

Q1a) Reinforcement Learning (Experiential learning)

- * Here the agent learns how to behave in an environment by performing actions and experiencing the result.

TYPES

Episodic Learning

Here we have starting & ending state

Thus an episode is created

→ Here the further action is based on feedback of result of earlier action.

Ex → Fear of dog after being bitten is Episodic learning

Continuous Task Learning

- * There is no terminal state here
- * Agent that does Automated stock tracking goes for Continuous learning.

Q1b)

Clustering

① Model will identify the pattern in data and will create clusters

Ex → Above given

Market Segment

Association

- ① We find dependencies of one data item on another.
- ② This dependency will help to map the relation and will enhance prediction
- ③ It's like "If - Then"

Clustering

- ① Model will identify the pattern in data and will create clusters

Ex → Above given

Market Segmentⁿ

Created Clusters of
high Sale }
Medium Sale }
Low Sale }
=====

Association

- ① We find dependencies of one data item on another.
② This dependency will help to map the relation and will enhance prediction
③ It's like "If - Then"

Ex Market Basket Analysis
Used by Amazon

If a person purchases Cell phone Then
the person has tendency to purchase
Screen guard & battery cover.

(Q1C)

Karl Pearson's Coefficient of Correlation

$$r = \frac{N \sum xy - \sum x \sum y}{\sqrt{[N \sum x^2 - (\sum x)^2] [N \sum y^2 - (\sum y)^2]}}$$

$r \Rightarrow$ quantifies the strength of relationship between two variables

The value of r will be betⁿ $\underline{\underline{+1}}$ and $\underline{\underline{-1}}$.

If $+1 \Rightarrow$ Total +ve correlation if $x \uparrow$ then $y \uparrow$

If $-1 \Rightarrow$ Total -ve correlation if $x \downarrow$, then $y \downarrow$

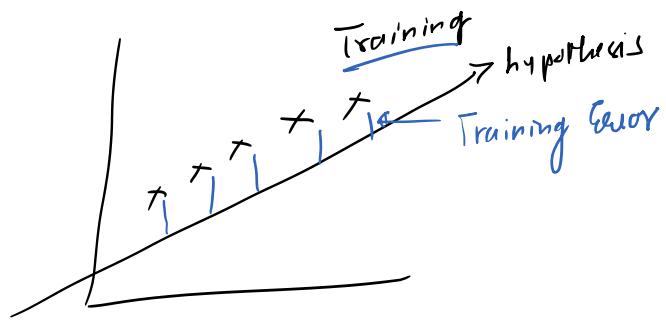
If $0 \Rightarrow$ No correlation.

\rightarrow It summarizes the degree & direction of correlation.

Old

① Training Error \Rightarrow

It is prediction error that we get when we apply the model on same data from which we trained.



$$E_{\text{Train}} = \frac{1}{n} \sum_{i=1}^n \text{Error}(\underline{f_D(x_i)}, y_i)$$

↑
prediction
on x_i ↑
actual
value

(for all samples)

② Test Error \Rightarrow It is prediction error we get when we apply model on altogether different data set (Test set) and not on the data it is trained.

$$E_{\text{test}} = \frac{1}{n} \sum_{i=1}^n \text{Error}(\underline{f_D(x_i)}, \underline{y_i})$$

over test set

① e) Performance Metric.

Consider a Skewed Data Set (One Sided)

In Such Cases to Evaluate Performance \Rightarrow

Confusion Matrix

		Actual	
		True	False
Predicted	True	True +ve (TP)	False +ve (FP)
	False	False -ve (FN)	True -ve (TN)

For 2 classes
True / False

Ex

		Actual	
		Has Heart Disease	Do not have Heart Disease
Predicted	Has Heart Disease	TP	FP
	Do not have Heart Disease	FN	T_N

* If we have more than Two classes.

Actual

	Chakde	KGF	DDLJ
Chakde	TP		
KGF		TP	
DDLJ			TP

- * Size of confusion matrix is determined by no of things we want to predict.
- * Confusion matrix summarizes what ML Algo did right and what it did wrong.

Tone +ve is evident easily

But for Tone -ve, FP, FN it depends on Question Asked -

	Actual			
	Chakde	KGF	DDLJ	
Predicted	Chakde	TP	FP	FP
	KGF	FN	TN	TN
	DDLJ	FN	TN	TN

Annotations:

- Red circles highlight TP, FP, and FN cells.
- A red arrow points from the 'Predicted Not watched Chakde' annotation to the FN cell in the KGF row.
- A blue circle highlights the TN cells in the DDLJ row.
- A blue arrow points from the 'not Actually watched Chakde But predicted to watch' annotation to the FP cell in the DDLJ row.
- A blue arrow points from the 'If question for Chakde' annotation to the FP cell in the DDLJ row.
- A blue arrow points from the 'if person saw Chakde' annotation to the TN cell in the DDLJ row.
- A blue arrow points from the 'not Actually watched Chakde & even not predicted Not watch' annotation to the TN cell in the DDLJ row.

Oif >

(1) Sensitivity \Rightarrow (TPR)

Ex \Rightarrow It tells us what percentage of patient with Heart Disease are correctly identified

$$= \frac{\text{Correctly Identified}}{\text{Total No of Patients with Heart Disease}} = \frac{TP}{\text{Actual +ve}}$$

(2) Specificity = Tells us what percentage of patients without Heart Disease were correctly identified.

$$= \frac{\text{Correctly Identified}}{\text{Total No of Patients without Heart Disease}} = \frac{TN}{\text{Actual -ve}}$$

FPR = $1 - \text{Specificity} = 1 - \frac{TN}{TN + FP} = \frac{TN + FP - TN}{TN + FP} = \frac{FP}{TN + FP}$ ✓

False +ve
Ratio

Ex

		Actual		
		Chakde	KLF	DDLJ
Predicted	Chakde	TP 12	102	93
	KLF	FN 112	23	77
	DDLJ	83	92	17

FP

TN

Calculate Sensitivity & Specificity for Chakde.

$$\text{Sensitivity} = \frac{\text{Correct +ve}}{n} = \frac{TP}{TP + FN} = \frac{12}{12 + 112 + 83} = \underline{\underline{0.06}}$$

$$\text{Sensitivity} = \frac{\text{Correct +ve}}{\text{Total +ve}} = \frac{TP}{TP+FN} = \frac{12}{12+112+83} = \underline{\underline{0.12}}$$

$$\text{Specificity} = \frac{\text{Correct -ve}}{\text{Total -ve}} = \frac{TN}{TN+FP} = \frac{23+77+92+17}{23+77+92+17+102+93} = \underline{\underline{0.52}}$$

Q2a)
Linear Regression Numerical Ex with multiple independent features (Two d).

