

Eg:-

x_0	x_1	x_2	y
1	1	2	3
2	1	1	4
3	3	3	5

- ① Compute prediction
- ② Compute gradient
- ③ Update Parameters

$$h_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2$$

$$\theta_0 = \theta_1 = \theta_2 = 0, \alpha = 0.1$$

$$J(\theta) > 25$$

→ First iteration

$$\textcircled{1} \quad h_\theta(x^0) = \theta_0 x_0^0 + \theta_1 x_1^0 + \theta_2 x_2^0 = 0 \times 1 + 0 \times 1 + 0 \times 2 = 0$$

$$h_\theta(x^1) = \theta_0 x_0^1 + " " = 0 \times 1 + 0 \times 2 + 0 \times 1 = 0$$

$$h_\theta(x^2) = \theta_0 x_0^2 + " " = 0 \times 1 + 0 \times 3 + 0 \times 3 = 0$$

\textcircled{2} Compute gradients

i) Compute errors

$$\begin{array}{c|c|c} e^{(0)} = h_\theta(x^0) - y & e^{(1)} = h_\theta(x^1) - y & e^{(2)} = h_\theta(x^2) - y \\ \hline 0 - 3 & 0 - 4 & 0 - 5 \\ -3 & -4 & -5 \end{array}$$

ii) Calculate Gradient

$$\frac{\partial J}{\partial \theta_0} = \sum_{i=0}^2 (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$= -3 \times 1 + -4 \times 1 + -5 \times 1$$

$$= -12$$

$$\frac{\partial J}{\partial \theta_0} = \sum_{i=1}^3 (\hat{y}_0(x_i) - y^{(i)}) x_0^{(i)}$$

= $-3 \times 1 + -4 \times 2 + -5 \times 3$

~~(Note) terms cancel out~~

~~$\therefore (0.02)^2 = -26$~~

$$\frac{\partial J}{\partial \theta_1} = \sum_{i=1}^3 (\hat{y}_1(x_i) - y^{(i)}) x_1^{(i)}$$

= $-3 \times 2 + -4 \times 1 + -5 \times 3$

~~(Note) terms cancel out~~

~~$\therefore -25 =$~~

③ Update Parameters

$$\theta_0 = \theta_0 - \alpha \frac{\partial J}{\partial \theta_0} = 0 - 0.1(-12) = 1.2$$

$$\theta_1 = \theta_1 - \alpha \frac{\partial J}{\partial \theta_1} = 0 - 0.1(-26) = 2.6$$

$$\theta_2 = \theta_2 - \alpha \frac{\partial J}{\partial \theta_2} = 0 - 0.1(-25) = 2.5$$

$$\theta_0 = 1.2 \quad \theta_1 = 2.6 \quad \theta_2 = 2.5$$

$$J = 94.95 \quad * \text{This may be because we took large steps}$$

Iteration - 2 :-

$$\theta_0 = -1.02, \theta_1 = -2.41, \theta_2 = -2.6$$

$$J = 25,678.56$$

$$\text{If: } \theta_0 = \theta_1 = \theta_2 = 0, \alpha = 0.01 \quad J = 10.81$$

After Iteration 1

$$\theta_0 = 0.19, \theta_1 = 0.43, \theta_2 = 0.41 \quad J = 5.41$$

After Iteration 2

$$\theta_0 = 0.25, \theta_1 = 0.55, \theta_2 = 0.53 \quad J = 3$$

① Defining the Distribution

$$L(\theta) = P(x_{1:n} | \theta)$$

$$= \prod_{i=1}^n P(x_i | \theta) \quad \# \text{ since } I.i.D \text{ is assumed}$$

② Computing the log-likelihood :-

$$l(\theta) = \log L(\theta)$$

$$= \log \prod_{i=1}^n P(x_i | \theta)$$

$$= \log [\theta^m (1-\theta)^{n-m}] \quad [\log(ab) = \log(a) + \log(b)]$$

$$= m \log \theta + (n-m) \log (1-\theta)$$

③ Differentiate and set to '0' to find estimate of θ :-

$$l(\theta) = m \log \theta + (n-m) \log (1-\theta)$$

$$\frac{d(l(\theta))}{d\theta} = \frac{d}{d\theta} [m \log \theta + (n-m) \log (1-\theta)]$$

$$= \frac{m}{\theta} + (n-m)(-1) \frac{1}{1-\theta} = \frac{1}{\theta(1-\theta)} \text{ pol } n -$$

$$= \frac{m(1-\theta) + (m-n)\theta}{\theta(1-\theta)} \quad \text{programme} \quad \# (103) \text{ pol programme}$$

\Rightarrow Equate to zero

$$\theta \cdot \theta(1-\theta) = m(1-\theta) + (m-n)\theta$$

$$\theta = m[1-\theta] + (m-n)\theta$$

$$= m - m\theta + m\theta - n\theta$$

$$\cancel{\text{cancel}} \quad \theta = \frac{m}{n}$$

$$\theta = \frac{m}{n}$$

⇒ One-Hot Encoding :- [part of EDA]

Colour = {Red, Blue, Green}

$$\begin{matrix} \cdot & \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} & \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} & \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \\ & 3 \times 1 & 3 \times 1 & 3 \times 1 \end{matrix}$$

- * Vector of size of the set is created and each bit is assigned for each variable respectively
- * In EDA we will determine that colour feature must be converted using One-hot Encoding.
- * But in target variable we do not use One-hot encoding
Instead ; $Y = 1^T N^{20}$
- * ∵ When we assign a vector for a categorical data it is like creating multiple columns within the column

$$\text{eg:- } 0_0 x_0 + 0_1 x_1 + 0_2 x_2 + \dots + 0_d x_d =$$

- * Disadv :- The number of features increases

green	Blue	Red
1	0	0
\rightarrow for green colour		

x_1	x_2	y
Red	35.2	Y
blue	38.0	N



x_{1-0}	x_{1-1}	x_2	y
0	1	35.2	Y
1	0	38	N

$$x_0 x_1^1 = \hat{y} = 0_0 x_0 + 0_1 x_{1-0} + 0_2 x_{1-1}$$

$$\Rightarrow + 0_3 x_2$$

$$= 0_0 x_0 + 0_1 x_{1-0} + 0_2 x_{1-1} + 0_3 x_2$$

Eg: k -fold cross validation

1
2
3
4
5

x_1	x_2	y
2	-1	1
0.5	1.2	0
1	2	1
-3	-2	1
4	0.1	0

$$F1 \{O_1, O_2\} = \{1.8, 2.8\}$$

$$F2 \{O_1, O_2\} = \{2.1, 3.1\}$$

$$F3 \{O_1, O_2\} = \{1.9, 4\}$$

Perform 3 fold CV

split must always remain same
so do not change the folds

\Rightarrow Fold 1 :

