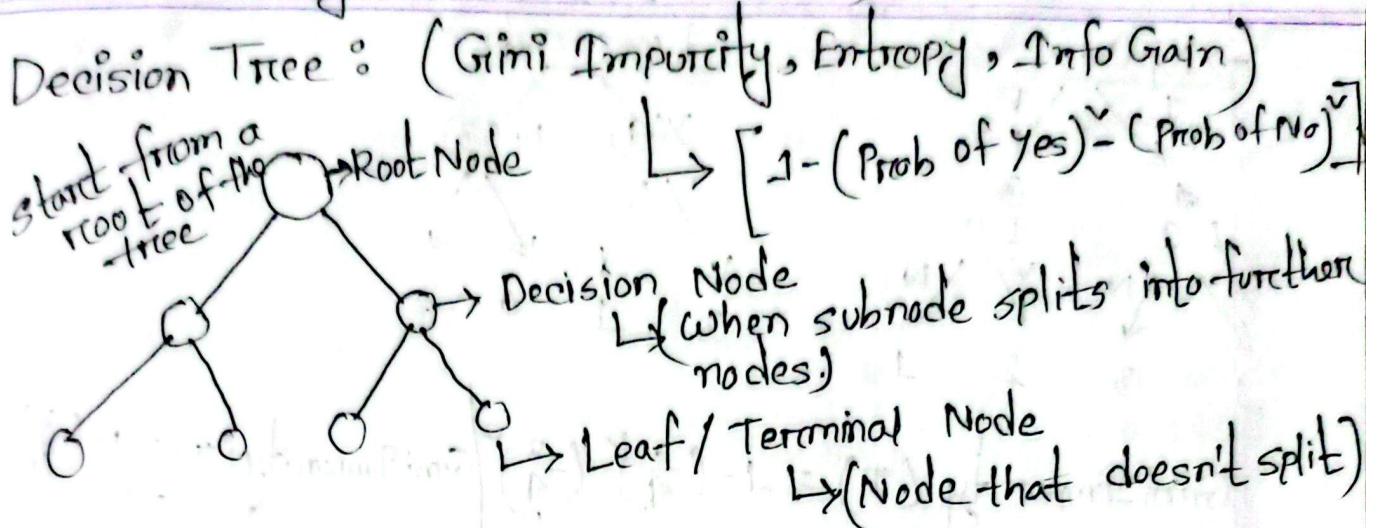


Machine Learning :

(Final Topic)



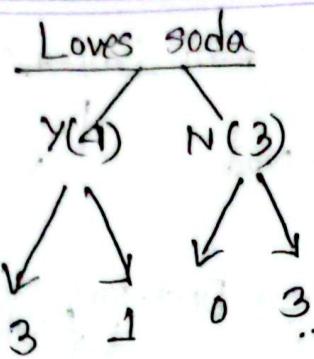
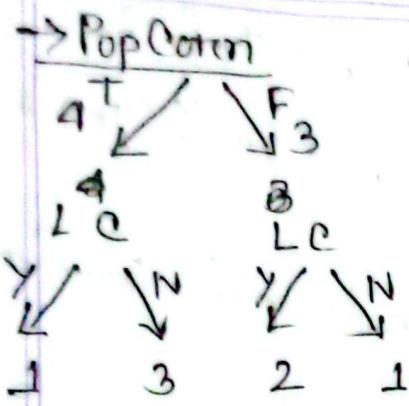
We stop, Splitting it there are no decision rules left to the group impurity.

Adv: Easy to interpret, straight-forward visualization.

Disadv: Single DT more prone to overfitting.
→ Early stop specifying the layers.

Loves Popcorn	Loves Soda	Age	Loves Cool as ice
Y	Y	7 → Y	N
Y	N	12 → Y	N
N	Y	18 → N	Y
N	Y	35 → N	Y
Y	Y	38 → N	Y
Y	N	50 → N	N
N	N	83 → N	N

For (Y, Y, 15) what should be the output?



$$\text{Gini Impurity } (T) = 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2$$

$$= 0.375$$

$$\text{u } (F) = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2$$

$$= 0.444$$

weighted avg of

$$\text{love popcorn} = \left(\frac{4}{7} \times 0.375\right) + \left(\frac{3}{7} \times 0.444\right)$$

$$= 0.405.$$

$$\text{Gini Impurity } (T) = 1 - \left(\frac{3}{7}\right)^2 - \left(\frac{4}{7}\right)^2$$

$$= 0.375$$

$$\text{u } (F) = 1 - \left(\frac{0}{3}\right)^2 - \left(\frac{3}{3}\right)^2$$

$$= 0$$

$$\text{w of loves-soda}$$

$$= \left(\frac{4}{7} \times 0.375\right) + \left(\frac{3}{7} \times 0\right)$$

$$= 0.214.$$

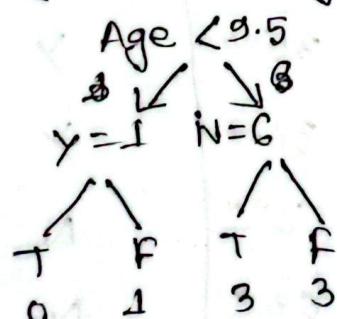
* For numerical value,

(i) Firstly, Row ঘোর হোর খেকে বড় সাজা।

(ii) Avg age of all adjacent people.

(iii) Calculate gini impurity for each avg age.