[For > 2 features we will need Matrix Algebra.]

Consider

Sr.ND	x_1	x_2	y
1	3	8	-3.7
2	4	5	3.5
3	5	7	2.5
4	6	3	11.5
5	2	1	5.7
6	3	2	2.4

given

$$y = \underline{\underline{0}_0} + \underline{\underline{0}_1} x_1 + \underline{\underline{0}_2} x_2$$

Independent features ↑ Dependent on x_1 & x_2 .

Linear Reg for 2 Independent Variable.

$$\underline{\underline{0}_0} = \bar{y} - \underline{\underline{0}_1} \bar{x}_1 - \underline{\underline{0}_2} \bar{x}_2$$

$$\underline{\underline{0}_1} = \frac{(\sum x_1^2)(\sum x_1 y) - (\sum x_1 x_2)(\sum x_2 y)}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2}$$

$$\underline{\underline{0}_2} = \frac{(\sum x_1^2)(\sum x_2 y) - (\sum x_1 x_2)(\sum x_1 y)}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2}$$

$$\text{Here } \sum x_1^2 = \underline{\underline{\sum x_1 x_1}} - \frac{(\sum x_1)(\sum x_1)}{N} = \sum (x_1)^2 - \frac{(\sum x_1)^2}{N}$$

$$\sum x_2^2 = \underline{\underline{\sum x_2 x_2}} - \frac{(\sum x_2)(\sum x_2)}{N}$$

$$\sum_{\text{S.No}} xy = \sum xy - \frac{(\sum x_1)(\sum y)}{N}$$

$$\sum x_2 y = \sum x_2 y - \frac{(\sum x_2)(\sum y)}{N}$$

$$\sum x_1 x_2 = \sum x_1 x_2 - \frac{(\sum x_1)(\sum x_2)}{N}$$

S.No	y	x_1	x_2	$(x_1)^2$	$(x_2)^2$	$x_1 y$	$x_2 y$	$x_1 x_2$
1	-3.7	3	8	9	64	-11.1	-29.6	24
2	8.5	4	5	16	25	14	17.5	20
3	2.5	5	7	25	49	12.5	17.5	35
4	11.5	6	3	36	9	69	34.5	18
5	5.7	2	1	4	1	11.4	5.7	2
Σ	19.5	20	24	90	148	95.8	45.6	99.

$$\sum x_1^2 = \sum (x_1)^2 - \frac{(\sum x_1)^2}{N} = 90 - \frac{(20)^2}{5} = 90 - 80 = \underline{\underline{10}}$$

$$\sum x_2^2 = 148 - \frac{(24)^2}{5} = \underline{\underline{32.8}}$$

$$\sum x_1 y = 95.8 - \frac{20 * 19.5}{5} = 17.8$$

$$\sum x_2 y = 45.6 - \frac{24 * 19.5}{5} = -48$$

$$\sum x_1 x_2 = 99 - \frac{20 * 24}{5} = 99 - 96 = 3$$

$$B_1 = \frac{(32.8)(17.8) - (3)(-48)}{(10)(32.8) - (3)^2} = \underline{\underline{2.28}}$$

$$Q_2 = \frac{(10)(-48) - (3)(17.8)}{(10)(32.8) - (3)^2} = \underline{\underline{-1.67}}$$

$$Q_0 = \bar{y} - Q_1 \bar{x}_1 - Q_2 \bar{x}_2 = \frac{19.5}{5} - \frac{2.28 \times 20}{5} - \frac{(-1.67)(24)}{5}$$

$$= 2.796$$

$\therefore \hat{y} = 2.796 + 2.28x_1 - 1.67x_2$

Substitute $x_1 = 3$ $x_2 = 2$ $\hat{y} = \underline{\underline{2.796}}$

Q2 b)

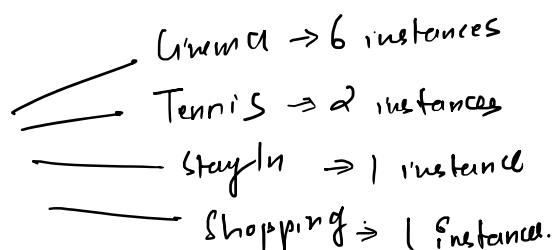
Create Decision Tree for following Using Gini Index
 (Classification & Regression Tree) \rightarrow CART.

Weekend	Weather	Parent	Money	Decision	(y)
w ₁	Sunny	Yes.	Rich	Cinema	-
w ₂	Sunny	No	Rich	Tennis	-
w ₃	Windy	Yes.	Rich	Cinema	-
w ₄	Rainy	Yes.	Poor	Cinema	.
w ₅	Rainy	No	Rich	Stay In	
w ₆	Rainy	Yes.	Poor	Cinema	.
w ₇	Windy	No	Poor	Cinema	.
w ₈	Windy	No	Rich	Shopping	
w ₉	Windy	Yes	Rich	Cinema	
w ₁₀	Sunny	No	Rich.	Tennis	-