Training Set

x_1	x_2	y
2	-1	1
0.5	1.2	0
1	2	1
-3	-2	1

Test set

x_1	x_2	y	$g(z)$
4	0.1	0	0

Accuracy

$$\theta^T x^1 = 1.8 \times 4 + 2.8 \times 0.1 = -6.9$$

$$g(z) = \frac{1}{(1 + e^{-z})}$$

$$= 0.001$$

$$0.001 < 0.5$$

\therefore Accuracy = 100%

\Rightarrow Fold 2 :

Training set

x_1	x_2	y
1	2	1
-3	-2	1
4	0.1	0

Test set

x_1	x_2	y	$g(z)$
2	-1	1	2
0.5	1.2	0	1

$$\theta^T x^1 = 2.1 \times 2 + 3.1 \times -1 = 1.1$$

$$\theta^T x^2 = 2.1 \times 0.5 + 3.1 \times 1.2 = 4.77$$

$$g(z)^1 = 0.75 \quad g(z)^2 = 0.9912$$

\therefore Accuracy = 50%

\Rightarrow Fold 3 :

Training set

x_1	x_2	y
2	-1	1
0.5	1.2	0
4	0.1	0

Test set

x_1	x_2	y	$g(z)$
1	2	1	1
-3	-2	1	0

$$\theta^T x^3 = 1 \times 1.9 + 2 \times 4 = 9.9$$

$$\theta^T x^4 = -3 \times 1.9 + -2 \times 4 = -13.7$$

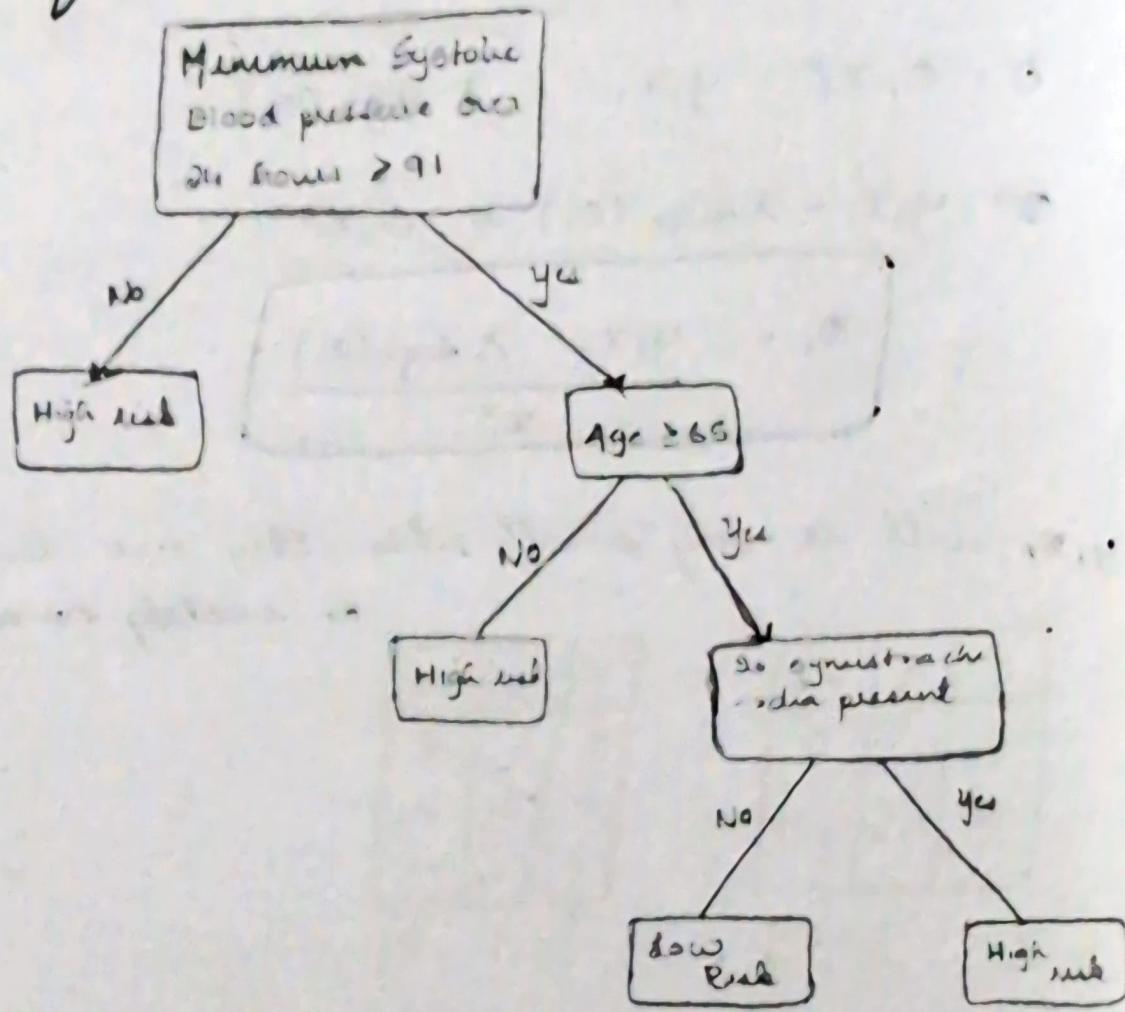
$$g(z)^3 = 0.999$$

$$g(z)^4 = 0.000002$$

\therefore Accuracy = 50%

⇒ Decision Tree Algorithm :-

1) Classification Tree :-



- * Here, the tree can be considered as the hypothesis function.
- * Internal nodes can be initialised at random.
- * For every internal node the parameters are feature & threshold
 - split dimension
 - split value

Eg:-
Features :-

x_1 : date

x_2 : age

x_3 : height

x_4 : weight

x_5 : Since Tachycardia

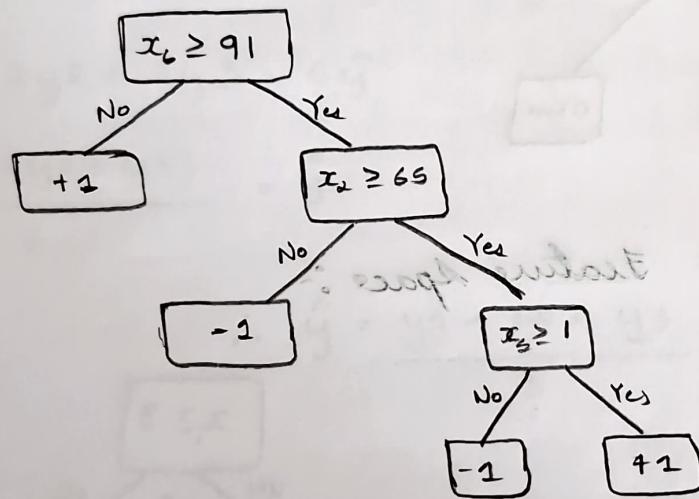
x_6 : minimum systolic bp 24 hrs

x_7 : latent systolic bp.