Age
7 → 9.5
12 → 15
18 → 26.5
35 → 36.5
38 → 41
50 → 66.5
83

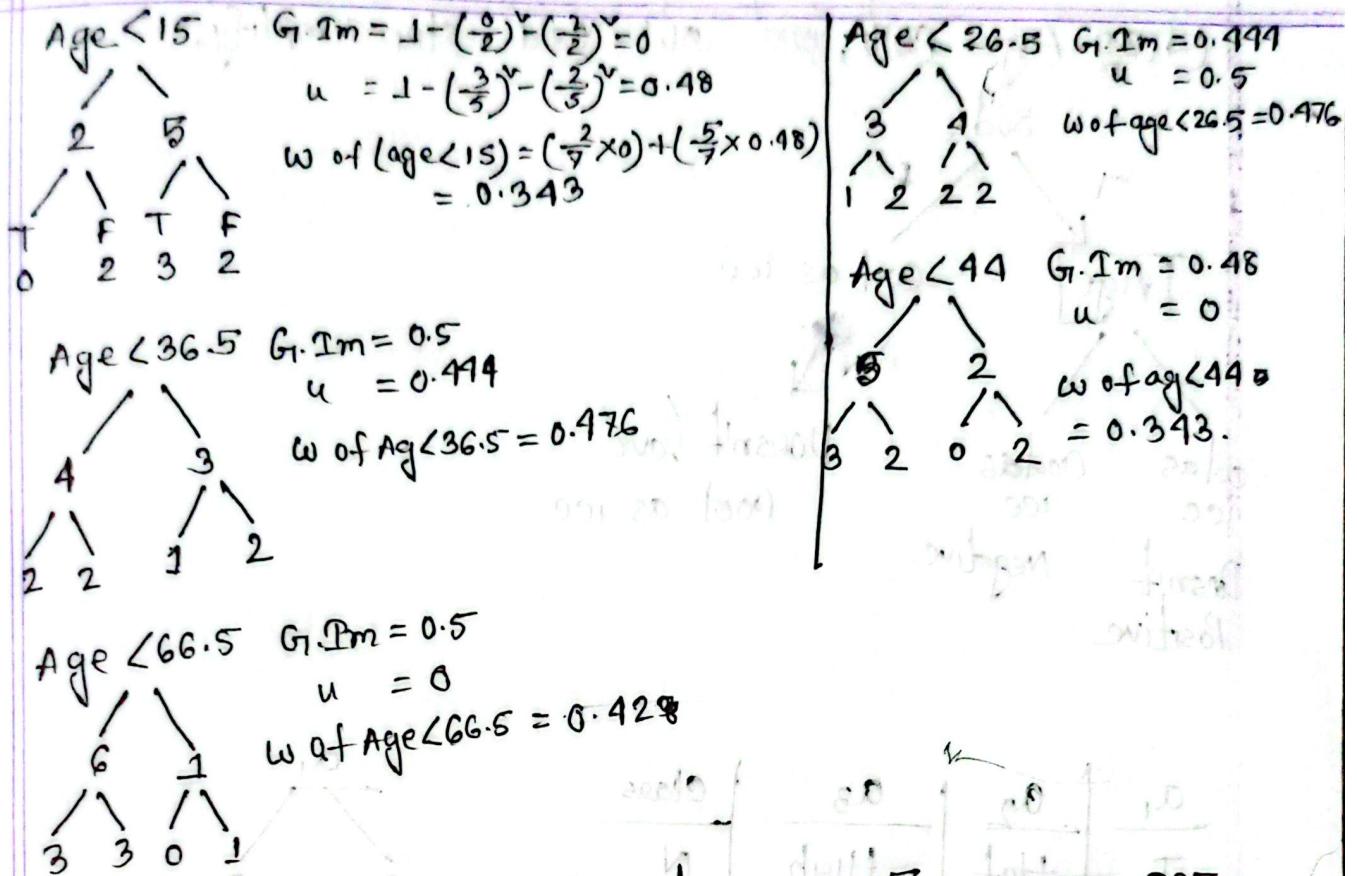


$$G.I.m = 1 - \left(\frac{0}{7}\right)^2 - \left(\frac{1}{7}\right)^2 = 0.$$

$$\text{u } = 1 - \left(\frac{3}{6}\right)^2 - \left(\frac{3}{6}\right)^2 = 0.5$$

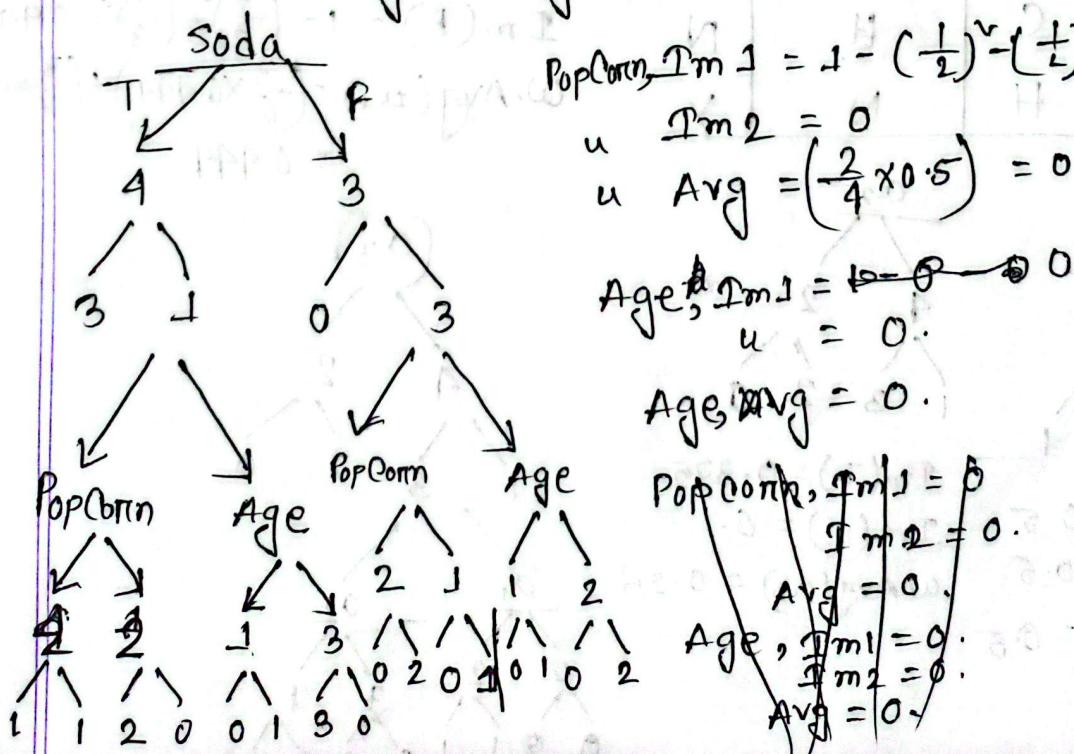
$$\text{w of age} < 9.5 = \left(\frac{1}{7} \times 0\right) + \left(\frac{6}{7} \times 0.5\right)$$

$$= 0.429$$

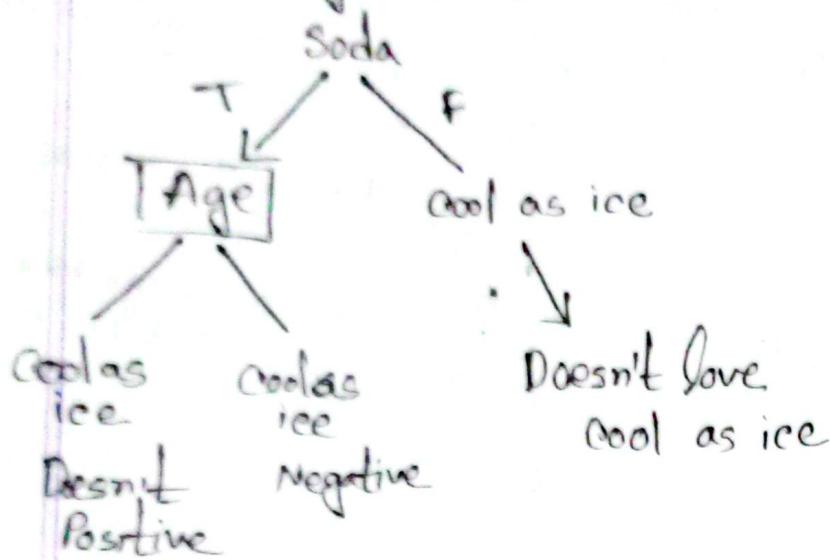


∴ Among those values, lowest. 0.343 [Same यहाँ जूँ है
लिखा गया है]

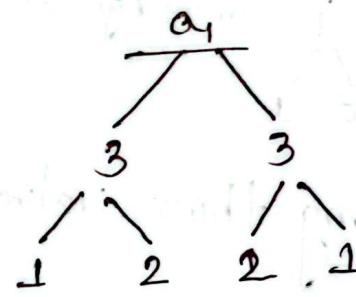
Lower impurity among all columns - Soda (Top of the tree).



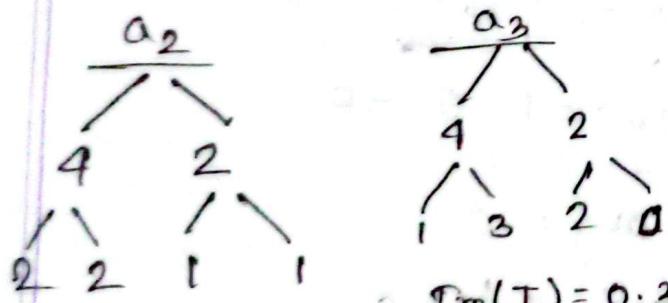
using Age attr, one out 2nd step simplify,



a_1	a_2	a_3	Class
T	Hot	High	N
T	H	H	N
F	H	H	Y
F	Cool	Normal	Y
F	C	H	N
T	H	N	Y

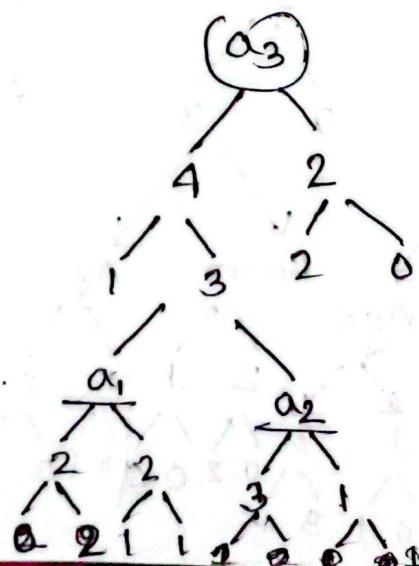


$$\begin{aligned} Im(T) &= 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = 0.444 \\ Im(F) &= 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = 0.444 \\ wAvg(a_1) &= \left(\frac{3}{6} \times 0.444\right) + \left(\frac{3}{6} \times 0.444\right) \\ &= 0.444 \end{aligned}$$



$$\begin{aligned} Im(T) &= 0.5 \\ Im(F) &= 0.5 \\ wAvg(a_2) &= 0.5 \end{aligned}$$

$$\begin{aligned} Im(T) &= 0.375 \\ Im(F) &= 0. \\ wAvg(a_3) &= 0.25. \end{aligned}$$