Soln

Calculate Gini Index for Overall collection of Outcomes
 of Training Example

There are 4 possible outcome Variable



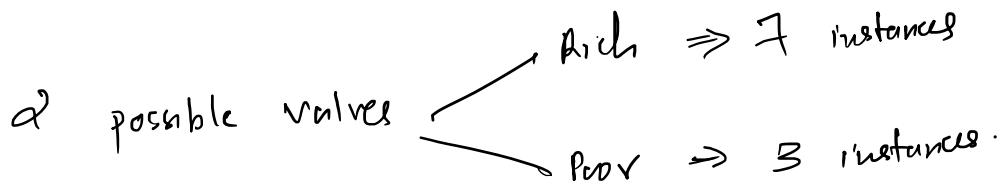
$$\begin{aligned}
 \text{Gini (Decision)} &= 1 - \left[\left(\frac{6}{10} \right)^2 + \left(\frac{2}{10} \right)^2 + \left(\frac{1}{10} \right)^2 + \left(\frac{1}{10} \right)^2 \right] \\
 &= \underset{\text{Cinema}}{\uparrow} \quad \underset{\text{Tennis}}{\mid} \quad \underset{\text{StayIn}}{\uparrow} \quad \underset{\text{Shopping}}{\uparrow}
 \end{aligned}$$

$$= 0.58.$$

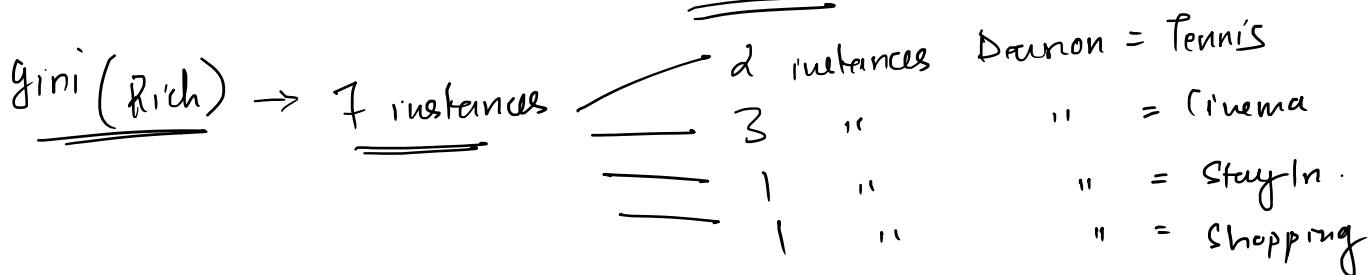
* We will calculate Gini Index for every attribute and based on value of gini for Attribute we will select root and start constructing Tree.

①

Let find Gini Index for Money feature.



We Need to find Gini for Both Rich & Poor



$$\begin{aligned} \text{Gini}(\text{Rich}) &= 1 - \left(\left(\frac{2}{7} \right)^2 + \left(\frac{3}{7} \right)^2 + \left(\frac{1}{7} \right)^2 + \left(\frac{1}{7} \right)^2 \right) \\ &= \underline{0.694} \end{aligned}$$

Gini(Poor) \rightarrow 3 instances \rightarrow Decision is Cinema.

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

$$= 1 - 1 = \underline{0}$$

Weighted Avg Gini for Money

\nwarrow \nearrow \downarrow \uparrow \dots Info available.

Weighted Avg Gini for Money

= (Gini Index * Proportion) for each possible values.

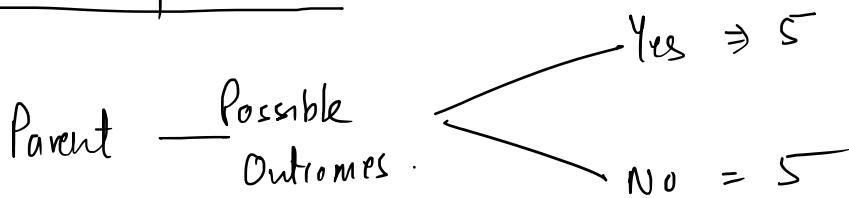
$$= \left(0.694 * \frac{7}{10} \right) + \left(0 * \frac{3}{10} \right)$$

Rich. Poor

$$= 0.486$$

$$\boxed{\text{Gini (Money)} = 0.486}$$

② Gini Index for Parent.



$$\text{Gini (Parent = Yes)} = 5 \text{ instances} \xrightarrow[\text{Possible Decision}]{} 5 \text{ Cinema}.$$

$$\text{Gini (Parent = No)} = 1 - [(5/5)^2] = 0.$$



$$\begin{aligned} \text{Gini (Parent = No)} &= 1 - [(2/5)^2 + (1/5)^2 + (1/5)^2 + (1/5)^2] \\ &= 0.72 \end{aligned}$$

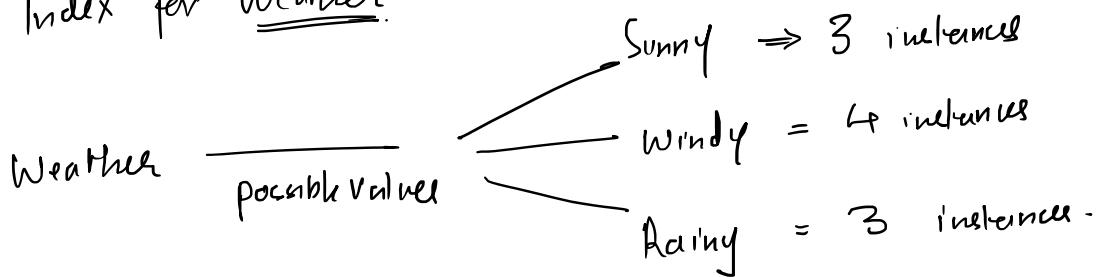
$$= \underline{0.72}$$

$$\text{Weighted Avg Gini for Parent} = \left(0 * \frac{5}{10}\right) + \left(0.72 * \frac{5}{10}\right)$$

$$= \underline{\underline{0.36}}$$

$$\therefore \boxed{G_{\text{ini}}(\text{Parent}) = 0.36}$$

③ Gini Index for Weather.



