.



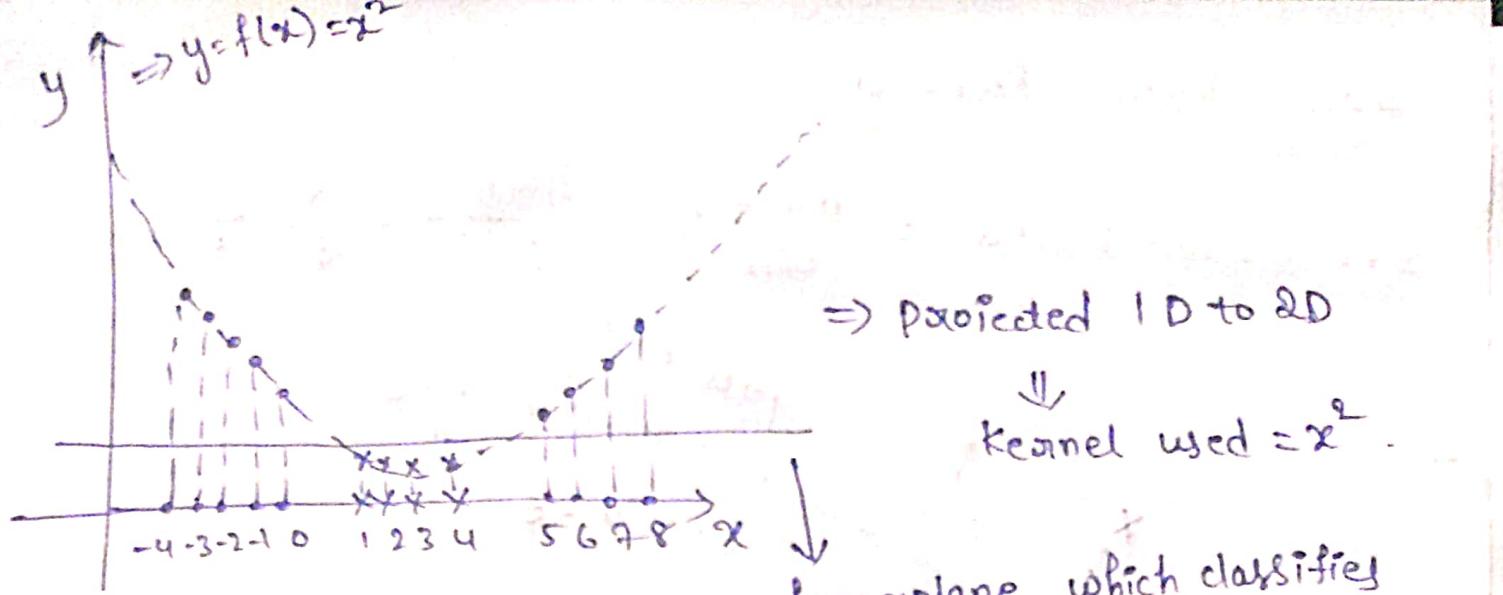
Not Exactly Separable



1d \rightarrow 2d

using some kernel

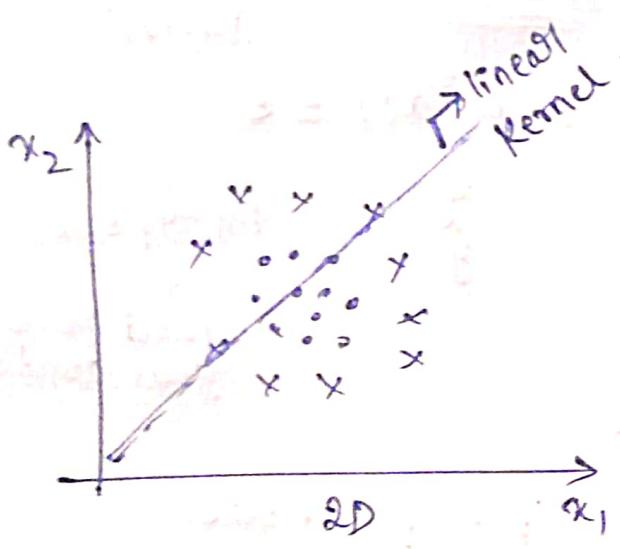
$$y = f(x) \\ = x^2$$



\Rightarrow projected 1D to 2D

\downarrow
Kernel used $= x^2$

hyperplane which classifies
data into classes.



\Rightarrow linear kernel for this 50% accuracy only.