T, H, H,

$$\alpha_1, \text{Im}(T) = 0$$

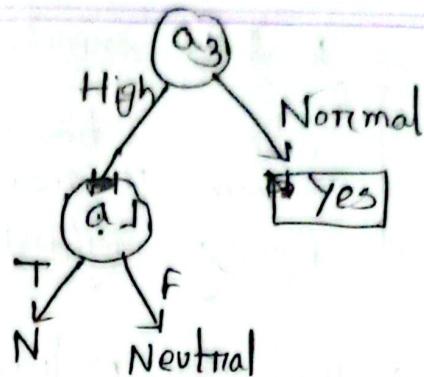
$$u(F) = 0.5$$

$$w. \text{Avg}(\alpha_1) = 0.25$$

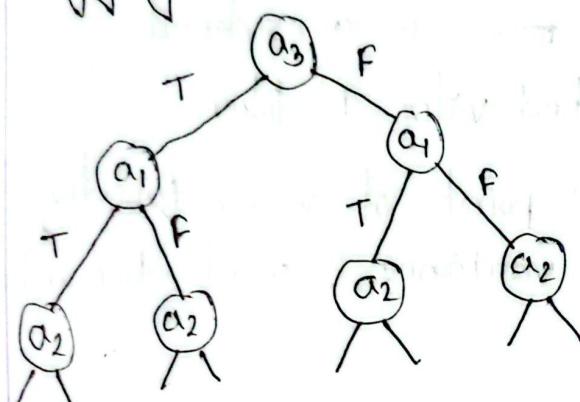
$$\alpha_2, \text{Im}(T) = 0.444$$

$$u(F) = 0$$

$$w. \text{Avg}(\alpha_2) = 0.33$$

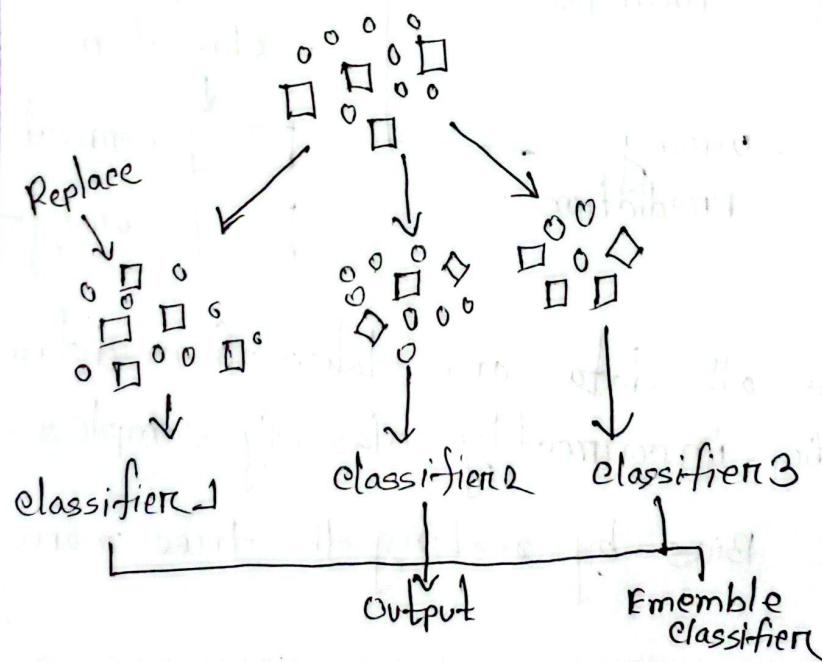


Bagging:

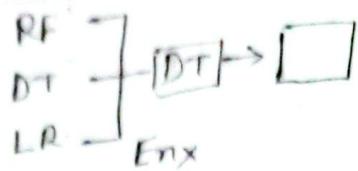


Bootstrap Aggregation:

sampling with replacement from the original data.



- * Bagging
 - ↳ parallelly execute
 - ↳ Takes less time to train the model
 - ↳ Equally weighted all

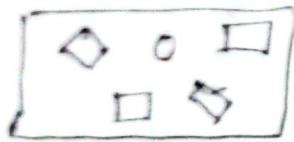


- * Bagging is used for reducing variance.
- * Prevents overfitting
- * Improves robustness.

- ✓ for classification majority voting is considered.
- ✓ for regression, avg weighted value is taken.

From bias variance trade off point of view, Boosting starts with low bias + high variance & work towards reducing variances.

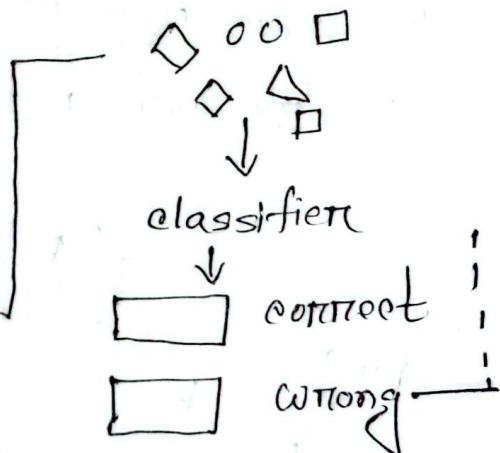
works best with strong & cplx model



Correct Prediction



wrong Prediction



- * No extra data all data are taken from dataset
- * weights up the incorrectly classify samples.

Boosting → Reduce Bias by making the tree more complicated.

✓ AdaBoost
✗ GBBoost

COMPA

$y = 9$

< 9

> 9

< 9

$y = 9$

$y = 9$

Predicted Job officer

Y

N

Y

N

Y

Y

Y

Actual Job officer

Y

Y

N

N

Y

Y

Adaboost steps:

- (i) Assigning weights
- (ii) Creating all decision stumps
- (iii) Choosing the best decision stumps.
- (iv) Calculating weight
- (v) Adjusting w
- (vi) Normalizing w
- (vii) Preparing data for next step.
- (viii) Assigning next weight.

* Initial weight assigned to each item = $\frac{1}{n} = \frac{1}{6}$.

① Compute weighted error,

$$\text{Error} \rightarrow \sum_{j=1}^N \delta^{*} h_j(d_j) w_t(d_j)$$

for correct prediction
+ " wrong "

$$\text{Error}_{\text{corr}} = 4 \times 0 \times \frac{1}{6}$$

$$+ \text{Neg} = 2 \times 1 \times \frac{1}{6} = 0.333$$

② Compute the weight of each weak learner

$$\alpha_{\text{CGPA}} = \frac{1}{2} \left(\frac{\ln(1 - \text{Error})}{\text{Error}} \right) = \frac{1}{2} \left(\frac{\ln(1 - 0.333)}{0.333} \right) = 0.347$$

③ Calculate Normalizing factor

$\Rightarrow \text{CGPA} = w_t(\text{correct}) \times \text{num of correct classification}$

$$+ w_t(\text{wrong}) \times \text{num of wrong classification}$$

$$= \left(\frac{1}{6} \times 4 \times e^{-0.347} \right) + \left(\frac{1}{6} \times 2 \times e^{0.347} \right)$$

$$= 0.0428$$

④ Update the weight of all data instances,

$$w_t(d_j)_{i+1} = \frac{w_t(d_j) \text{ of correct instance} \times e^{-\alpha}}{\frac{1}{6} \times e^{-0.347} \times \text{CGPA}}$$

$$= \frac{0.0428}{0.0428} = 0.1249$$

$$w_t(d_j)_{i+1} = \frac{w_t(d_j) \text{ of incorrect instance} \times e^{\alpha}}{\frac{1}{6} \times e^{0.347}}$$

$$\Theta = \frac{\frac{1}{6} \times e^{0.247}}{0.5928} = 0.2501.$$

Sample	Actual	Predicted
1	0	1
2	0	0
3	1	0
4	1	1
5	1	0

① Initially, we are taking random weight $= \frac{1}{n} = \frac{1}{5} = 0.2$

② Compute weighted error

always 0 for correct prediction as $H_i(d_j) = 0$,

$$\text{Error}_{\text{wrong}} = 3 \times 1 \times 0.2 = 0.6.$$

③ Compute learners weight

$$\alpha = \frac{1}{2} \frac{\ln(1 - \text{err})}{\text{err}} = \frac{1}{2} \ln \frac{1 - 0.6}{0.6} = -0.203.$$