Labels :-

high risk :- +1

low risk :- -1



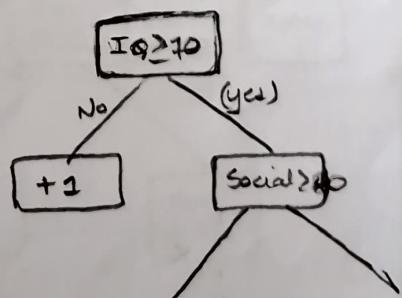
Eg
check for this data :-

$$x^{(i)} = \{2020/11/17, 49, \dots\}$$

Eg:-

Features

IQ	Social	Verbal	Risk for autism
100	79	86	Low
30	40	70	High
80	85	65	Low
90	82	45	High



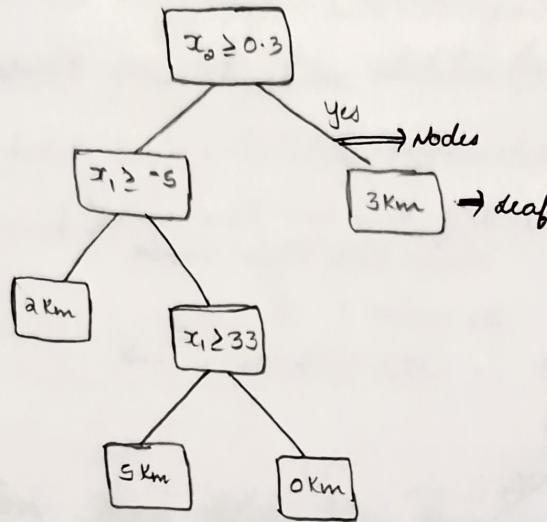
\Rightarrow Regression Tree :-

Eg:-

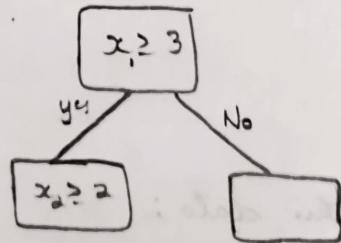
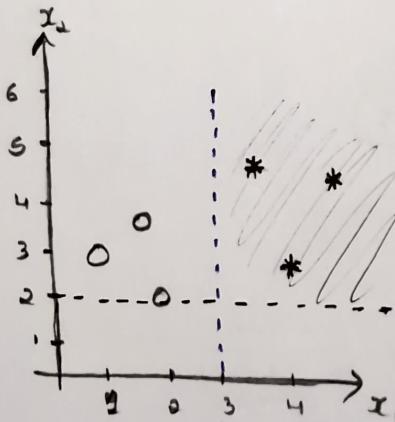
x_1 = Temperature

$y = 1 \text{ Km}, 2 \text{ Km}, 3 \text{ Km}$

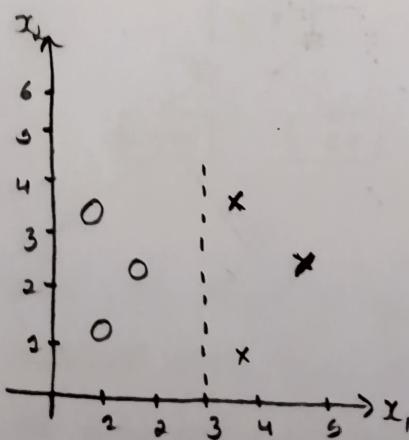
x_2 = precipitation



\Rightarrow Partitioning the Feature Space :-



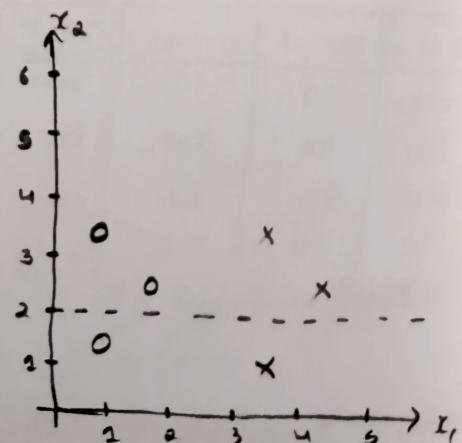
* every node will allow us to split the space



$x_1 >= 3$

Q) Which split is better?

Sol: $x_1 >= 3$



$x_1 >= 2$

$$y = \begin{bmatrix} -2 \\ 3 \\ 4 \end{bmatrix} \begin{matrix} \rightarrow y_1 \\ \rightarrow y_2 \\ \rightarrow y_3 \end{matrix}$$

$\hat{y} = ?$

$$f = (\hat{y} - y_1)^2 + (\hat{y} - y_2)^2 + (\hat{y} - y_3)^2$$

(min)

To find min, we must compute derivative and set it to 0

$$\frac{\partial f}{\partial \hat{y}} = 2(\hat{y} - y_1) \cdot 1 + 2(\hat{y} - y_2) \cdot 1 + 2(\hat{y} - y_3) \cdot 1$$