\Rightarrow Now, polynomial Kernel apply

$$p \cdot K = (x_1^T \cdot x_2 + 1)^d \rightarrow \text{dimensions}$$

$$f(x_1, x_2)$$

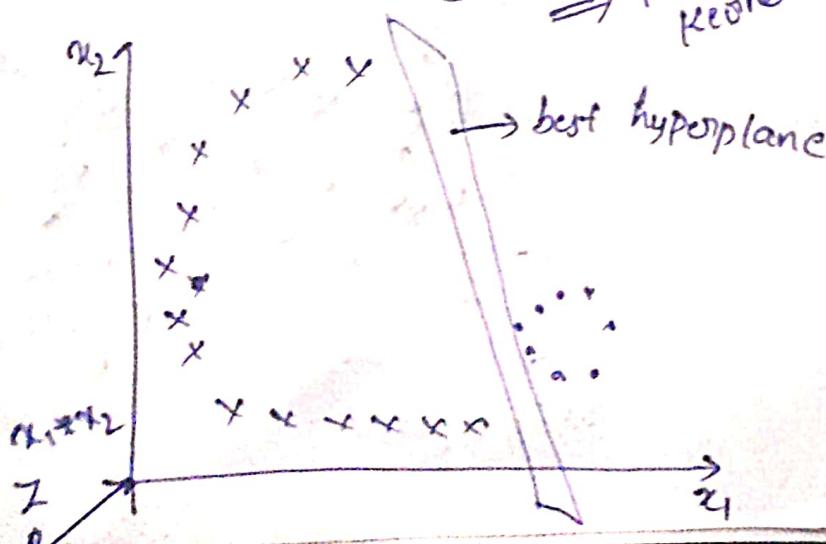
$$x_1 \quad x_2 \quad y_0 \quad y = f(x_1) x_2$$

$$f(x_1, y_2) = (x_1^T \cdot x_2 + 1)^d$$

$$\Rightarrow x_1 \quad x_2 \quad y_0 \quad x_1^2 \quad x_2^2 \quad x_1 x_2$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \cdot \begin{bmatrix} x_1 & x_2 \end{bmatrix}$$

\Downarrow 3D \Rightarrow polynomial kernel



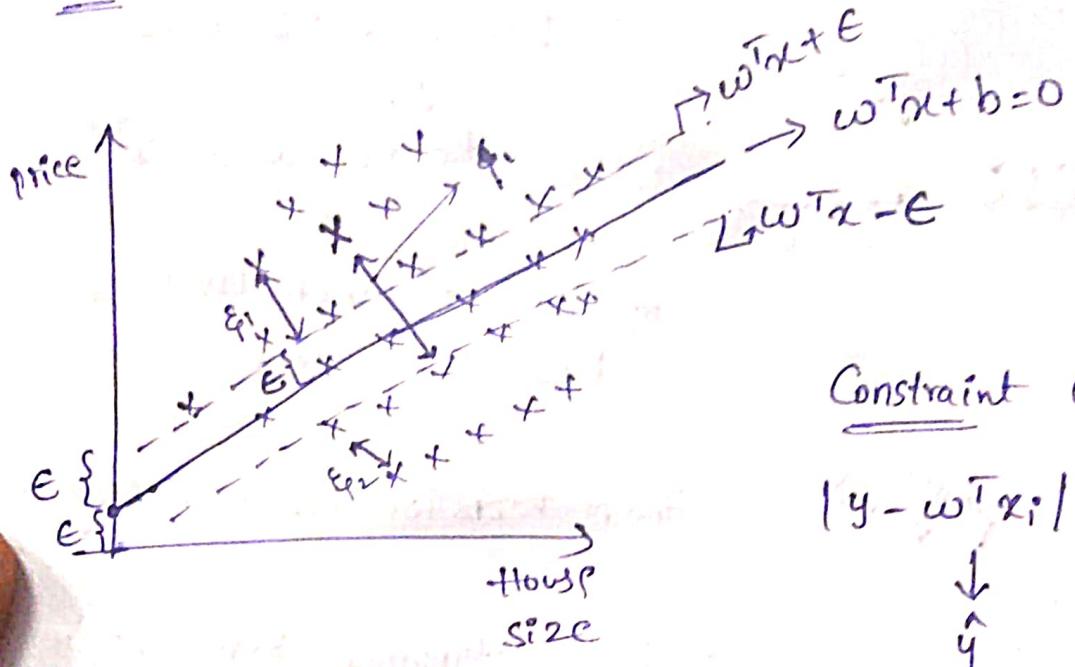
$$= \begin{bmatrix} x_1^2 & x_1 x_2 & x_2^2 \\ x_2 x_1 & \end{bmatrix}$$

\Rightarrow Type of kernel

- 1) Polynomial } determined
 - 2) RBF } by
 - 3) Sigmoid } hyperparameters
- Tuning.