④ Compute Normalizing factor

$$= 0.2 \times 2 \times e^{-(-0.203)} + 0.2 \times 3 \times e^{-0.203} = 0.9798.$$

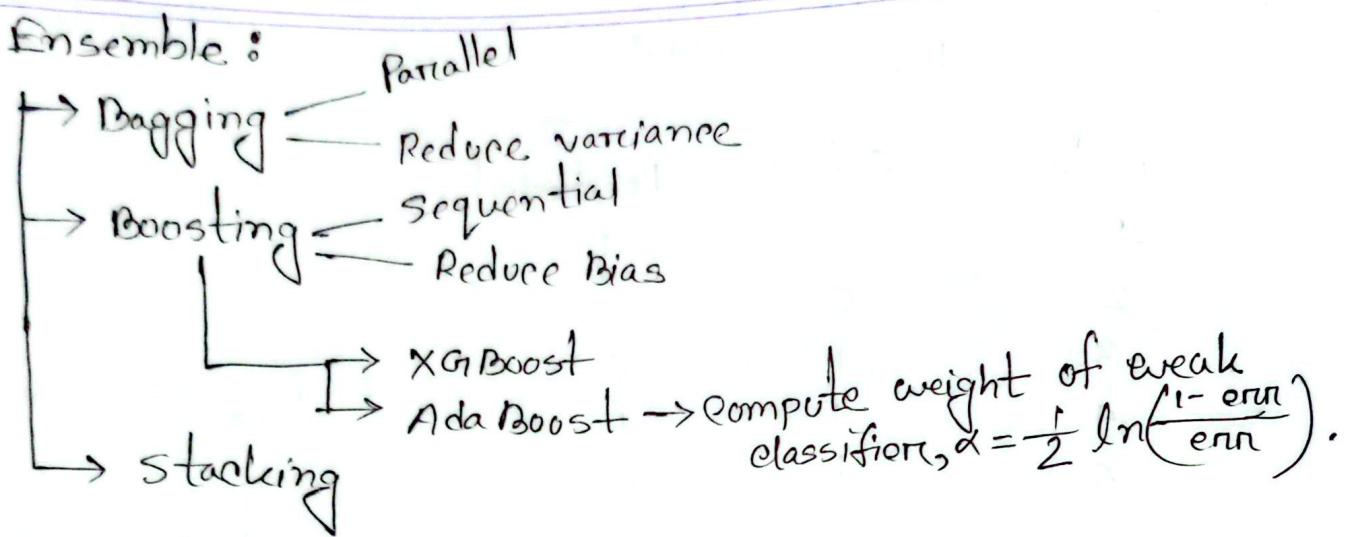
⑤ Update the weight of all data instances.

$$\oplus \rightarrow \frac{0.2 \times e^{0.203}}{0.9798} = 0.245$$

$$\ominus \rightarrow \frac{0.2 \times e^{-0.203}}{0.9798} = 0.167$$

In this math $\alpha \ominus$

This means the weak learner performs worse than guessing error. So, AdaBoost will effectively invert its predictions in future & down weight its import.



XGBoost :

For classification:

- ① Probability based on binary classification
- ② Classification \rightarrow Base value = $\log\left(\frac{\text{Prob}}{1 - \text{Prob}}\right)$
- ③ Residual = $\text{Act} - \text{Pred}$
- ④ Similarity weight = $\sum (\text{Residual})^2 / \sum \text{Pred} (1 - \text{Pred})$
- ⑤ Calculate gain = All leaf similarity - similarity of Root.

salary	credit	Approval	Res
$<= 50$	Bad	0	-0.5
$<= 50$	Good	1	0.5
> 50	Bad	0	-0.5
> 50	Good	1	0.5
$<= 50$	Normal	0	-0.5

For salary %

① Probability = 0.5 as Binary.

② Residual = $(Ae - P_{pred})$

③ Base model = $\log \left(\frac{P_{pred}}{1-P_{pred}} \right) = \log \left(\frac{0.5}{1-0.5} \right) = \log(1) = 0.$

④ Similarity weight

⑤ for $<= 50$, $SW = \frac{(-0.5 + 0.5 - 0.5)^2}{\{0.5(1-0.5)\}^2} \times 3 = 0.33.$

⑥ For > 50 , $SW = \frac{(0.5 - 0.5)^2}{\{0.5(1-0.5)\}^2} \times 2 = 0.$

⑦ SW for root = $\frac{(-0.5 + 0.5 - 0.5 + 0.5 - 0.5)^2}{\{0.5(1-0.5)\}^2} \times 5 = 0.2$

⑧ Grain = All leaf similarity - similarity of Root
 $= (0.33 + 0) - 0.2 = 0.13.$

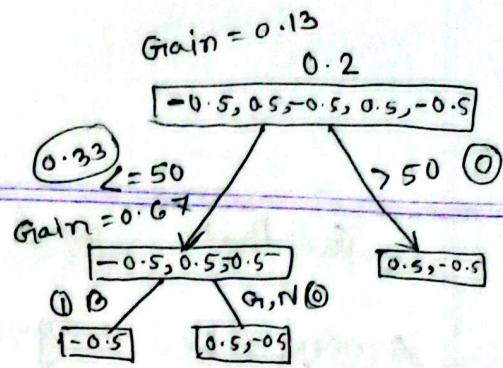
for credit %

① SW of Bad = $\frac{(-0.5)^2}{\{0.5(1-0.5)\}^2} \times 1 = 1.$

sw of Good, Normal = $\frac{(0.5 - 0.5)^2}{\{0.5(1-0.5)\}^2} \times 2 = 0.$

② Info Grain, $(1+0) - 0.33 = 0.67.$

Among Bad, Good, Normal, we took Bad onesite, others are on the other side as we have to do binary classification. Similarity, we will calculate the combination.



$G_1 - B, N$ and $N \rightarrow G_1, B$.

among the 3 gain, the highest value is kept for next step.

considering the calculated one as highest gain,

New update prob of given point ($\angle = 50$, Bad) $\rightarrow 0.1 / 0.2 / 0.5$ (let)

$$= \sigma(\text{Base model} + \lambda \times \text{sw})$$

$$= \sigma(0 + 0.1 \times 1)$$

$$= \sigma(0.1)$$

$$= \frac{1}{1+e^{-0.1}} = 0.6.$$

∴ New feed residual $= A_e - P_{\text{fit}} = 0 - 0.6 = -0.6$.
(for 1st point)

New value becomes

-0.6 instead of -0.5. similarly, the process will continue with all new values.

Parameters

λ - Regularization Parameter

↳ Value λ \rightarrow ∞ , impact of outliers \rightarrow ∞ ,

↳ Pruning \rightarrow Aggressive/careful

2. Gamma \rightarrow Threshold - Auto pruning, etc

overfitting classification $\Rightarrow \gamma = \text{Prob}(1 - \text{Prob})$

$$= 0.5 (1 - 0.5) = 0.25$$

Ensemble (1.5)

→ AdaBoost \rightarrow S
→ XGBoost \rightarrow S
→ theory \rightarrow S.

if gain < 0.25 , we can prune that part from the tree, ε - how fast convergence happen (speed)
Prune (child count)

XGBoost:

Just an extremely optimized gradient boosting model which by itself is a powerful ML Algo,

why XGBoost is better

→ Missing values automatically handled.

→ Cross validation.

→ Interpolates both L₁ and L₂ regularization technique to prevent overfitting and enhance model generalization.

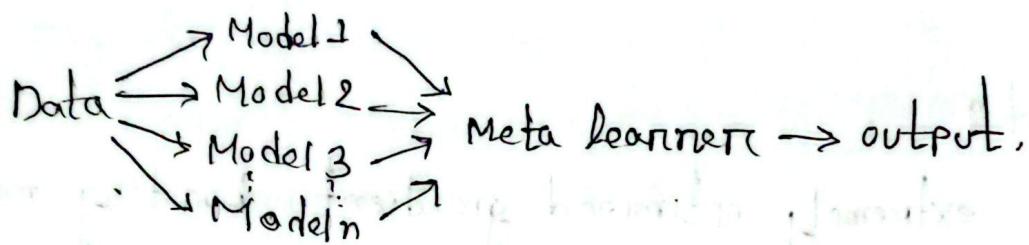
→ Parallelization: At a time, more work possible.