$$G_{\text{ini}}(\text{Weather} = \text{Sunny}) \xrightarrow{3 \text{ instances}} \begin{array}{l} 1 \text{ Cinema} \\ 2 \text{ Tennis} \end{array}$$

$$G_{\text{ini}}(\text{Weather} = \text{Sunny}) = 1 - \left(\left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 \right) = \underline{\underline{0.444}}$$

$$G_{\text{ini}}(\text{Weather} = \text{Windy}) = 4 \text{ instances} \xrightarrow{\begin{array}{l} 3 \text{ Cinema} \\ 1 \text{ Shopping} \end{array}}$$

$$G_{\text{ini}}(\text{Weather} = \text{Windy}) = 1 - \left(\left(\frac{3}{4}\right)^2 + \left(\frac{1}{4}\right)^2 \right) = \underline{\underline{0.375}}$$

$$G_{\text{ini}}(\text{Weather} = \text{Rainy}) = 3 \text{ instances} \xrightarrow{\begin{array}{l} 2 \text{ Cinema} \\ 1 \text{ StayIn} \end{array}}$$

$$\therefore G_{\text{ini}}(\text{Weather} = \text{Rainy}) = 1 - \left(\left(\frac{2}{3}\right)^2 + \left(\frac{1}{3}\right)^2 \right) = \underline{\underline{0.444}}$$

$$\text{weighted Avg Gini for weather} = \left(0.444 \times \frac{3}{10}\right) + \left(0.375 \times \frac{4}{10}\right) + \left(0.444 \times \frac{3}{10}\right)$$

Sunny Windy Rainy

$\boxed{\text{Gini (weather)} = \underline{\underline{0.416}}}$

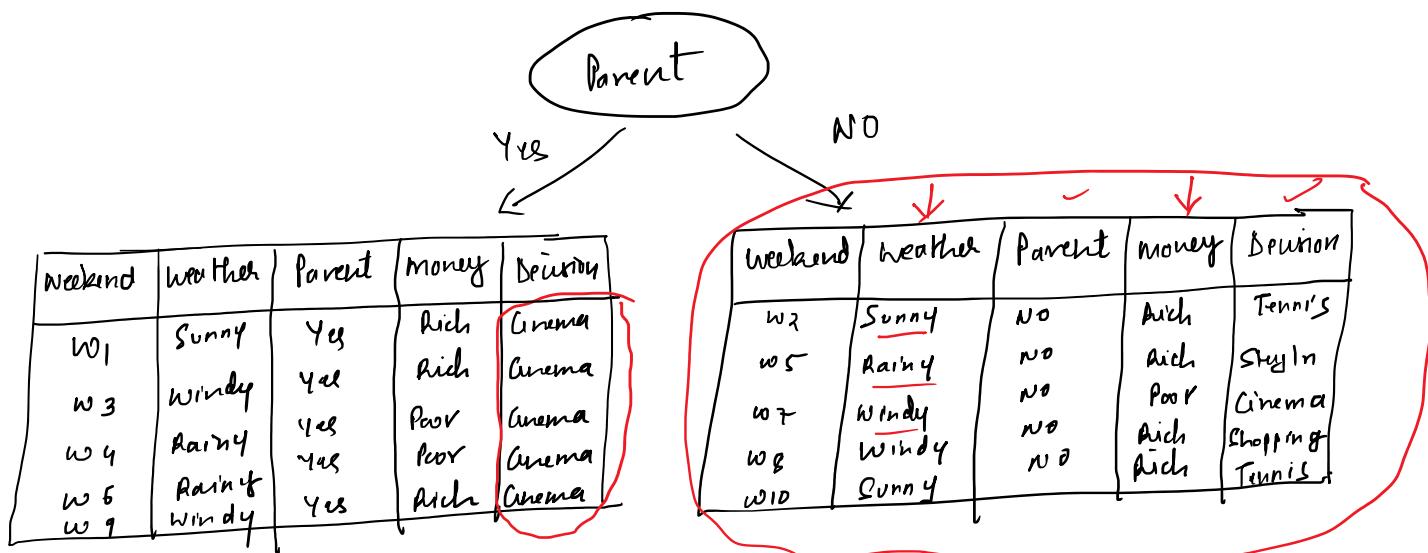
$$\therefore \text{Gini (weather)} = 0.416$$

$$\text{Gini (Parent)} = 0.36$$

$$\text{Gini (Money)} = 0.486.$$

Gini Index for Parent is Minimum So it is said to have maximum information.

So Select Parent as Root Node of Decision Tree



* Computation of Gini Index for Parent = Yes Not needed.
as In All cases Decision = Cinema

* For parent = No we do not have single Decision
So again we have to calculate Gini for.

$$\underline{\text{gini} (\text{Parent} = \text{NO} \& \text{Weather})}$$

possible
value

Sunny - 2 instance
Windy - 2 instance
Rainy - 1 instance

$$\text{gini} (\text{Parent} = \text{NO} \& \text{Weather} = \text{Sunny}) \Rightarrow 2 \text{ instance} \leftarrow ^2 \text{ tennis}$$

$$= 1 - ((\frac{1}{2})^2) = 0.$$

$$\text{gini} (\text{Parent} = \text{NO} \& \text{Weather} = \text{Windy}) = 2 \text{ instance} \leftarrow ^1 \text{ Cinema} \\ ^1 \text{ Shopping}$$

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

$$\text{gini} (\text{Parent} = \text{NO} \& \text{Weather} = \text{Rainy}) = 1 \text{ instance} \rightarrow 1 \text{ Stayin} \\ = 1 - [(\frac{1}{1})^2] = 0.$$

$$\therefore \text{Avg weighted gini} (\underline{\text{Parent} = \text{NO} \& \text{Weather}}) = \left(0 \times \frac{2}{5} \right) + \left(0.5 \times \frac{2}{5} \right) + \left(0 \times \frac{1}{5} \right) \\ = \underline{\underline{0.2}}$$

Parent

Weekend	Weather	Parent	Money	Decision
w1	Sunny	Yes	Rich	Cinema
w3	Windy	Yes	Rich	Cinema
w4	Rainy	Yes	Poor	Cinema
w6	Rainy	Yes	Poor	Cinema
w9	Windy	Yes	Rich	Cinema

Weekend	Weather	Parent	Money	Decision
w2	Sunny	No	Rich	Tennis
w5	Rainy	No	Rich	StayIn
w7	Windy	No	Poor	Cinema
w8	Windy	No	Rich	Shopping
w10	Sunny	No	Rich	Tennis