Support Vector Regression

$$\text{SVC} \Rightarrow \text{Cost Function} = \underset{w,b}{\text{minimize}} \quad \frac{\|w\|}{2} + C \sum_{i=1}^n |\xi_i| \rightarrow \text{Error}$$



Model
Constraint (without outside points) -

$$|y - w^T x_i| \leq \epsilon$$

$$\downarrow y$$

$$|y - \hat{y}| \leq \epsilon$$

good point
good prediction

Cost Function - for outside points.

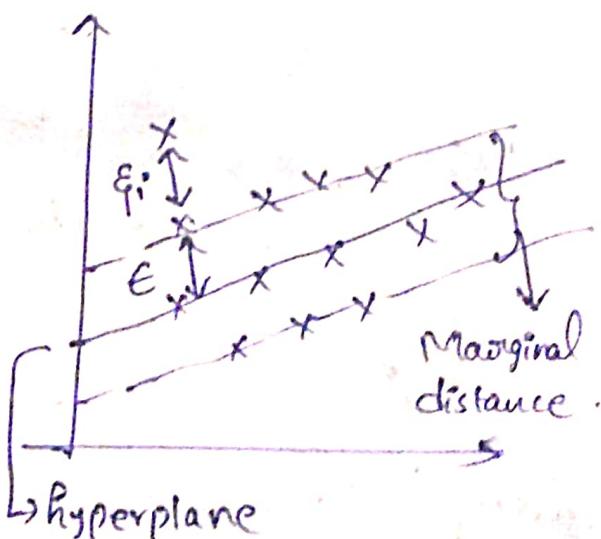
$$\underset{(w,b)}{\text{minimize}} \quad \frac{\|w\|}{2} + \left[C \sum_{i=1}^n |\xi_i| \right] \rightarrow \text{hyperparameter}$$

\downarrow outside points distance

How many points
Can have outside
Marginal plane

Constraint

$$|y_i - w^T x_i| \leq \epsilon + |\xi_i|$$

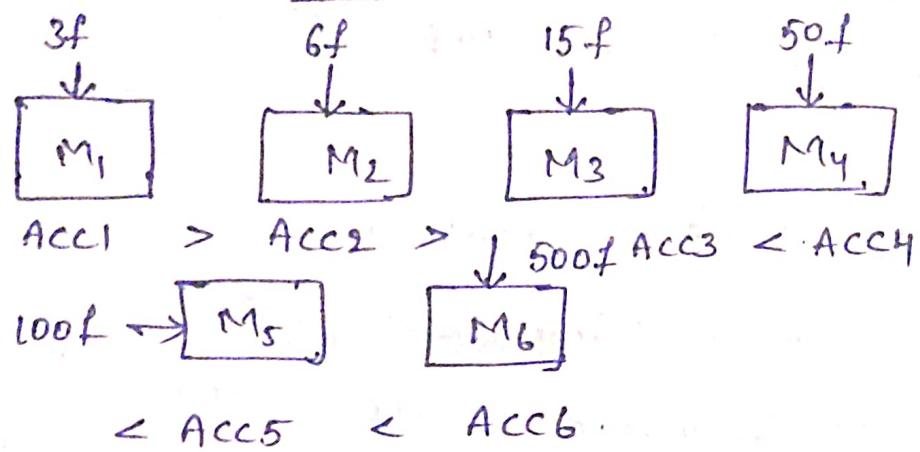


Principal Component Analysis (PCA) [Dimensionality Reduction]

Why should we use PCA?