→ Tree pruning: If possible, it prunes to make the tree smaller.

Stacking:

- Ensemble method that combines heterogeneous weak learners, combined using a meta learning algorithm.
- Bagging and boosting can be used as ensemble at start of stacking.

Though ensemble is not interpretable, it can give highest accuracy. So, it is used when accuracy is our main target.



stacking steps :

Individual estimators feed their prediction to the meta learn.

Meta learners take them as input features

Work on how to best combine these predictions to improve the overall performance.

Their goal : correct the weakness of individual models by finding optimal weight.

thus, ensures higher accuracy and generalization

XG Boost Regression:

- I Base = Avg of output column.
- II Residual = Act - Base.
- III Similarity weight = $\frac{\sum \text{Residual}^2}{\text{no of Residual} + \lambda}$

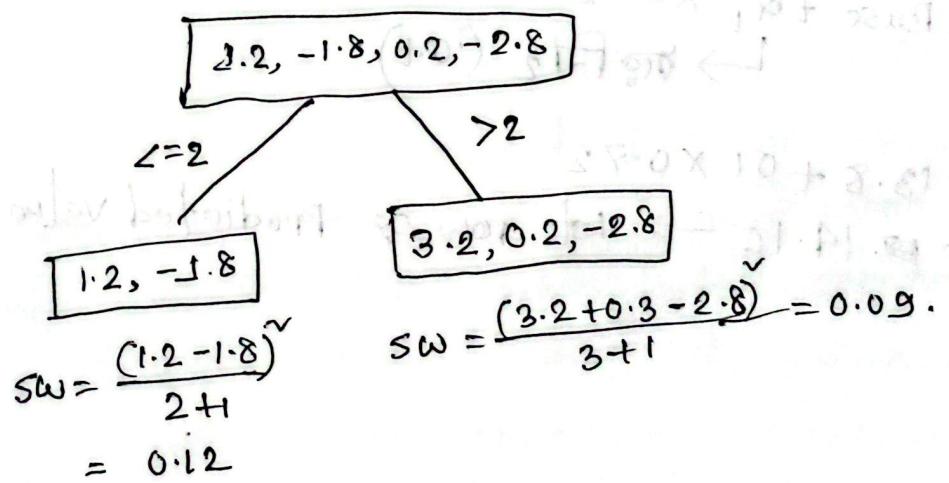
- IV Info Gain = SW of a2 leaf - Root.
- V Calculate new probability

Ball	color	size	Price	Residual
1	Red	3	15	$15 - 13.8 = 1.2$
2	Blue	25	12	$12 - 13.8 = -1.8$
3	Red	4	17	$17 - 13.8 = 3.2$
4	Blue	35	14	$14 - 13.8 = 0.2$
5	Red	2	11	$11 - 13.8 = -2.8$

$$\text{Base value} = \frac{\sum \text{Price}}{n}$$

$$= \underline{\underline{13.8}}$$

102

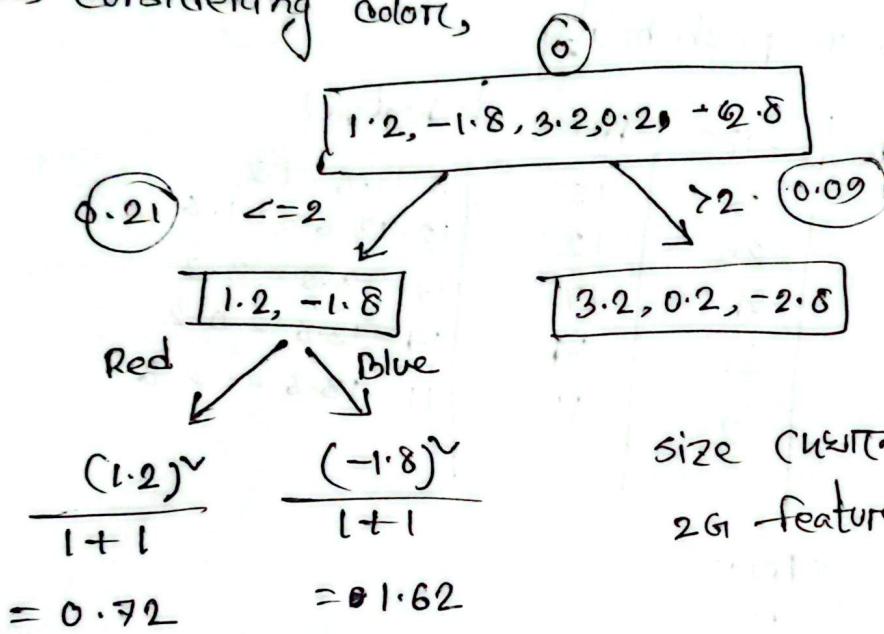


$$\text{Root} = \frac{(1.2 - 1.8 + 3.2 + 0.2 - 2.8)}{5+1}^{\vee} = 0$$

$$\therefore \text{Gain} = (0.12 + 0.09) - 0 = 0.2$$

for making the calculation easier, we are considering,
the above one gives the highest gain.

Now, considering colors,



size (छानोड़ पॉलो वाले)
2G feature freq (मुज्जरे 2G)

considering the fast cost,

$$= \text{Base} + \alpha_1 \times 0.72 \rightarrow \text{Possible case value} \rightarrow \begin{matrix} \text{first row} \\ \geq 0.1 \end{matrix}$$

$$= 13.8 + 0.1 \times 0.72$$

$$= 14.16 \rightarrow \text{1st row } \geq \text{Predicted value.}$$

Performance Matrix :

TP : The positive class is predicted as positive.

TN : " Negative " " " " " Negative "

FN : " Positive " " " " " " " "

FP : " Negative " " " " " " " " Positive . "

Accuracy : Overall how often.

$$Acc = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$\text{F1-score} = \frac{2 \times P \times R}{P + R}$$

$$\text{False Neg} = 1 - R$$

$$\text{u Pos} = 1 - \text{specificity}$$

$$\text{u Neg} = \frac{FP}{FP + TN}$$

		Predicted Class	
		No	Yes
Actual Class	No	45	5
	Yes	5	95

$$TN = 45$$

$$FP = 5$$

$$FN = 5$$

$$TP = 95$$

$$Acc = \frac{95 + 45}{95 + 45 + 5 + 5} = 0.933$$

$$R = \frac{95}{95 + 5} = 0.95$$

$$P = \frac{95}{95 + 5} = 0.95$$

$$f_1 = \frac{2 \times 0.95 \times 0.95}{0.95 + 0.95} = 0.95$$

→ 10,000 patient were tested to check corona effect. Among them, 9000 were healthy & 1000 were actually sick. Those who are sick, when tested, 620 were found \oplus & 380 where \ominus similarly, when checking the healthy patient, 180 were found \oplus corona & 8820 were \ominus found precision, recall, F1.

		Predicted
		Yes No
Actually	Yes	620 380
	No	180 8820

True class	Predicted class		
	1	2	3
1	8	2	0
2	1	9	0
3	1	2	7

TP	FN
FP	TN

1 → 2, 3

2 → 1, 3

3 → 1, 2

B	2
2	18

9	81
81	16

7	3
0	20

$$\text{Total Acc} = \frac{A_1 + A_2 + A_3}{3}$$

$$P = \frac{P_1 + P_2 + P_3}{3}$$

$$R = \frac{R_1 + R_2 + R_3}{3}$$

$$f1 = \frac{P_1 + R_1 + f1_3}{3}$$

Predicted Class

	A	B	AB	O
A	350	11	30	9
B	15	359	21	5
AB	35	35	320	10
O	20	20	27	333

find the acc & f1 score of the example.