$$0 = 2(\hat{y} - y_1) + 2(\hat{y} - y_2) + 2(\hat{y} - y_3)$$

$$\Rightarrow 2\hat{y} - 2y_1 + 2\hat{y} - 2y_2 + 2\hat{y} - 2y_3$$

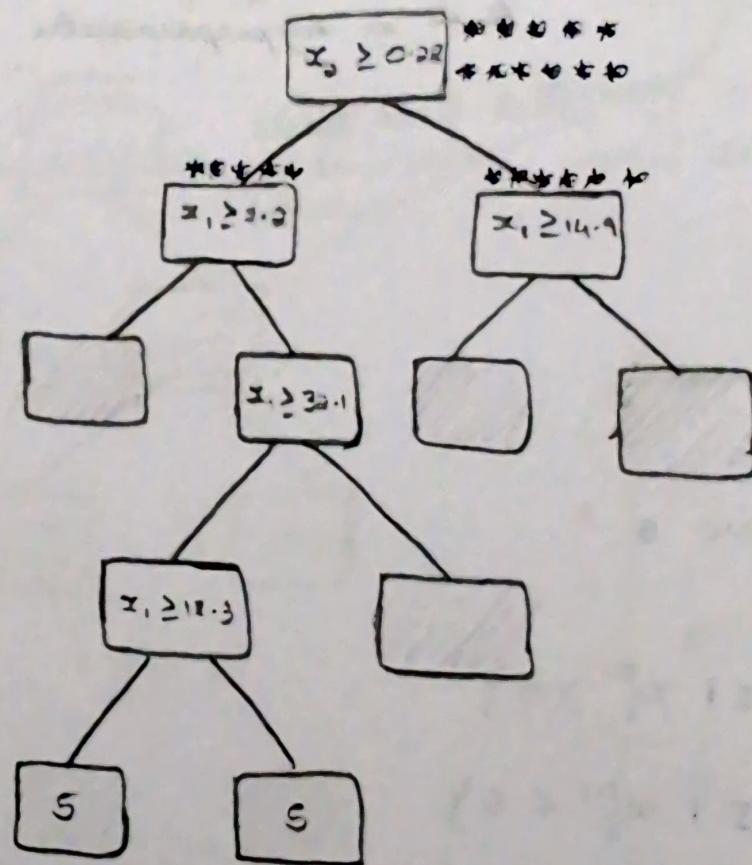
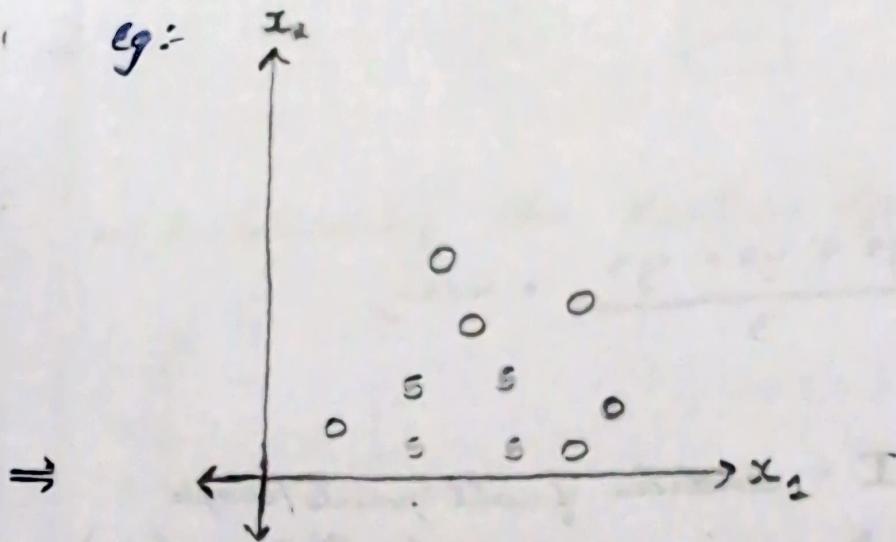
$$\Rightarrow 6\hat{y} - 2y_1 - 2y_2 - 2y_3$$

$$2y_1 + 2y_2 + 2y_3 = 6\hat{y}$$

$$\frac{2(y_1 + y_2 + y_3)}{6} = \hat{y}$$

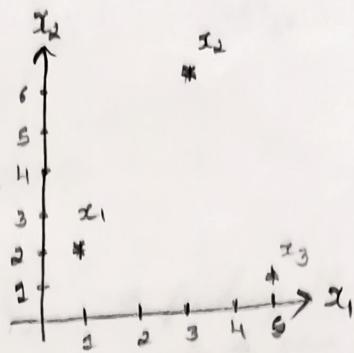
$$\therefore \hat{y} = \frac{y_1 + y_2 + y_3}{3} = \underline{\text{mean}}$$

e.g.:



⇒ Example Question (Decision Tree)

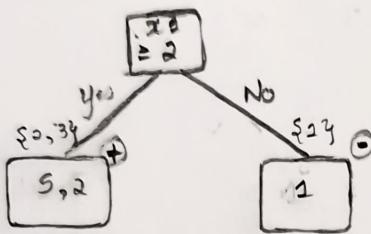
x_1	x_2	y
1	2	1
3	6	5
5	2	2



for $x_1, S = \{2, 4\}$
 $x_2, S = \{5\}$

$$I = \{1, 2, 3\}$$

$$k = 2$$



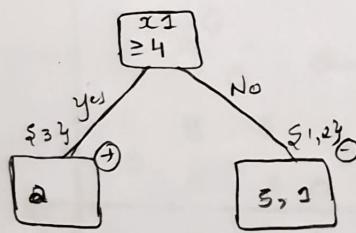
$$\hat{y}^+ = \frac{5+2}{2} = \frac{7}{2} = 3.5$$

$$\hat{y}^- = \frac{1}{2} = 0.5$$

$$E^+ = (5-3.5)^2 + (2-3.5)^2 = 4.5$$

$$E^- = (1-0.5)^2 = 0$$

$$E = 4.5$$



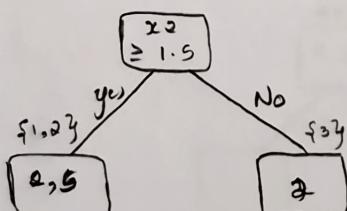
$$\hat{y}^+ = \frac{2}{2} = 1$$

$$\hat{y}^- = \frac{5+1}{2} = \frac{6}{2} = 3$$

$$E^+ = (2-1)^2 = 0$$

$$E^- = (3-3)^2 + (1-3)^2 = 4+4 = 8$$

$$E = 8 - 0 = 8$$



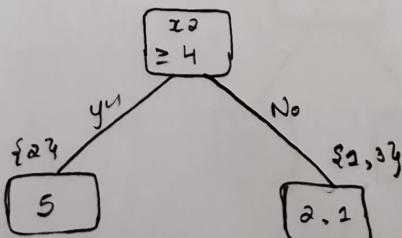
$$\hat{y}^+ = \frac{1+5}{2} = \frac{6}{2} = 3$$