① Curse of Dimensionality

Dimensions is Features



Dataset :- 500 features

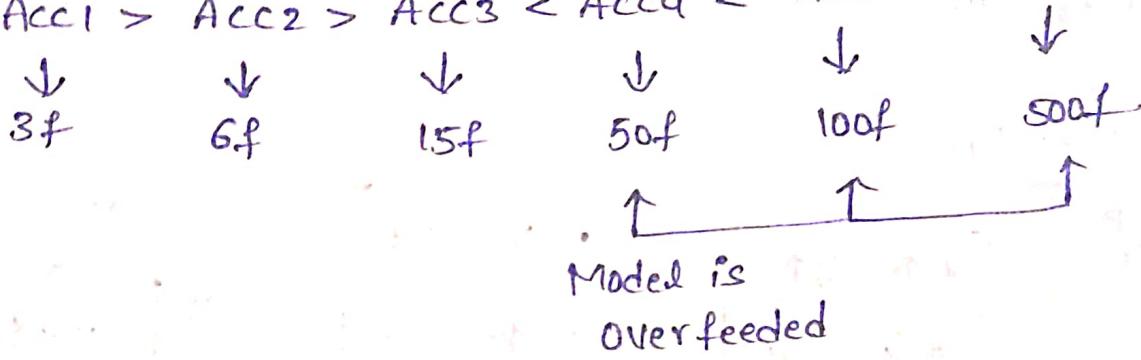
Price of house

House size

No. of bedrooms

No. of bathrooms

* ACC₁ > ACC₂ > ACC₃ < ACC₄ < ACC₅ < ACC₆.



② Model performance Degrade as no. of features increasing

Model performance reduces because many mathematical formulas get behind at one point of time give too confused data.

* Two different ways to remove Curse of Dimensionality.

1) Feature Selection.

⇒ Take only imp features → Train our model.

2) Dimensionality Reduction (PCA) ⇒ Feature Extraction.

⇒ Derive a feature from a set of features.

$f_1, f_2, f_3, f_4 \Rightarrow O/P$

II Feature Extraction

$Df_1, Df_2 \Rightarrow O/P$

- * Feature selection & extraction helps to reduce the no. of features or to extract imp features from older features.

④ why we use dimensionality reduction?

→ prevent → Curse of Dimensionality.

→ Improve the performance of the model. (reduce model time of execution)

→ Visualize the data → understand the data.

Feature Selection

I/P
 x

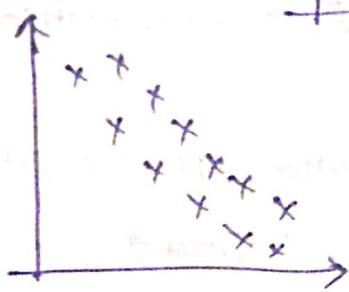
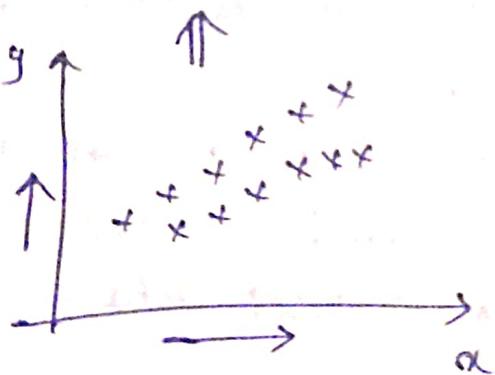
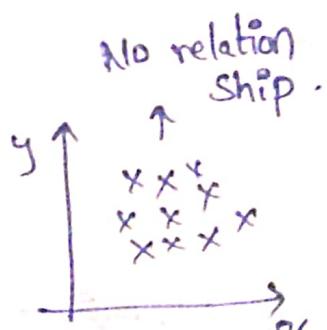
O/P
 y

-

-

$x \uparrow$	$y \uparrow$
$x \downarrow$	$y \downarrow$

$x \downarrow$	$y \uparrow$
$x \uparrow$	$y \downarrow$



$$\text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$n-1 \rightarrow$ for sample data.

+ve → linear Relationship follows $\Rightarrow x \uparrow y \uparrow$

-ve → inverse linear Relationship follows $\Rightarrow x \uparrow y \downarrow$

$x \uparrow y \uparrow$

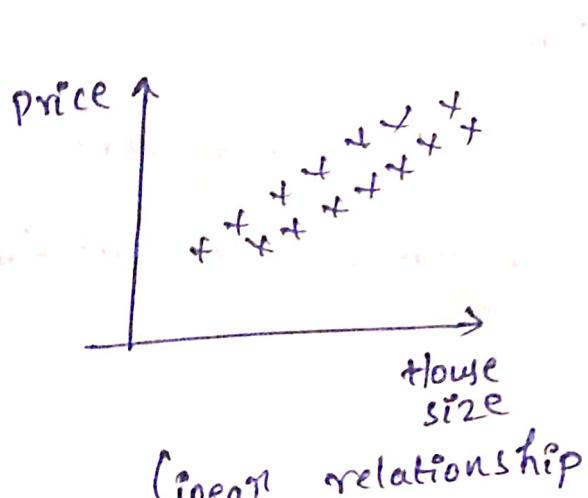
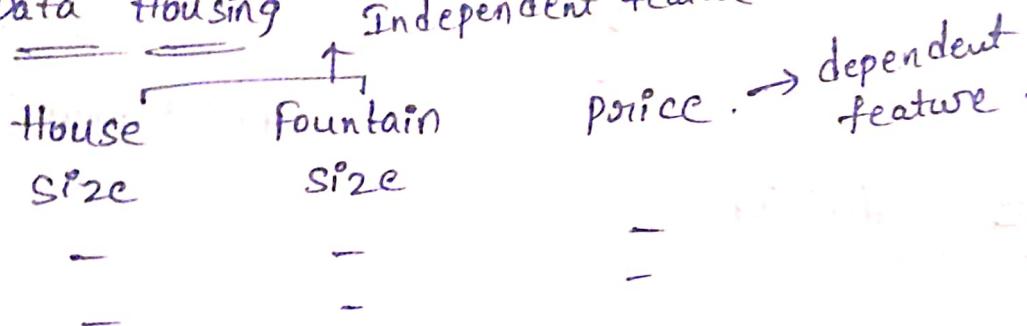
Covariance $\approx 0 \Rightarrow$ No linear relationship b/w x & y .