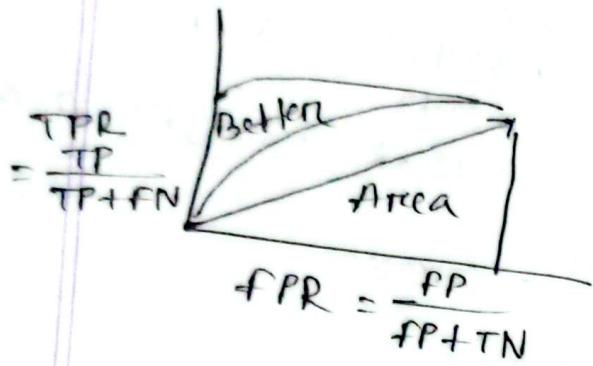
$A \rightarrow B, AB, O$	\rightarrow	<table border="1"> <tr><td>350</td><td>50</td></tr> <tr><td>70</td><td>1130</td></tr> </table>	350	50	70	1130	$ACC = \frac{1480}{1600} = 0.925$
350	50						
70	1130						
$B \rightarrow A, AB, O$	\rightarrow	<table border="1"> <tr><td>359</td><td>91</td></tr> <tr><td>66</td><td>1131</td></tr> </table>	359	91	66	1131	$P = 0.833, R = 0.875, f1 = 0.854$
359	91						
66	1131						
$AB \rightarrow A, B, O$	\rightarrow	<table border="1"> <tr><td>320</td><td>80</td></tr> <tr><td>78</td><td>1122</td></tr> </table>	320	80	78	1122	$ACC = 0.933$
320	80						
78	1122						
$O \rightarrow A, B, AB$	\rightarrow	<table border="1"> <tr><td>333</td><td>67</td></tr> <tr><td>24</td><td>1176</td></tr> </table>	333	67	24	1176	$P = 0.845, R = 0.898, f1 = 0.87$
333	67						
24	1176						
			$ACC = 0.901$				
			$P = 0.804, R = 0.8, f1 = 0.802$				
			$ACC = 0.943$				
			$P = 0.933, R = 0.832, f1 = 0.88$				
$ACC = \frac{0.925 + 0.933 + 0.901 + 0.943}{4} = 0.925$							
$f1 = \frac{0.854 + 0.87 + 0.802 + 0.88}{4} = 0.852$							

Confusion / Error Matrix - Used to decide the performance of a classification model on test data.

Instead of being overwhelmed with confusion matrix, ROC graph provides a simple way to summarize all of the info.

ROC Plot TPR vs fPR at diff classification threshold it shows how well classifier distinguishes \oplus & \ominus classes it can be 0 to 1. Higher ROC AUC indicates better performance. A perfect model would have AUC of 1, while a random make would have 0.5.

ROC - Receiver Operating characteristic curve.
 AUC - Area under the curve.



যদি কোন ফল পুরো স্বতন্ত্র হয়ে থাকে এবং এটি কোন স্বতন্ত্র পুরো নয়, তবে এটি একটি অসুবিধা।

Our idea point $(0,1)$
 which is rarely can be
 done by the classifier.

ROC

Plot the ROC curve

Class	Prob	TP	FP	TPR	FPR
P	0.9	1	0	1/5	0/5
P	0.8	2	0	2/5	0/5
N	0.7	2	1	2/5	1/5
P	0.6	3	1	3/5	1/5
P	0.55	4	1	4/5	1/5
N	0.54	4	2	4/5	2/5
N	0.53	4	3	4/5	3/5
N	0.51	4	4	4/5	4/5
P	0.5	5	4	5/5	4/5
N	0.4	5	5	5/5	5/5



<u>x</u>	<u>y</u>	<u>Ac</u>	<u>Pred</u>
l_1	80	80	75
l_2	90	75	85
l_3	100		
l_4	110		
l_5	120		

suppose, now predicted values for l_6 and l_7 are 75 & 85. The actual value was 80 and 75. calculate Mean Absolute error & mean square error.

$$\textcircled{I} \text{ Mean Absolute Error} = \frac{1}{n} \sum_{i=0}^{n-1} |(y_i - \hat{y}_i)|$$

$$= \frac{1}{2} (|80 - 75| + |75 - 85|)$$

$$= \frac{15}{2} = 7.5$$

$$\textcircled{II} \text{ Mean Square Error, } MSE = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \hat{y}_i)^2$$

$$= \frac{1}{2} ((80 - 75)^2 + (75 - 85)^2)$$

$$= 62.5$$

clustering:

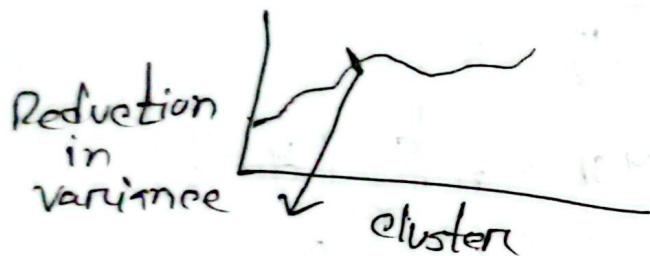
→ kmeans clustering

→ DBScan.

\textcircled{I} kmeans:

- (I) select num of clusters you want to identify in your data. This is the k in kmeans clustering.
- (II) Randomly select k values - these are initial clusters.
- (III) Measures the distances bet'n the first point & 3 int initial clusters.

- IV Assign the first point to nearest cluster.
similarly, continue for all.
 - V Calculate the mean of 3 clusters.
 - VI Iterate I to V until clusters no longer change.
- Selecting the value of k .
- I Total variation within each clusters অথবা
সর্ব তার, k এর মান নির্দেশ করে।
 - II Reduction in variance for value k .



Hierarchical clustering:

- No need to predefine the value of k , builds a tree like structure of clusters.
- Similar data points are in same clusters, otherwise in different clusters.

Tree can be built in 2 ways:

- I Agglomerative Bottom up Approach.
- II Divisive Top Down

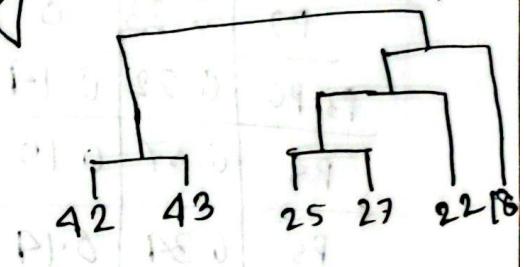
Consider the following sets of 6, dimensional data points $18, 22, 25, 27, 42, 43 \rightarrow$ sort করে ২ সাইde দ্বয়ারা,

apply the agglomerative hierarchical clustering to build the clustering dendrogram (tree like structure)

step:

- Merge the clusters using min distance & update the proximity matrix accordingly,
- Clearly present the proximity matrix in each iteration.

	18	22	25	27	42	43
18	0	4	7	9	24	25
22	4	0	3	5	20	21
25	7	3	0	2	17	18
27	9	5	2	0	15	16
42	24	20	17	15	0	1
43	25	21	18	16	1	0



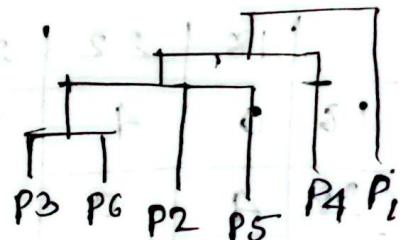
18	0	4	7	9	24
22	4	0	3	5	20
25	7	3	0	2	17
27	9	5	2	0	15
42,43	24	20	17	15	0

18	0	4	7	9	24	18	0	4	24
22	4	0	3	5	20	22	4	0	20
25	7	3	0	2	17	25	7	3	17
27	9	5	2	0	15	27	9	5	15
42,43	24	20	17	15	0	42,43	24	20	0

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

	<u>x</u>	<u>y</u>	P1	P2	P3	P4	P5	P6	PG
P1	0.40	0.53	0	0.23	0.22	0.37	0.39	0.24	
P2	0.22	0.38	0.23	0	0.14	0.19	0.19	0.24	
P3	0.35	0.32	0.22	0.14	0	0.16	0.28	0.10	
P4	0.26	0.19	0.37	0.19	0.16	0	0.28	0.22	
P5	0.08	0.91	0.39	0.14	0.28	0.10	0.28	0	0.39
PG	0.45	0.30	0.29	0.24	0.28	0.22	0.39	0	

	P1	P2	P3, PG	P4	P5
P1	0	0.23	0.22	0.37	0.39
P2	0.23	0	0.14	0.19	0.14
P3, PG	0.22	0.14	0	0.16	0.28
P4	0.37	0.19	0.16	0	0.28
P5	0.39	0.14	0.28	0.28	0



	P1	P2, P3, PG, P5	P4
P1	0	0.23	0.37
P2, P3, PG, P5	0.23	0	0.19
P4	0.37	0.19	0

	P1	P2, P3, P4, PG, P6
P1	0	0.23
P2, P3, P4, PG, P6	0.23	0

3. (1,5)	4. (6,8)
2. (1,4)	5. (5,9)
1. (2,2)	

* Ensemble 10/15

- XG Boost
- Adaboost

Feature selection (5)
(PCA, UMAP, t-SNE)

Performance Matrix (5)

Clustering (5)

→ Kmeans

→ Hierarchical

Deep learning (5)

DT *

Reinforcement learning:

Optimal behaviour is learnt through interaction & observations of how it responds.