$$\hat{y}^- = \frac{2}{2} = 1$$

$$E^+ = (2-3)^2 + (5-3)^2 = 4+4 = 8$$

$$E^- = 2-1 = 1$$

$$E = 8 - 1 = 7$$



$$\hat{y}^+ = \frac{5}{2} = 2.5$$

$$\hat{y}^- = \frac{2+1}{2} = \frac{3}{2}$$

$$E^+ = (5-5)^2 = 0$$

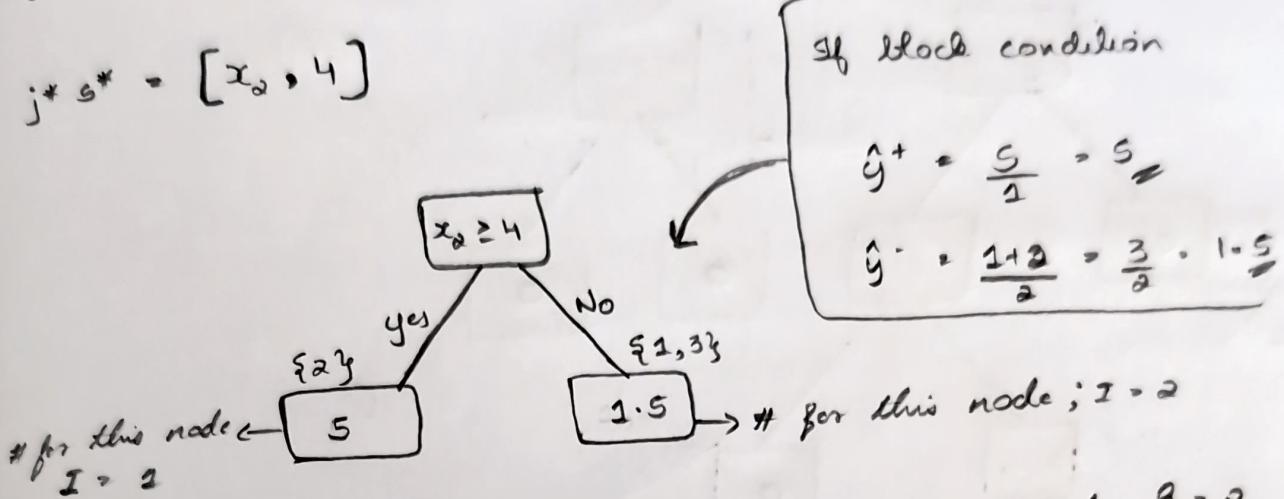
$$E^- = \left(2-\frac{3}{2}\right)^2 + \left(1-\frac{3}{2}\right)^2 = \frac{1}{4} + \frac{1}{4}$$

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

$$j^* \text{ } S^* = \min \{ E_{(x_1, 0)}, E_{(x_1, 1)}, E_{(x_0, 0)}, E_{(x_0, 1)} \} \Rightarrow \arg \min (\text{expand}_d)$$

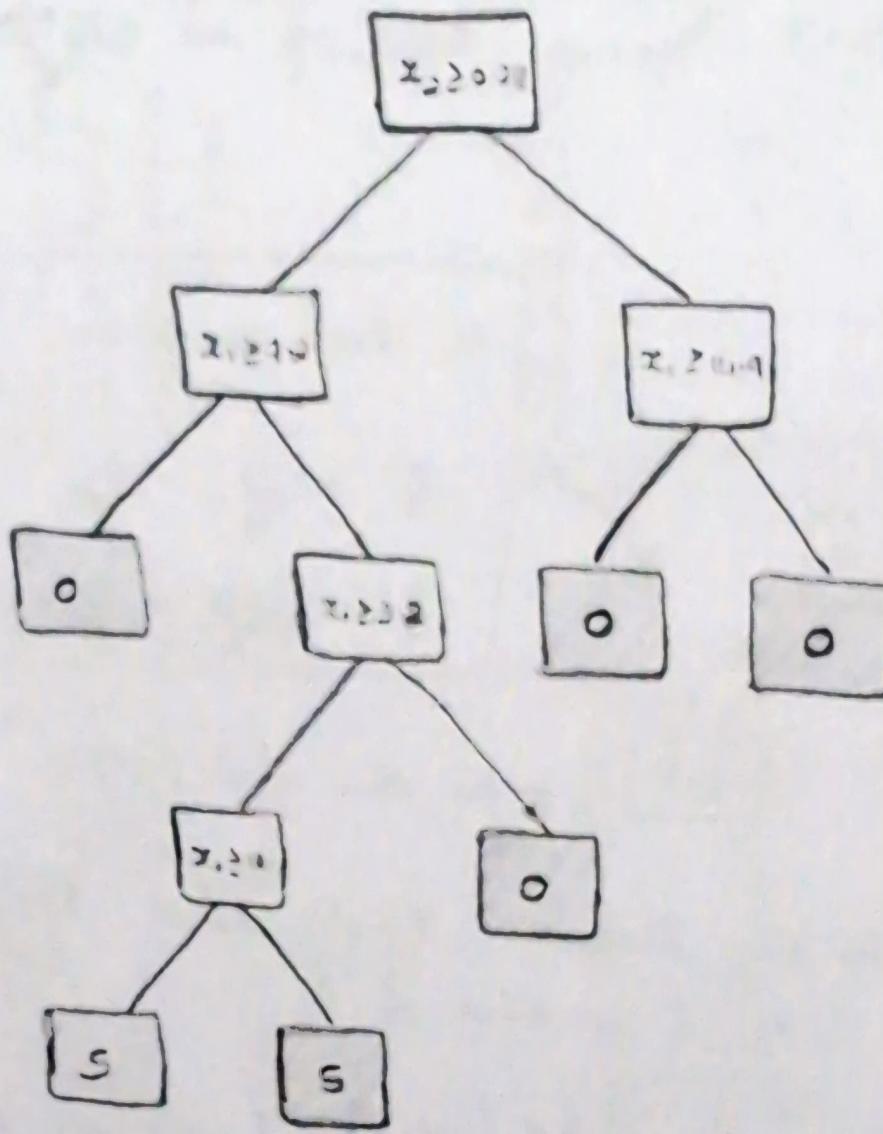
$$j^* \text{ } S^* = x_{(0,1)}$$

$$j^* \text{ } S^* = [x_0, 1]$$

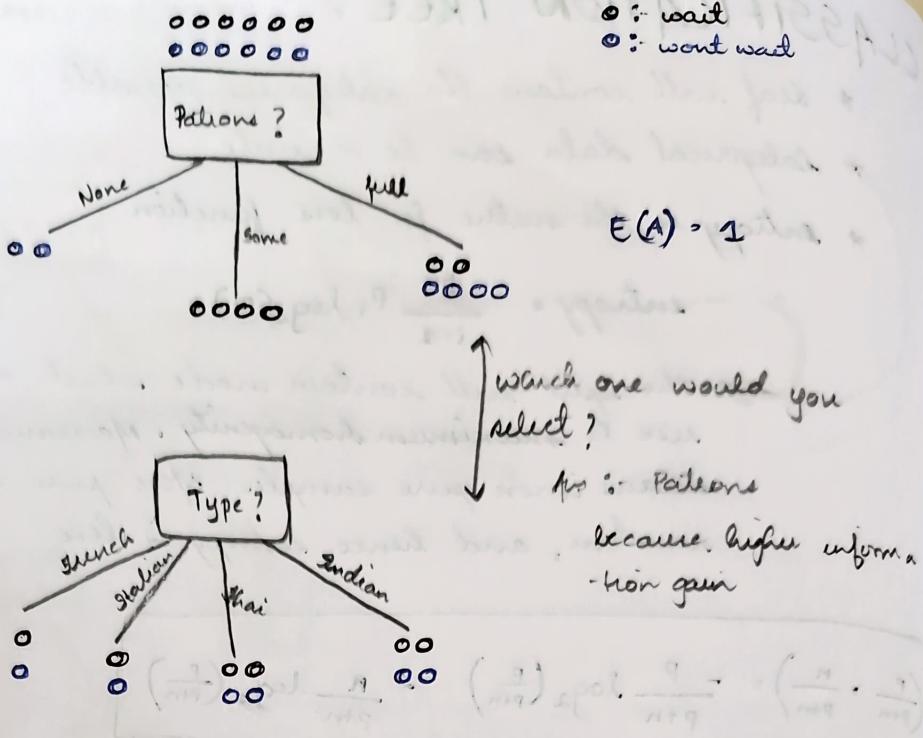


stop here because $I = 2$ and $\hat{g} = 1.5$
ie I is now $\leq k$

- * When we invoke build-tree for right node [No'], we see that $I = 2$ which satisfies the if condition \therefore doesn't enter else loop.
- * When we invoke build-tree for left node [Yes], we see that $I = 1$, which satisfies if condition, \therefore doesn't enter else loop.
- * The leaf contains predicted \hat{y} values for the corresponding data-set ie \hat{y} predicted for 1st and 3rd dataset is $\frac{1+5}{2} = 3$ and 2 dataset is $\frac{5}{2} = 2.5$