* Pearson Correlation = $\frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} = -1 \rightarrow 1$

\Rightarrow The more towards the value +1 \Rightarrow The more +ve Correlated & vice versa.

$\Rightarrow \approx 0 \Rightarrow$ No correlation.

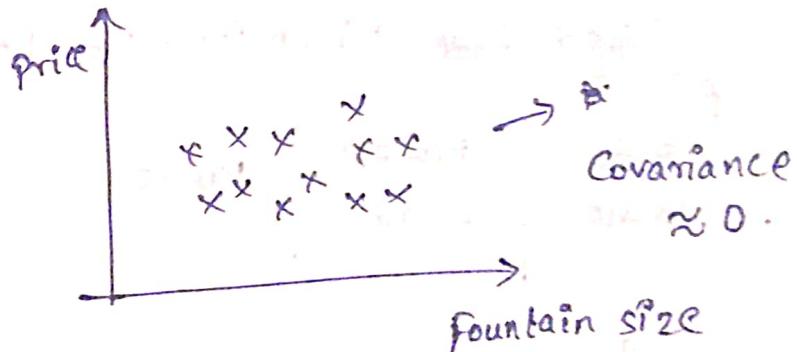
Data Housing Independent feature



Linear relationship

Covariance = +ve

↓ Selected



No relationship.

↓
No need to select
as a feature.

* This is called Feature Selection.

↓
Using Covariance & Correlation
we do this.

Feature Extraction

Room
House size No. of rooms \Rightarrow Price

↓
Apply Transformation to Extract New feature.

↓
House size \Rightarrow price

Dimensionality reduction

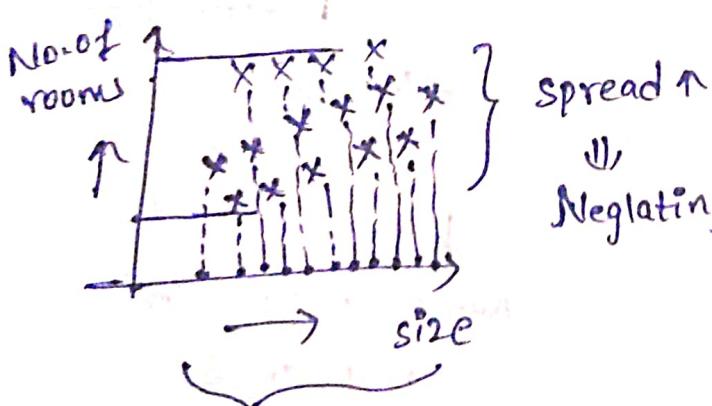
- ↓
- 2 features
- ↓
- 1 feature

* Hence feature selection couldn't helpful bcoz 2 features are super imp to predict house price.

PCA Geometric Intuition

↓
Used for Dimensionality Reduction.