Gini Index Parent = NO and Money:

- Rich → 4 instances
- Poor → 1 instances

$$\text{Gini Index} (\text{Parent} = \text{NO} \& \text{Money} = \text{Rich}) = \frac{2}{5} \text{ Tennis} + \frac{1}{5} \text{ StayIn} + \frac{1}{5} \text{ Shopping}$$

$$\text{Gini Index} (\text{Parent} = \text{NO} \& \text{Money} = \text{Rich}) = 1 - \left(\frac{(2/5)^2 + (1/5)^1 + (1/5)^2}{5} \right) = 0.625$$

$$\text{Gini Index} (\text{Parent} = \text{NO} \& \text{Money} = \text{Poor}) = 1 \text{ instance Cinema} \\ = 1 - ((1/1)^2) = 0$$

Weighted Average Gini Index for (Parent = NO & Money) =

$$= 0.625 \times \frac{4}{5} + 0 \times \frac{1}{5}$$

and

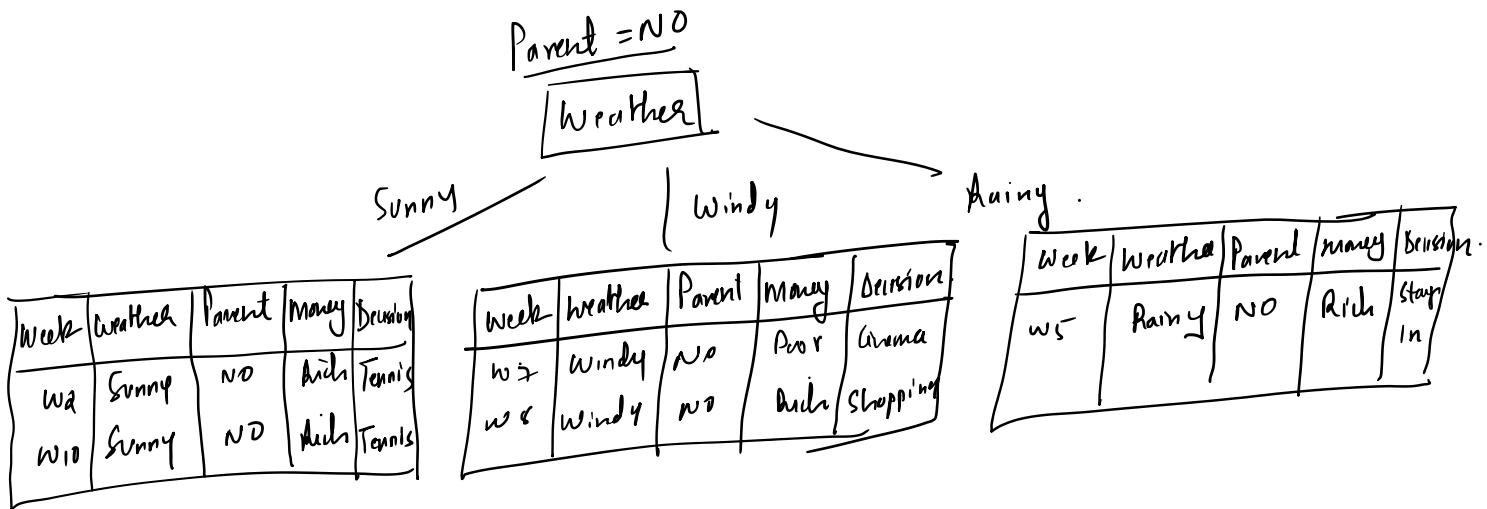
$$= \underline{\underline{0.5}}$$

$$\therefore \text{Gini Index} (\text{Parent} = \text{NO} \& \text{weather}) = 0.2$$

$$\text{Gini Index} (\text{Parent} = \text{NO} \& \text{money}) = 0.65$$

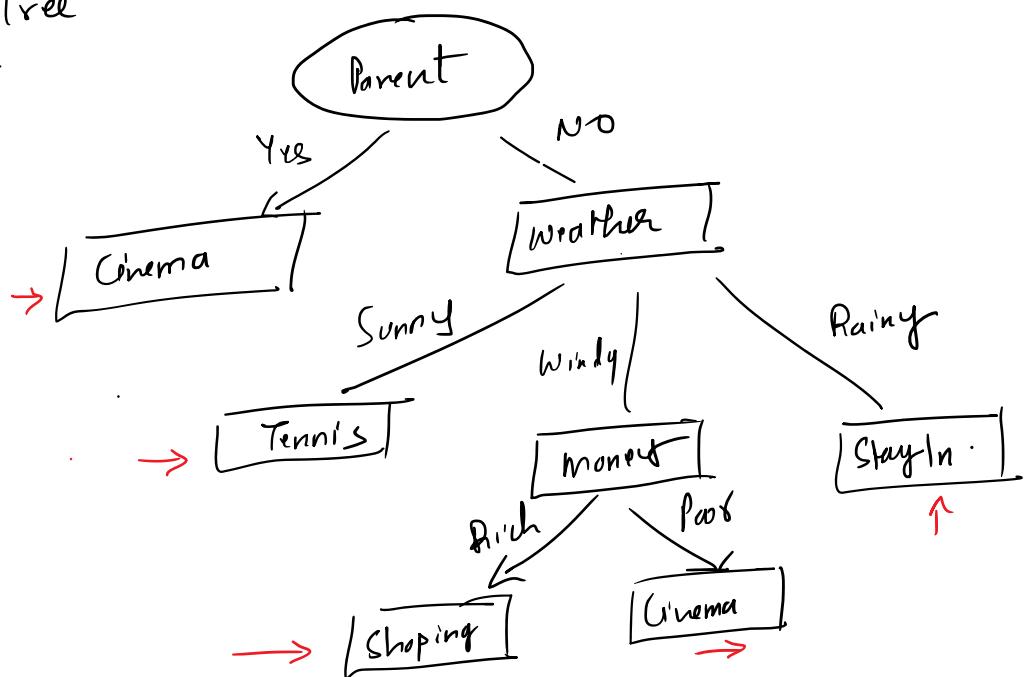
Here Weather Attribute has smaller Gini Index

Select Weather as Node.