(i) Agent or learner

Environment the
agent interacts with

(i) observation

(ii) Reward

(iii) Action

* How exactly it chooses action in general?

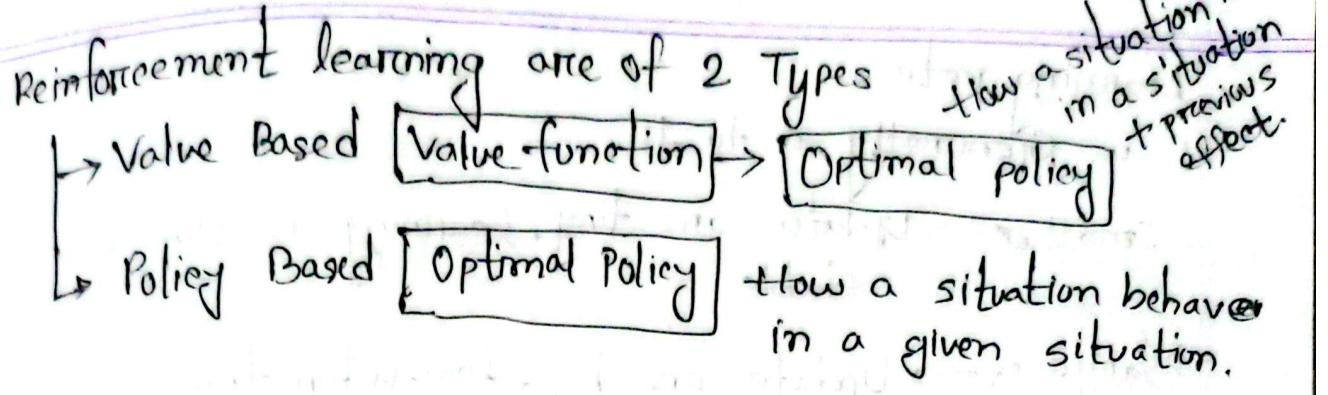
① Exploration: Act of exploring the environment to find out info about it.

② Exploitation: Act of exploiting info in order to maximize return.

Balance: Epsilon Greedy strategy. $\text{Epsilon} = \frac{\text{Exploration}}{\text{Exploration} + \text{Exploit}}$

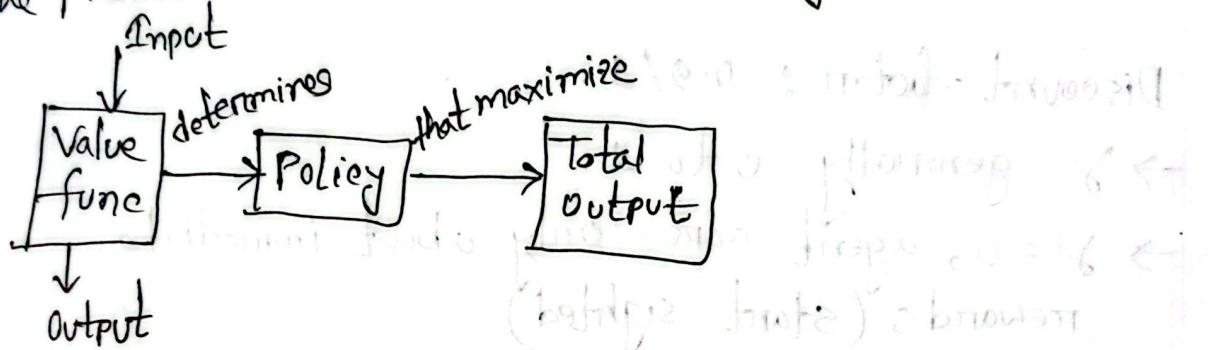
Reinforcement learning:

- System, at first explores then, $\epsilon = 1$.
- Gradually, likelihood of the system decreases & it knows enough about the environment.
- It takes to exploit the possible cases which one gives the highest accuracy.
- Based on its prediction, a penalty or reward is given & it tries to further improve itself.
- Each time, it generates a random number. If it is greater than the current ϵ value, agent will choose exploitation.
- Update the state using Bellman eqn so that Q will eventually optimize.



Deep Q Learning :

→ Value based reinforcement learning method to solve the problem.



It handles both exploration & exploitation.

↳ Exploration → Variety situation knowledge.

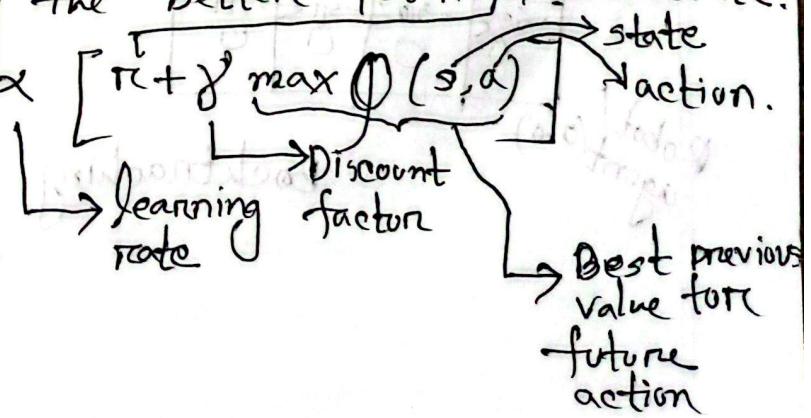
Exploitation → Better Accuracy.

Both are needed, otherwise better option might miss.

It tries to learn something new even if it seems bad before discovering the better path. → Immediate.

$$Q(s,a) = Q(s,a) + \alpha [r + \gamma \max Q(s',a')]$$

↖ Previous state of that situation



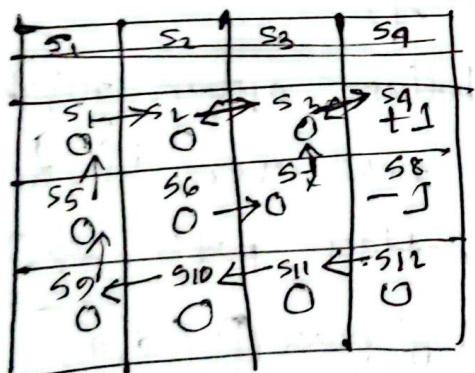
Learning rate

- Generally 0 to 1
- Small $\alpha \rightarrow$ Update are tiny, learning is slow but stable
- Large $\alpha \rightarrow$ Update are big, fasten training but can be noisy.
- Practically, we start with higher α & reduce it over time.

Discount factor: 0.9/1

- γ generally 0 to 1
- $\gamma = 0$, agent cares only about immediate rewards (short sighted)
- $\gamma = 1$, agent cares long-term rewards (plans ahead)

Usually starts with higher value.