size of house No. of rooms price



spread of data points

spread ↑
variance ↑

PCA Aim Here

2 dimension \rightarrow 1 dimension

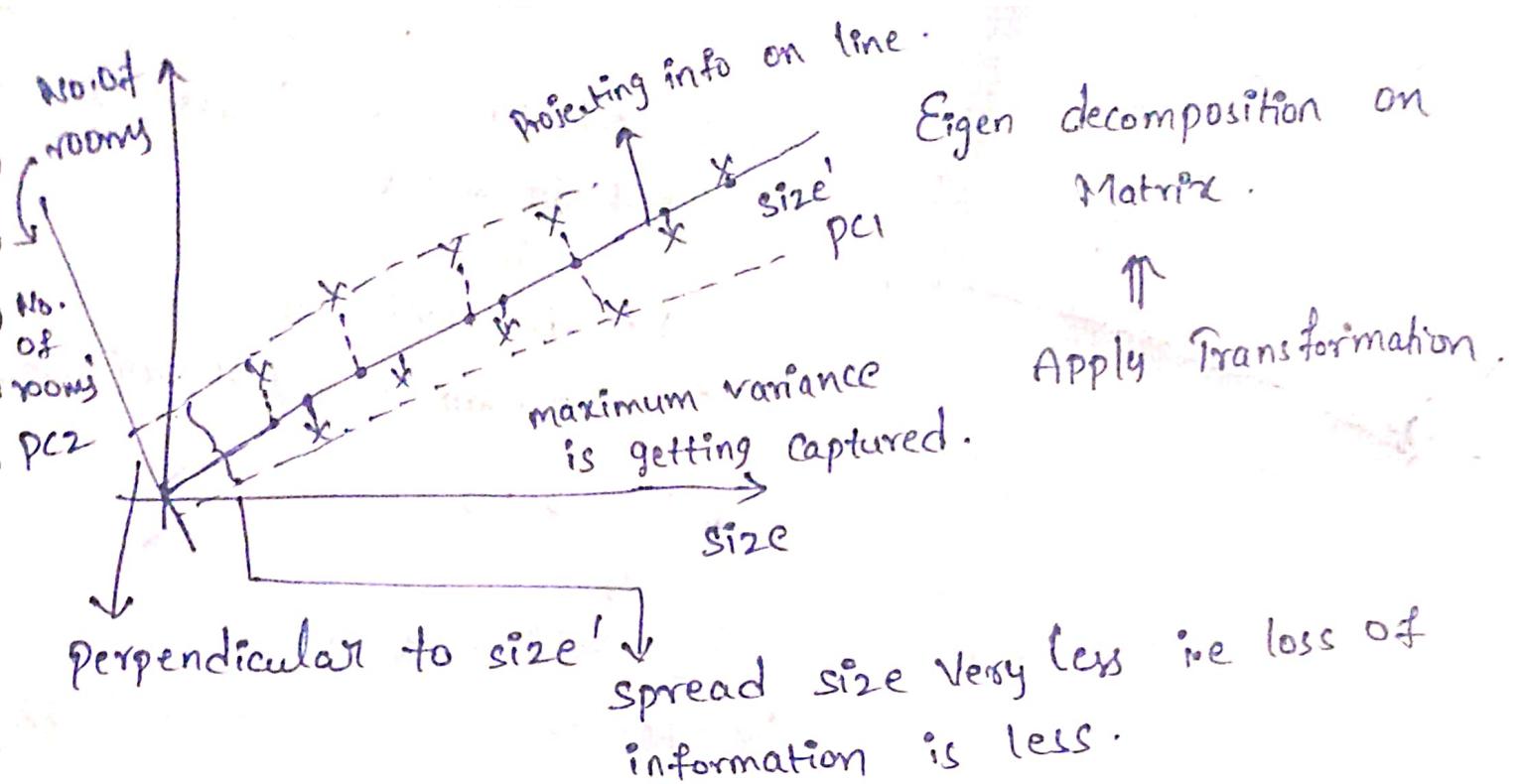
spread ↑
Neglecting this points by projecting in 1D.

2D \rightarrow 1D

* loss of information
(No. of rooms)



feature Extraction
with lot of information lost.



* Here $2D \rightarrow 1D$.



Such that much information is not lost.

- \Rightarrow PC1 \rightarrow captures Maximum Amount of Variance.
- \Rightarrow PC2 \rightarrow captures Second maximum Amount of Variance.

Goal of PCA Algorithm

* To find out the principal Components in such a way that maximum variance will be Captured.

PC1, PC2 \leftarrow 2 dimension

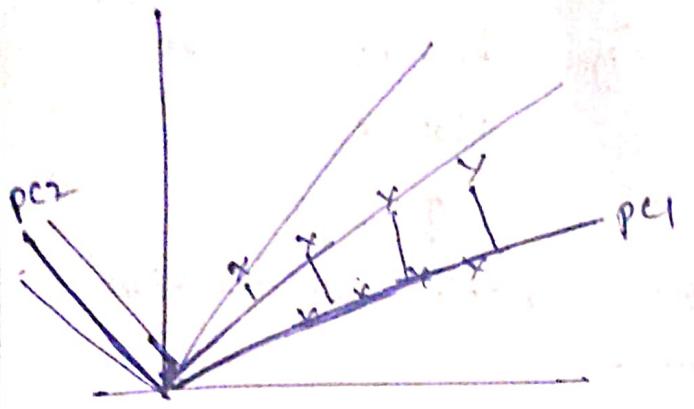
$\text{var}(\text{PC1}) > \text{var}(\text{PC2})$

3 dimension.