Decided Decision Tree

For Given Dataset



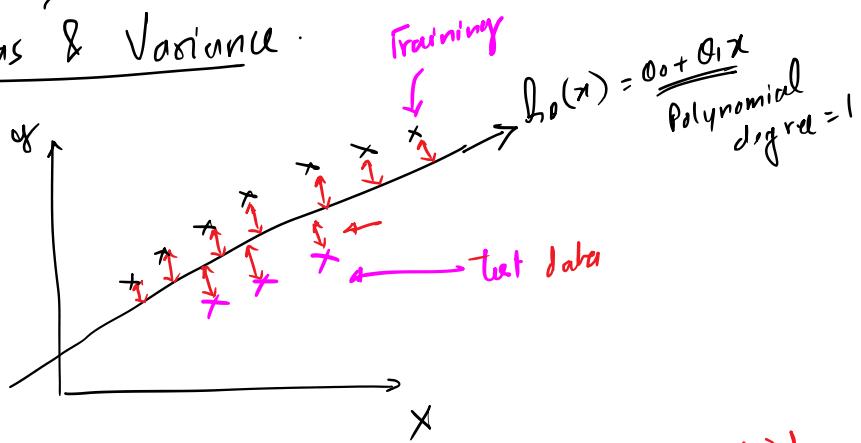
MW

Construct an Optimal Decision Tree for following.

outlook	Temperature	Humidity	Wind	Play (Decision)
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Mild	High	F	Y
Rainy	Mild	High	F	Y
Rainy	Cool	Normal	F	Y
Rainy	Cool	Normal	T	No
overcast	Cool	Normal	T	Y
Sunny	Mild	High	F	No
Sunny	Cool	Normal	F	No
Rainy	Mild	Normal	F	Y
Sunny	Mild	Normal	T	Y
overcast	Mild	High	T	Y
overcast	Hot	Normal	F	Y
Rainy	Mild	High	T	No

Q3a)

Bias & Variance



Normally
100% Training
70% Training
30% Test

Here line is of
lower degree Polynomial

We can have Higher Training Error \Rightarrow Bias
& Higher Test Error \Rightarrow High Variance.

Underfitting \Leftrightarrow Whatever data we have trained the model the error is quite high.

However It's Normally Found That Underfitting Model

has + Bias & ↓ Variance

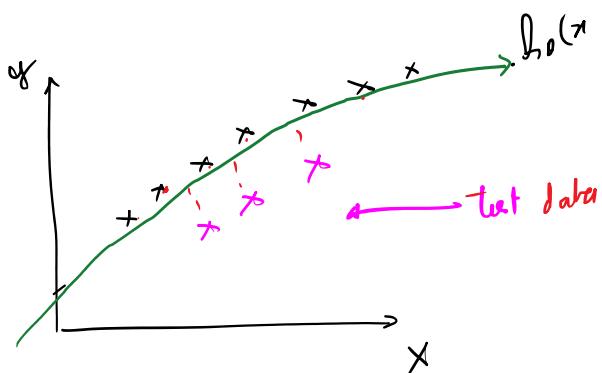
In Comparison to Earlier Model

degree of Polynomial 2

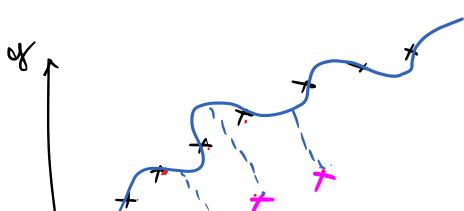
Here Training Error \downarrow
(Bias)

Test Error \downarrow
(Variance)

(as compared to overfit)

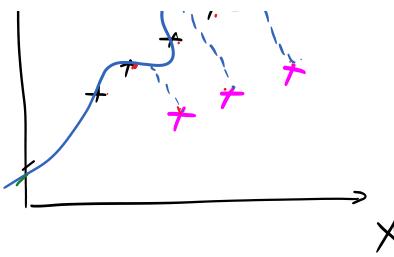


Higher degree Polynomial



Perfectly Fitting the Training Data

... \rightarrow Perfect Fit



Training Data
 Training Error (Bias) \rightarrow Very Low
 Test Error (Variance) \rightarrow Very High

Overfitting \rightarrow Here for Training
Each & Every point
is Satisfied by
Regression Line.

\rightarrow Here training data perfectly fits the line
 But Test Data may show error (High).

Note

Normally Degree of Polynomial low \rightarrow Underfit \rightarrow high Bias
low Variance

if Degree of Polynomial is very
high (complex \rightarrow Overfit) \rightarrow low Bias
High Variance

Training Error \Rightarrow Bias \rightarrow Represent Error when Regression line
fits the Training Data.

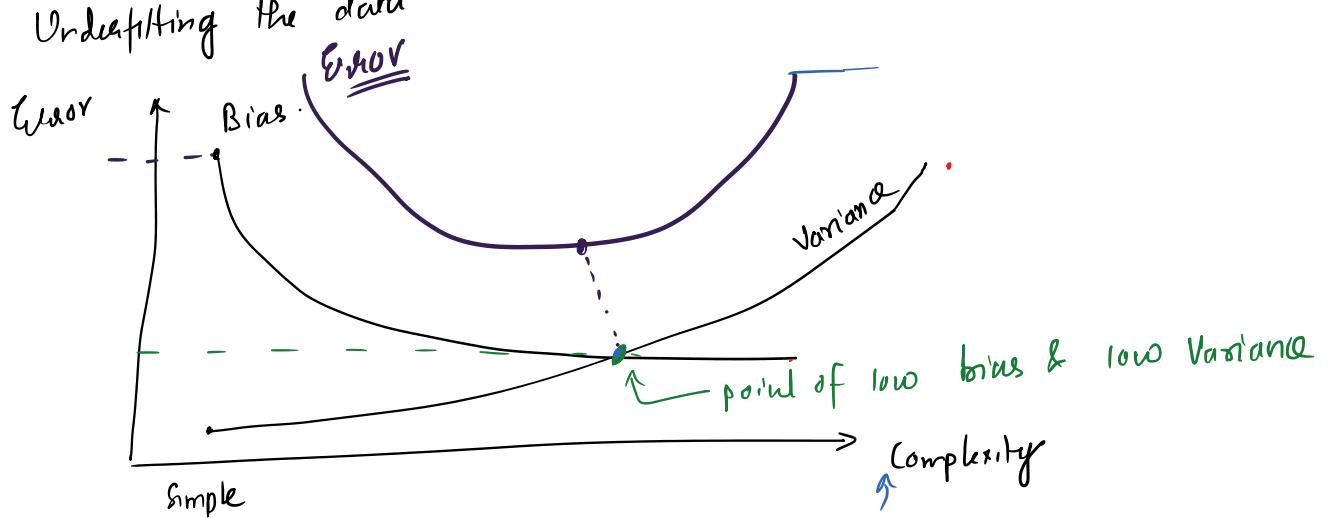
Test Error \rightarrow Variance = Represent Error when Regression line
(Cross Validation) fits the Test Data.