- * When the tree is potentially large and/or has very high number of leaves there is very high chance that there can be overfitting of the model.
- * Hence we do regularization to avoid this where we add a penalty term.



$$EH(A) = \sum_{i=1}^k - \left(\frac{p_i + n_i}{p+n} \right) H \left(\frac{p_i}{p_i + n_i}, \frac{n_i}{p_i + n_i} \right) \quad p := 0 \\ n := 0$$

$$I(A) = H \left(\frac{p}{p+n}, \frac{n}{p+n} \right) - EH(A) \quad \left. \begin{array}{l} H(\text{patrons}) - EH(\text{patrons}) \\ \text{Information gain} \end{array} \right.$$

* Pick value of entropy which gives highest information gain

$$(parent)H(\text{patrons}) = H \left(\frac{p}{p+n}, \frac{n}{p+n} \right)$$

$$= H \left(\frac{6}{12}, \frac{6}{12} \right)$$

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

$$(all \ 3 \ children)EH(\text{patrons}) = \sum_{i=1}^k - \frac{p_i + n_i}{p+n} H \left(\frac{p_i}{p_i + n_i}, \frac{n_i}{p_i + n_i} \right)$$

$$= \frac{p_1 + n_1}{12} H \left(\frac{0}{2}, \frac{2}{2} \right) + \frac{p_2 + n_2}{12} \left[H \left(\frac{4}{4}, \frac{0}{4} \right) \right]$$

$$+ \frac{p_3 + n_3}{12} H \left(\frac{2}{6}, \frac{4}{6} \right)$$

$$H\left(\frac{0}{2}, \frac{2}{2}\right) = \frac{-P}{P+n} \log_2 \left(\frac{P}{P+n}\right) - \frac{n}{P+n} \log_2 \left(\frac{n}{P+n}\right)$$

$$= \frac{-0}{2} \log_2 \left(\frac{0}{2}\right) - \frac{2}{2} \log_2 \left(\frac{2}{2}\right)$$

$$= 0 \log_2 (0) - 1 \log_2 (1)$$

$$EH(\text{patients}) = \frac{2}{12} (-0 \log_2 (0) - 1 \log_2 (1)) +$$

$$\frac{4}{12} (-1 \log_2 (1) - 0 \log_2 (0)) +$$

$$\frac{6}{12} \left(-\frac{2}{6} \log_2 \left(\frac{2}{6}\right) - \frac{4}{6} \log_2 \left(\frac{4}{6}\right) \right)$$

$$= \frac{1}{6} (0 - 0) + \frac{1}{3} (0 - 0) + \frac{6}{12} \left(-\frac{2}{6} (-0.585) - \frac{4}{6} \right)$$

(-0.585)

$E[H]$

~~$0.00 + 0.00 + 0.585 = 0.585$~~
 $= 0.455$

~~$0.00 + 0.585 = 0.585$~~
 $\therefore I(A) = 1 - 0.455$

$I(A) = 0.55$

$I(\text{patients}) = 0.55$

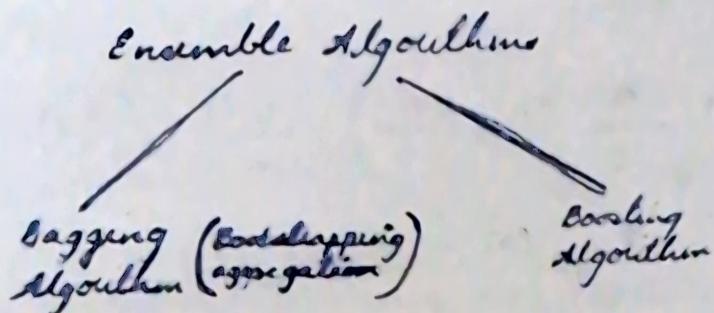
using similar steps

Gini index

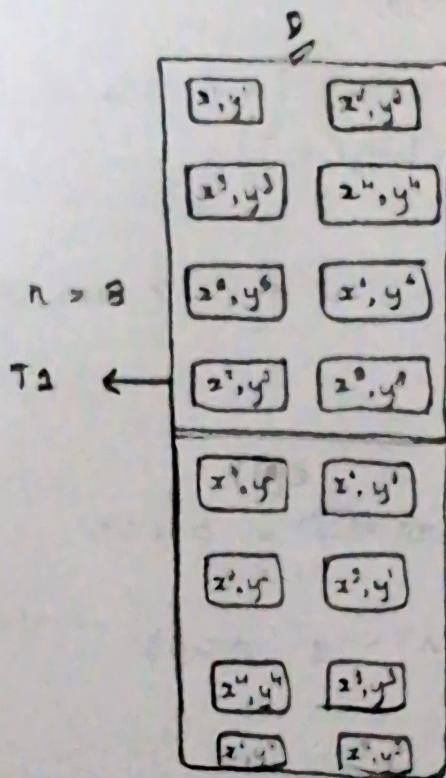
$I(\text{types}) = 0.00$

∴ we select patients as it have more information gain.

⇒ Ensemble Methods :-



Ensemble Algorithms :



- * Randomly pick one dataset
 - * Then again randomly pick another dataset (with replacement)
 - * Continue this for n times
- * Bootstrapping :-
- $D_2, D_2 \dots, D_B$
 - $T_2, T_2 \dots, T_B$
 - * i.e. this B datasets can be constructed.
 - * Each dataset will be a tree
 - * ∴ we will get B number of trees
 - * decision tree algorithm will be applied for each bag.

→ Inference :-

$$x \longrightarrow \hat{y}$$

$$x \longrightarrow T_1 \longrightarrow \hat{y}_1$$

$$x \longrightarrow T_2 \longrightarrow \hat{y}_2$$

$$\vdots$$

$$x \longrightarrow T_B \longrightarrow \hat{y}_B$$

$$\hat{y} = \frac{1}{B} \sum_{i=1}^B \hat{y}_i$$

* Aggregation :-

→ Why is ensemble better than decision tree?
→ What learns?