\Downarrow
 $\text{var}(\text{PC1}) > \text{var}(\text{PC2}) > \text{var}(\text{PC3})$

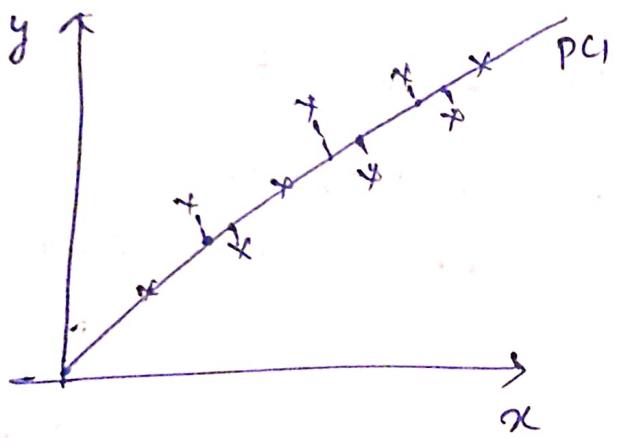
Eigen decomposition on Matrix.

Apply Transformation.



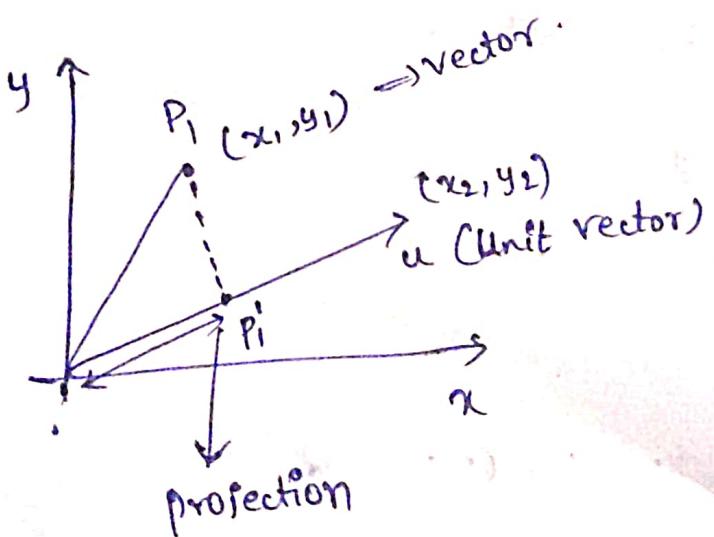
- * 2 Best principal Components need to select
 - ↳ depends which Captures Maximum amount of variance.
- $3D \rightarrow 1D$
- ↳ PC_1, PC_2, PC_3
- $\text{Var}(PC_1) > \text{Var}(PC_2) > \text{Var}(PC_3)$

Math Intuition behind PCA Algorithm.



PCA comes to conclusion about Best PCA based on

- ① projection
- ② cost function.



$$\text{proj}_{P_1} u = \frac{P_1 \cdot u}{\|u\|}$$

$\|u\| = 1 \Rightarrow \text{unit vector.}$

$$\text{prod}_{P_1} u = P_1 \cdot u$$

↳ scalar value.

$$P_0^{'}, P_1^{'}, P_2^{'}, P_3^{'}, \dots, P_n^{'}$$

↓
scalar values.

$$x_0^{'}, x_1^{'}, x_2^{'}, x_3^{'}, x_4^{'}, \dots, x_n^{'}$$

$$\text{Max Variance} = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n} \quad [\text{goal is to find best unit vector which captures maximum variance}]$$

↓
Cost function.

* It's not possible to find Every unit vector for each point & choose best.

↓
Eigen Decomposition helps here.

↓
Eigen vectors And Eigen values

Covariance Matrix between features need to find.

① Covariance Matrix between features need to find.

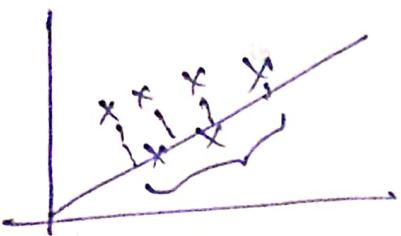
② Eigen vectors and Eigen values will found out from this Covariance Matrix.

③ Eigen vector → Eigen value → highest magnitude of the eigen vector.