Robot/
agent (0,0)

Backtracking

$$v(s_3) = 0 + \gamma [-1 + 0.9 \times 0] = -1$$

$$v(s_2) = 0 + \gamma [0 + 0.9 \times 1] = 1 \times 0.9 = 0.9$$

$$v(s_1) = 0 + \gamma [0 + 0.9 \times 0.9] = 0.81$$

$$v(s_5) = 0 + \gamma [0 + 0.9 \times 0.81] = 0.73$$

$$v(s_9) = 0 + \gamma [0 + 0.9 \times 0.73] = 0.66$$

$$v(s_{10}) = 0 + 0.9 \times 0.66 = 0.59$$

$$v(s_{11}) = 0 + 0.9 \times 0.59 = 0.53$$

$$v(s_{12}) = 0 + 0.9 \times 0.53 = 0.43$$

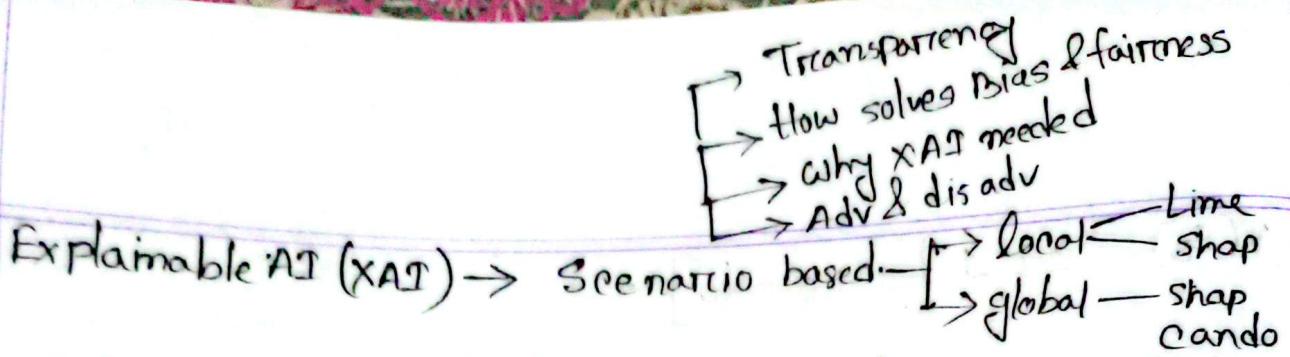
$$v(s_7) = 0 + 0.9 \times 1 = 0.9 \rightarrow \text{New } (\alpha, s, a)$$

$$\alpha \sim (s, a) = 0 + 0.9 \times 0.9 = 0.81$$

$\alpha = 0$ (Scenario) → Same path

→ Iteration fixed.

A robot in a 5×5 grid must learn to reach a goal $(5, 5)$ while avoiding obstacles. Explain how α learning updates its α table over time & what role do α & γ play in this.



→ Refers to AI method's decision making process.

* ~~Process~~ SHAP (shapely Additive Explanation).

→ Based on game theory.

→ Shows how much each feature adds & subtracts from the model's final prediction.

SHAP has added advantage that it can compute global feature importance by averaging.

SHAP values accuracy across many samples.

* LIME (Local Interpretable Model Agnostic Explanation).

→ Builds a semi simple, local model around each individual prediction.

- ① Sampling & obtaining a surrogate dataset.
- ② feature selection } Presented the output
- ③ Ridge Regression model } through a sample.

XAI help to

- ① Build Trust - Users see the reasoning behind prediction.
- ② Ensure fairness - Detect & Locate correct bias, reveals the abnormal behaviour.
- ③ Improve Accountability - Regulators & humans verify AI behaviour.

(IV) Support debugging: fix wrong logic.

- Transformer:

- ① Attention: how model knows which were to attend to is all learned while training with back propagation. Attention mechanism has an infinite references window.
- ② encoder: Maps an input sequence into an abstract continuous representation that holds learned info of that input.
- ③ Decoder: Takes that continuous representation & step by step generate a single output, while are also being fed to previous input.

How Encoder (BERT, RoBERT, ALBERT) works?

- ① Input embedding - Vector representation of given information.
- ② Inject positional information through encoding (using sine & cosine func).
- ③ Encoder layers (Multiple layers) - To map all input sequence into an abstract continuous representation that holds the learned info for that entire sequence. It contains 2 submodule.
 - ④ output goes to a linear layer, each output are concatenated to a single vector.

① Multihead Attention → We fed 3 distinct fully connected layers to create query key & value vectors.
At first, query \times key = score.

The higher score has higher priority.
Then, $\text{softmax} \left(\frac{\text{score}}{\sqrt{\text{dimension}}} \right) \times \text{value} = \text{output}$.

④ Residual Connection, Layer Normalization & Pointwise FFNN.

Multitread all output + original input goes through these layers & final output are again fed to input. Continuous further for reaching representation.

→ How decoder (GPT works)?

① Input embedding.

② Inject positional info

③ Decoder layer

The main difference lies in multihead attention. Hence, it can access the current position with all previous info. we cover the future words & gives it to predict.

This is called Masking.

Score + Look Ahead Mask = Masked score

[\rightarrow for previous works, α -for future works]

here, when goes to

sigmoid func, the future values become 0, then it continues & predict 1. output.

(iv) In second layer, final output & all input are passed and continued.

Math (25/30)

Ensemble - Compare

feature selection

XAT, Trax, VAE, GAN

NN, Gradient Descent.

10/15



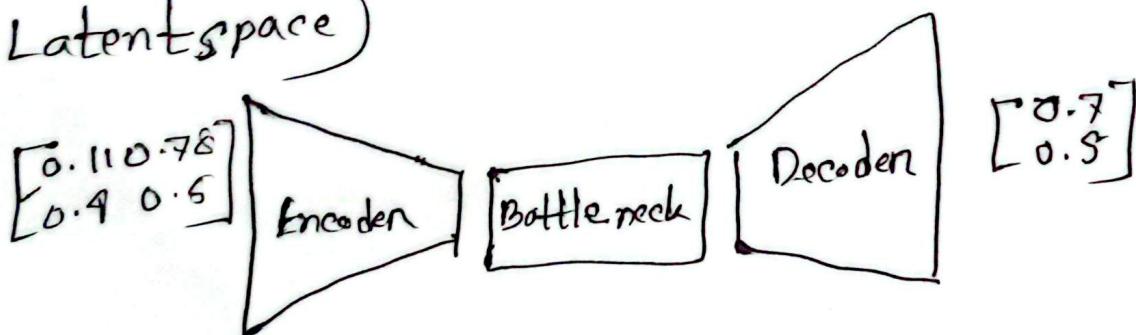
Generative Adversarial Network (GAN)

- takes random noise as input & tries
- ① Generator → Takes random noise as input & tries to create fake data.
 - ② Discriminator → Receives real data, fake data from the generator & learns to distinguish bet'n real and fake data.
 - ③ Training is like competition. Generator gets better at producing realistic data. The discriminator can no longer tell real from fake. Thus GAN converges.

Variational Auto Encoder (VAE)

→ Unlike normal encoders, that outputs a fixed latent vector, a variational encoder outputs a distribution using mean (μ) & variance (σ^2).

Basic concept of auto encoders is based on 3 component — Encoder, Decoder, Bottleneck (Latent space)



- ① Encoder compresses input data into latent space representation. Latent space is low dimensional space that captures the essential features of the input data.
- ② Bottleneck layer holds the compressed representation
- ③ Decoder reconstructs the input data from compressed representation produced by encoder & bottleneck.

Minimize the diff betⁿ original & reconstructed data.

Improves decoders ability to accurately reconstruct the original data from its compressed representation. At the same time, Encoder becomes good at compressing data in a way that preserves critical info.

The diff of 2 images are measured in pixel value. Loss func $\rightarrow \frac{1}{n} \sum_{i=1}^n (x_i - y_i)^2$.

Try to reduce mean square error for each image it processes. At first, Reconstruction weight \rightarrow High, Auto Encoder weight \rightarrow Random. As training process, error decrease. Auto Encoder gets better at recreating.

* VAE learns a probabilistic latent space & reconstructs data using reconstruction loss, while GAN learns to generate data through Generator Discriminator adversarial network.

* VAE - Output smooth but slightly blurry samples, GAN - More sharper, realistic, image.

* VAE - Anomaly detection, feature reduction

GAN - Produce New data.

* GAN - is harder to train.

Hyper parameter tuning:

Grid Search: Tries all combination of hyper parameters. Bayesian

Bayesian Optimization: Instead of trying all combinations, it predicts which parameters are promising.

Best for example problems, when search space is large.

Transfer learning

↳ federated learning

→ Model Train
→ Data Secure

✓ কোর্সের কাছ
✓ কামনা প্রয়োজনীয় Use

- * গোলি Activation func ব্যবহার করা হয় } NN
- * Gradient Descent

{ Ensemble
feature selection
XAI, GAN, VAE,
Gradient descent / perceptron.