(Cross Validation)

Random forest Eg:-

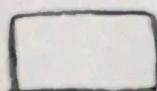
B = 2

	Age	Salary	Spend pattern	Buy Product
1	16	0	high	yes
2	30	90,000	low	No
3	31	90,000	high	No
4	65	6,20,000	Medium	No

2, 3, 3, 4 (random values)	1, 1, 3, 1 (random values)
-------------------------------	-------------------------------

Age	Sal	spat	Buy
30	90	high	No
31	90	low	No
31	90	high	No
65	500	M.d	No

Tree 1



→ Selecting first node

for feature in (Age, Salary)

{

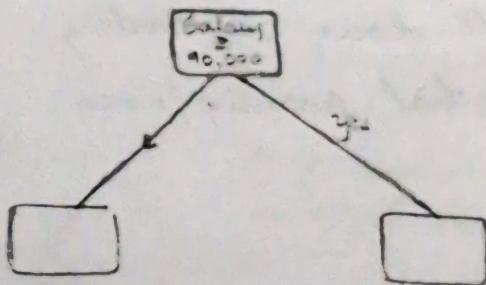
for each value in feature

{

IG = (IG of value). append

IG of Age = 1.2

IG of Salary = 1.2



for feature in (spend pattern, salary)

{

for each value in feature

{

IG = (IG of value). append

IG of SP = 0.2

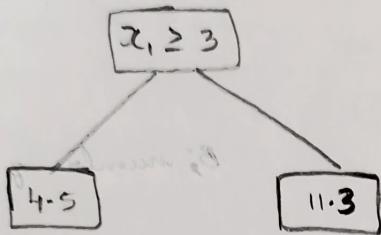
IG of Salary = 1 Append to the tree

x_1	y	\hat{y}
1	3	3
2	6	6
3	8	8
4	11	11
5	15	15

Step 1

$\hat{y} = y$

a) fit a model



b, c)

x_1	y	\hat{y}
1	3	3.45
2	6	6.45
3	8	9.13
4	11	10.13
5	15	16.13

$\lambda = 0.1$

$$\hat{y} = 3 + 0.1(4.5) \quad \text{for } x_1 = 4.5$$

$$\hat{y} = 8 + 0.1(11.3) \quad \text{for } x_1 = 11.3$$

c)

x_1	y	\hat{y}	\hat{y}_{new}
1	3	3.45	-0.45
2	6	6.45	-0.45
3	8	9.13	-1.13
4	11	10.13	-1.13
5	15	16.13	-1.13

Influence:

$x_1 = 5$ → scaling signed

$$= \sum_{i=2}^5 \lambda f_i(x)$$

$$= \lambda f_2(5) + \lambda f_3(5) + \lambda f_4(5) \dots \lambda f_{14}(5)$$

$$= 0.8$$

- 1) Boosted Regression Tree (AdaBoost Regression)
- 2) Boosted Classification Tree (AdaBoost classification)

$$f = f_0(x) + \lambda f_1(x) + \lambda f_2(x) + \dots + \lambda f_n(x)$$

Eg :-

x_1	x_2	y
1	2	4
2	3	5
3	4	6
4	5	7

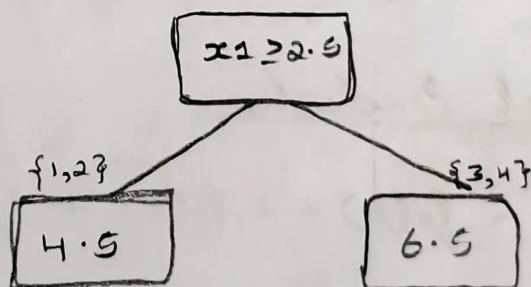
→ Iteration - 1

* $\alpha = y$ Step 1

x_1	x_2	γ
1	2	4
2	3	5
3	4	6
4	5	7

Step 2 :-

a) Fit a model $f_0(x) \rightarrow \gamma$



b) update the model ($\lambda = 0.5$) → c) Update the residuals

$$F(x) = f_0(x) + \lambda f_1(x)$$

$$\gamma^2 = \gamma^1 - \lambda f_0(x^0)$$

$$\gamma^2 = \gamma^0 - \lambda f_0(x^0)$$

$$\gamma^3 = \gamma^3 - \lambda f_0(x^3)$$

$$\gamma^4 = \gamma^4 - \lambda f_0(x^4)$$

c)

$$\hat{z}^{(1)} = 4 - [0.1 * 4.5] = 3.55$$

$$\hat{z}^{(2)} = 5 - [0.1 * 4.5] = 4.55$$

$$\hat{z}^{(3)} = 6 - [0.1 * 6.5] = 5.35$$

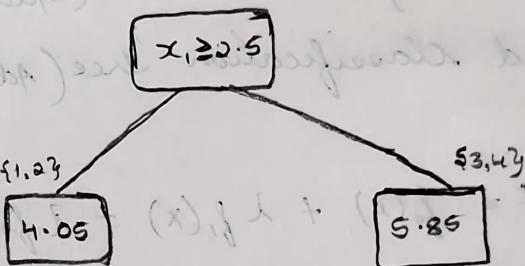
$$\hat{z}^{(4)} = 7 - [0.1 * 6.5] = 6.35$$

\Rightarrow Situation 2

x_1	x_2	z
1	2	3.55
2	3	4.55
3	4	5.35
4	5	6.35

Step 2

\Rightarrow fit a model



b)

c) Update residuals

$$z' = z^1 - \lambda f_2(x^1) = 3.55 - [0.1 * 4.05] = 3.145$$

$$z^2 = z^2 - \lambda f_2(x^2) = 4.55 - [0.1 * 4.05] = 4.145$$

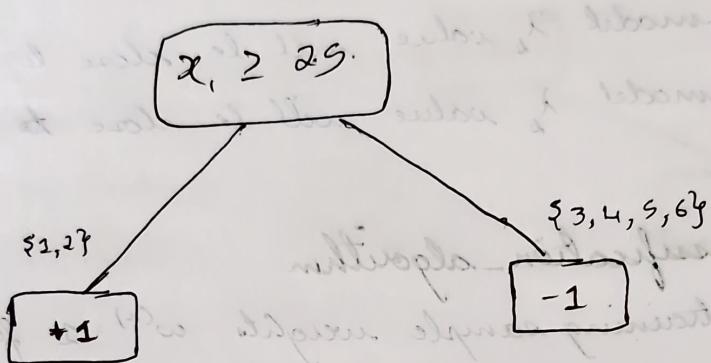
$$z^3 = z^3 - \lambda f_2(x^3) = 5.35 - [0.1 * 5.85] = 4.765$$

$$z^4 = z^4 - \lambda f_2(x^4) = 6.35 - [0.1 * 5.85] = 5.765$$

③ \Rightarrow Final Prediction

$$F(x; \theta) = f_0(x) + \lambda_1 f_1(x) + \lambda_2 f_2(x) + \dots + \lambda_B f_B(x)$$

x_1	x_2	y	w
1	2	1	1
2	3	1	1
3	3	1	1
4	5	-1	2
5	5	-1	1
6	6	-1	1



$$E = \frac{1}{6} \Rightarrow \frac{\text{mismatch}}{\text{total}} = 0.1667$$

$$\lambda = \frac{1}{2} \log_e \left[\frac{1 - 0.1667}{0.1667} \right] \approx 0.808$$

$$\omega^{(1)} = 1 \cdot e^{-0.808} \quad (\text{Correct}) = 0.447$$

$$\omega^{(0)} = 1 \cdot e^{0.808} = 2.236$$

x_1	x_2	y	w
1	2	+1	0.447
2	3	1	0.447
3	3	1	2.236
4	5	-1	0.447
5	5	-1	0.447
6	6	-1	0.447