↓
Captures the maximum variance.

Eigen vectors & Eigen values [Linear Transformation]

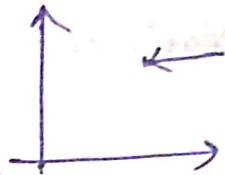
Eigen decomposition of Covariance Matrix



↓
Eigen vector & Eigen values.

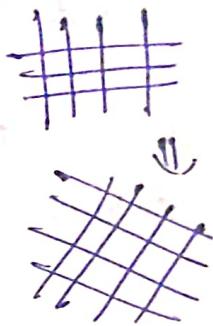
$$[\quad] * [v] = \lambda * v$$

↓
Eigen value



$$A * v = \lambda * v$$

↓
Matrix ↓
linear



Transformator on vector.



Among Eigen vectors \Rightarrow which have maximum Magnitude.



That principal component Captures
maximum variance.

* Eigen Vector \rightarrow with Max magnitude. (Max Eigen Value)



Max Eigen vector

↓
Selected as Best principal component

↓
bcz captures. max variance.

Steps to calculate Eigen value & vectors:

① Covariance of features.

x, y

\downarrow
 x'

Σ

$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{N-1}$$

$$\Rightarrow \text{Cov}(x, x) = \text{Var}(x)$$

x	y
x	$\text{Var}(x)$
y	$\text{Cov}(x, y)$

$A =$

x	y
x	$\text{Var}(x)$
y	$\text{Cov}(x, y)$

2×2

x	y	z	
x	$\text{Var}(x)$	$\text{Cov}(x, y)$	$\text{Cov}(x, z)$
y	$\text{Cov}(y, x)$	$\text{Var}(y)$	$\text{Cov}(y, z)$
z	$\text{Cov}(z, x)$	$\text{Cov}(z, y)$	$\text{Var}(z)$

3×3

② Eigen values.

$$A \cdot \varphi = \lambda \cdot \varphi$$

\Downarrow

$\lambda_1, \lambda_2 \Rightarrow$ Eigen values.

f_1, f_2

\Downarrow

λ_1

\Downarrow

PC_1

\Downarrow

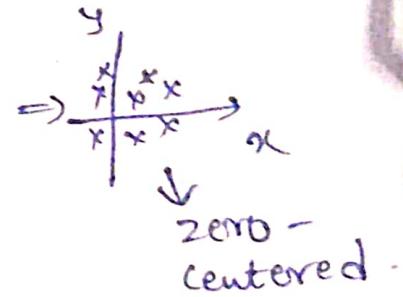
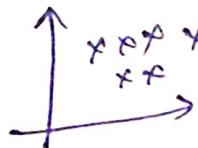
λ_2

\Downarrow

PC_2

\rightarrow Eigen values.

① Standardize data. \Rightarrow means



② Covariance of x & y matrix

x	y
x	$\text{Var}(x)$
y	$\text{Cov}(x, y)$

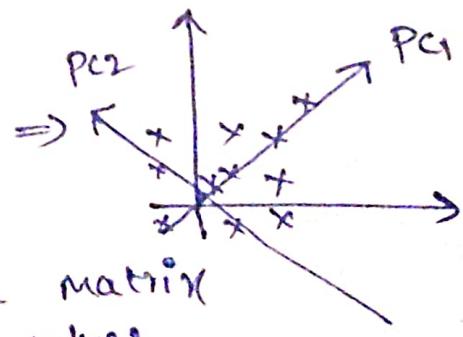
2×2

③ Find out Eigen vectors & values.

$$A \cdot \varphi = \lambda \varphi$$

\Downarrow

$\lambda_1, \lambda_2 \Rightarrow$ for 2×2 matrix
Eigen values.



$\Rightarrow 3D \rightarrow 1D$

$$\begin{matrix} \lambda_1, \lambda_2, \lambda_3 \\ \downarrow \quad \downarrow \quad \downarrow \\ PC_1 \quad PC_2 \quad PC_3 \end{matrix}$$

$$\begin{matrix} \checkmark & | \\ ID & 1 \\ ID \rightarrow 2D \end{matrix}$$

$3D \rightarrow 1D$

$$\begin{matrix} \lambda_1, \lambda_2, \lambda_3 \\ \backslash \quad / \\ PC_1 \quad PC_2 \end{matrix}$$

ID

$2D \rightarrow 1D$

$$\begin{matrix} \lambda_1, \lambda_2 \\ \downarrow \quad \downarrow \\ PC_1 \quad PC_2 \end{matrix}$$

$$\begin{matrix} \checkmark \\ ID \end{matrix}$$