Bias - Variance Tradeoff

① If model is too simple and few parameters \rightarrow Underfit
 \rightarrow High Bias
Low Variance

② If model has large no of Parameters and
are complex → Overfit
→ High Variance
→ Low Bias

③ we need to find good balance without overfitting &
Underfitting the data.



Bias Variance Tradeoff Says
we need a model that gives

- (1) low Bias
- (2) low Variance

i.e. we need to find point of low Bias & low Variance.

* The method are:

* regularization (penalize θ (parameters)).

* Boosting

* Bagging

We need to minimize the total Error = $\underbrace{\text{Bias}^2 + \text{Variance}}_{\text{Irreducible Error}}$

Note

model 1

Training Error 1% \Rightarrow low Bias

} Overfitting
(complex model)

Now

Model 1

Training Error 1% \Rightarrow low Bias } Overfitting
 Test Error 20% \Rightarrow High Variance } (complex model)

Model 2

Training Error = 25% \Rightarrow High Bias } Underfitting
 Test Error = 26% \Rightarrow High Variance

Model 3

Training Error < 10% low Bias } Recommended
 Test Error < 10% low Variance }

Q3b)

Kappa Statistic →

- * Kappa statistic or Cohen's Kappa is a statistical measure of inter-rater reliability for categorical variables.
- * It is used when two raters apply a criteria based on a tool to assess whether or not some condition occurs.
- Eg Two doctors rate whether or not each of 20 patients has diabetes based on Symptoms
- * If two raters use same criteria on same target to evaluate and their agreement is very high then we will have evidence of reliable ratings.
- * If their agreement is not very high then
 - either criterion tool is not useful
 - or raters are not trained enough.
- * Kappa statistics corrects for chance agreement and not percent agreement.

		Rater A	
		Yes	No
Rater B	Yes	35	20
	No	15	40

→ 35 times they agree on Yes
40 times they agree on No
20 time Rater A is No But Rater B is Yes
15 time Rater A is Yes But Rater B is No.

Cohen Suggested following Statistics.

value ≤ 0 \rightarrow No agreement

$0.01 \rightarrow 0.20$ \rightarrow as none to slight

$0.21 \rightarrow 0.40$ \rightarrow as fair

$0.41 \rightarrow 0.60$ \rightarrow as moderate

$0.61 \rightarrow 0.80$ \rightarrow as substantial

$0.81 \rightarrow 1.00$ \rightarrow perfect agreement

* Rather than calculating the percentage of items raters agreed on, Cohen's kappa attempts to account the fact that raters may happen to agree on some items purely by chance.

Ex Two curators asked to rate 70 paintings

		C2	
		Yes	No
C1	Yes	25	10
	No	15	20

Step 1) calculate relative Agreement betⁿ Raters

$$P_o = \frac{\text{Both said Yes} + \text{Both said No}}{\text{Total}} = \frac{25 + 20}{70} = \underline{\underline{0.6429}}$$

Step 2) calculate hypothetical probabilities of chance agreement between raters

Pe?

$$\underline{\underline{P_e}} = \frac{C_1}{\text{Total Yes}} * \frac{C_2 (\text{Yes})}{\text{Total Yes}} = \left(\frac{25}{70} \right) * \left(\frac{25+10}{70} \right) = \underline{\underline{0.285714}}$$

$$P(N_0) = \frac{C_1(N_0)}{\text{Total Res}} + \frac{C_2(N_0)}{\text{Total Res}} = \left(\frac{15+20}{70} \right) * \left(\frac{20+10}{70} \right) = \underline{\underline{0.214285}}$$

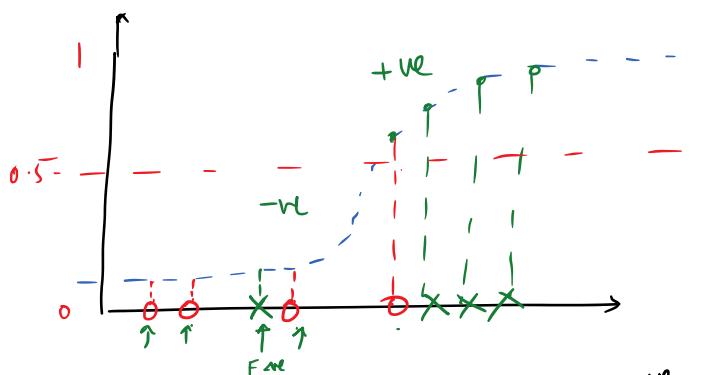
$$P_e = P(Y_{\text{Res}}) + P(N_0) = \underline{\underline{0.5}}$$

$$\text{Calculate Cohen's Kappa} = K = \frac{P_o - P_e}{1 - P_e} = \frac{(0.6429 - 0.5)}{1 - 0.5} = \underline{\underline{0.2857}}$$