\Rightarrow Kernel:

$$K(x, y) = \langle \phi(x), \phi(y) \rangle, \quad \langle \cdot \rangle = \text{dot product}$$

$x = (x_1, x_2, x_3)$ * x & y are just 2 random vectors, not features & ground truth
 $y = (y_1, y_2, y_3)$

~~eg~~ $\phi x = \begin{bmatrix} x_1, x_1 \\ x_1, x_2 \\ x_1, x_3 \\ x_2, x_1 \\ x_2, x_2 \\ x_2, x_3 \\ x_3, x_1 \\ x_3, x_2 \\ x_3, x_3 \end{bmatrix}$

$$x = [1, 2, 3]$$

$$y = [4, 5, 6]$$

$$\langle \phi(x), \phi(y) \rangle$$

$$\phi y = \begin{bmatrix} y_1, y_1 \\ y_1, y_2 \\ \vdots \\ y_1, y_3 \end{bmatrix}$$

$$\phi(x) = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 2 \\ 4 \\ 6 \\ 3 \\ 6 \\ 9 \end{bmatrix}$$

$$\phi(y) = \begin{bmatrix} 16 \\ 20 \\ 24 \\ 20 \\ 32 \\ 24 \\ 30 \\ 36 \end{bmatrix}$$

$$\begin{aligned} \langle \phi(x), \phi(y) \rangle &= 16 + 40 + 72 + 40 \\ &\quad + 100 + 160 + 72 \\ &\quad + 180 + 324 \\ &= 1024 \end{aligned}$$

* Can we do something on original space which is equivalent to higher space?

* This is where Kernel helps

$$K(x, y) = (\langle x, y \rangle)^2$$

$$= (32)^2$$

$$= 1024$$

* We are able to mimic what happened in higher dimension in lower dimensional space using kernel function

Eg:

x_1	x_2	y
1	2	0
3	4	0

$$\phi(x) = \begin{bmatrix} x_1^2 \\ x_2^2 \\ x_1 x_2 \end{bmatrix}$$

$$x = [x_1, x_2]^T \quad \# 2 \text{ dimensional space}$$

$$x = \mathbb{R}^2 \quad \phi(x) \in \mathbb{R}^3 \quad o \in \mathbb{R}^3 \Rightarrow [o_1, o_2, o_3]^T$$

$$o = \sum_{i=1}^n \beta_i \phi(x^i)$$

$$o = \beta_1 \sum_{i=1}^2 \beta_i \phi(x^i)$$

① Compute each sample to higher dimensional space

$$\phi(x^1) = \begin{bmatrix} 1 \\ 4 \\ 2 \end{bmatrix}_{3 \times 1}$$

$$\phi(x^2) = \begin{bmatrix} 9 \\ 16 \\ 12 \end{bmatrix}_{3 \times 1}$$

② Compute the o

$$o = \beta_1 \phi(x^1) + \beta_2 \phi(x^2)$$

$$\begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} = \begin{bmatrix} \beta_1 + 9\beta_2 \\ 4\beta_1 + 16\beta_2 \\ 2\beta_1 + 12\beta_2 \end{bmatrix}$$

\Rightarrow SVM

1) Maximal Margin classifier (Constrained optimization problem)

Maximize M

$$\beta_0, \beta_1, \beta_2, \dots, \beta_p, M$$

— ①

Subject to

$$\sum_{j=1}^p \beta_j^2 = 1 \quad — ② \text{ # similar to } l_2 \text{ regularization}$$

$$y_i(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p) \geq M \quad — ③$$

for all $i = 1 \dots n$

$$d = \frac{|\beta_0 + \beta_1 x_1 + \beta_2 x_2|}{\sqrt{\beta_1^2 + \beta_2^2}}$$

Eg: Three lines (hyperplane)

$$2x_1 + 3x_2 - 5 = 0$$

$$-x_1 + 4x_2 + 7 = 0$$

$$5x_1 - 12x_2 + 10 = 0$$

Data set

x_1	x_2	y
3	4	1
2	3	1
1	-1	-1
-2	1	-1

$$M_2 = \frac{+7 - 1(3) + 4(4)}{\sqrt{2^2 + 4^2}}$$

$$M_2 = \frac{-5 + 2(4) + 3(3)}{\sqrt{4 + 9}}$$

$$= \frac{-5 + 4 + 9}{\sqrt{13}}$$

$$\frac{8}{\sqrt{13}}$$

$$M_3 = \frac{1 - 5 + 2(2) + 3(-2)}{\sqrt{13}}$$

$$= \frac{1 - 5 + 2 - 3}{\sqrt{13}}$$

$$= \frac{+6}{\sqrt{13}}$$

$$M_4 = \frac{-5 + 2(2) + 3(0)}{\sqrt{13}}$$

$$= -5 - 4 + 3$$

$$= \frac{+6}{\sqrt{13}}$$

$$\min = \frac{6}{\sqrt{13}} = 1.664$$

\Rightarrow Hyperplane - 2

$$M_1 = \frac{+7 - 1(3) + 4(4)}{\sqrt{2^2 + 4^2}}$$

$$M_2 = \frac{+7 - 1(2) + 4(3)}{\sqrt{17}}$$

$$= \frac{7 - 3 + 16}{\sqrt{17}}$$

$$= \frac{+7 - 2 + 12}{\sqrt{17}}$$

$$= \frac{20}{\sqrt{17}}$$

$$= \frac{17}{\sqrt{17}}$$

$$M_3 = \frac{+7 - 1(1) + 4(-2)}{\sqrt{17}}$$

$$M_4 = \frac{+7 - 1(-2) + 4(2)}{\sqrt{17}}$$

$$= \frac{+7 - 1 - 8}{\sqrt{17}} = \frac{-2}{\sqrt{17}}$$

$$= \frac{7 + 2 + 4}{\sqrt{17}} = \frac{13}{\sqrt{17}}$$