It's in range $0.21 \rightarrow 0.40$ so the Agreement bet' two Curator is fair

Q3c) ROC & AUC →

ROC [Receiver Operator Characteristic]



for threshold 0.5

		+ve	-ve
+ve	3.	1	
	1.	3-	

$$\text{Sensitivity} = \frac{3}{4}$$

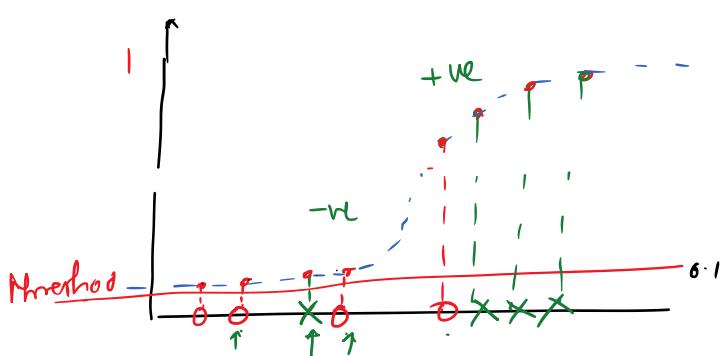
$$\text{Specificity} = \frac{3}{4}$$

$$\overline{\text{TPR}} = \text{Sensitivity} = \frac{4}{4} = 1$$

$$\text{Specificity} = \frac{0}{4} = 0 \quad \text{FPR} = 1 - 0 = 1$$

here if threshold 0.1

then the confusion matrix changes and according to the Sensitivity and Specificity changes -



- * Consider for logistic regression where we identified a threshold point and we prepared confusion matrix and calculated sensitivity & specificity.

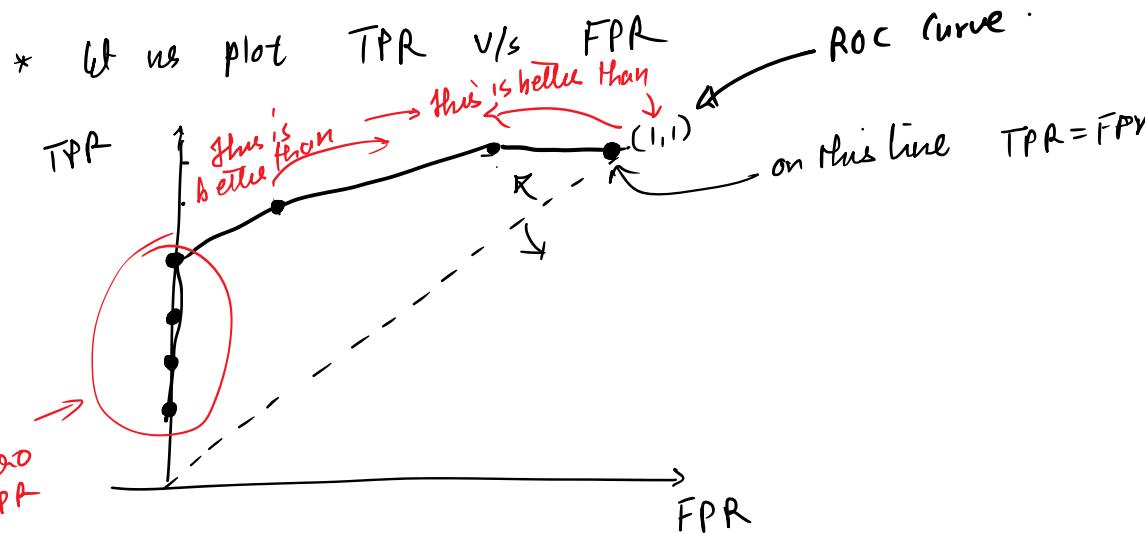
- * If threshold changes then the confusion matrix and accordingly the sensitivity and specificity changes.

- * We can have many such thresholds bet" 0 → 1 -

- * We want to analyse the performance at different threshold and want to identify the best of it

- * Instead of using so many confusion matrix we can use

ROC Curve



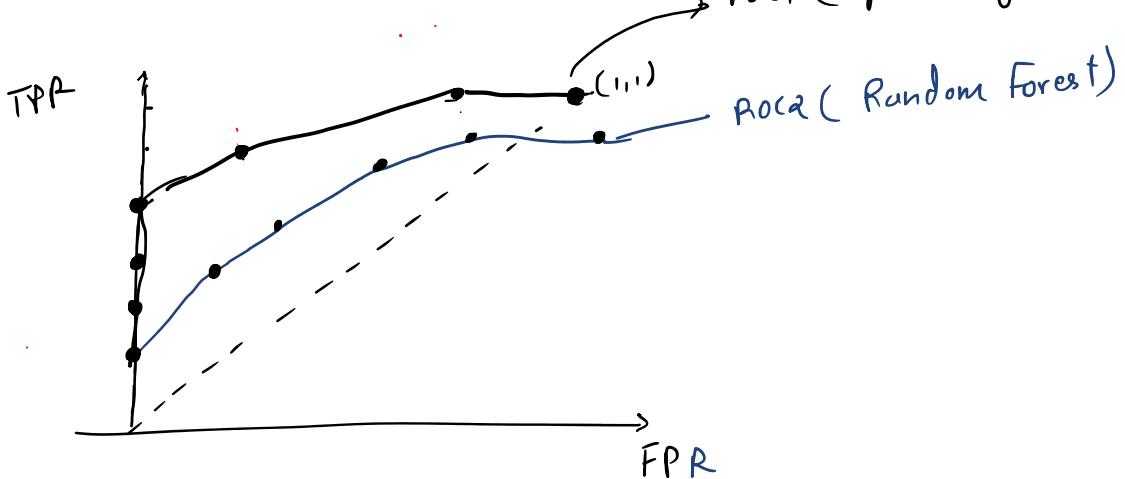
For diff threshold we calculated TPR and FPR and plotted accordingly \rightarrow in form of ROC curve.

\rightarrow From ROC curve we can directly determine which threshold is better.

AUC \Rightarrow Area Under Curve \Rightarrow It is a method to compare ROC for more than one method and will help us to judge which one is better.

ROC₁ (logistic Regression)

ROC₂ (Random Forest)



Here AUC for ROC₁ $>$ AUC for ROC₂

So method for ROC₁ (logistic Regression) is better than method for ROC₂ (Random Forest)