



**Vidyalankar Institute of Technology**  
**Semester 7 – CMPN - Mid Semester Assessment – 1**

|                  |                  |                  |
|------------------|------------------|------------------|
| Date: 05/08/2024 | Machine Learning | 30 Marks /1 hour |
|------------------|------------------|------------------|

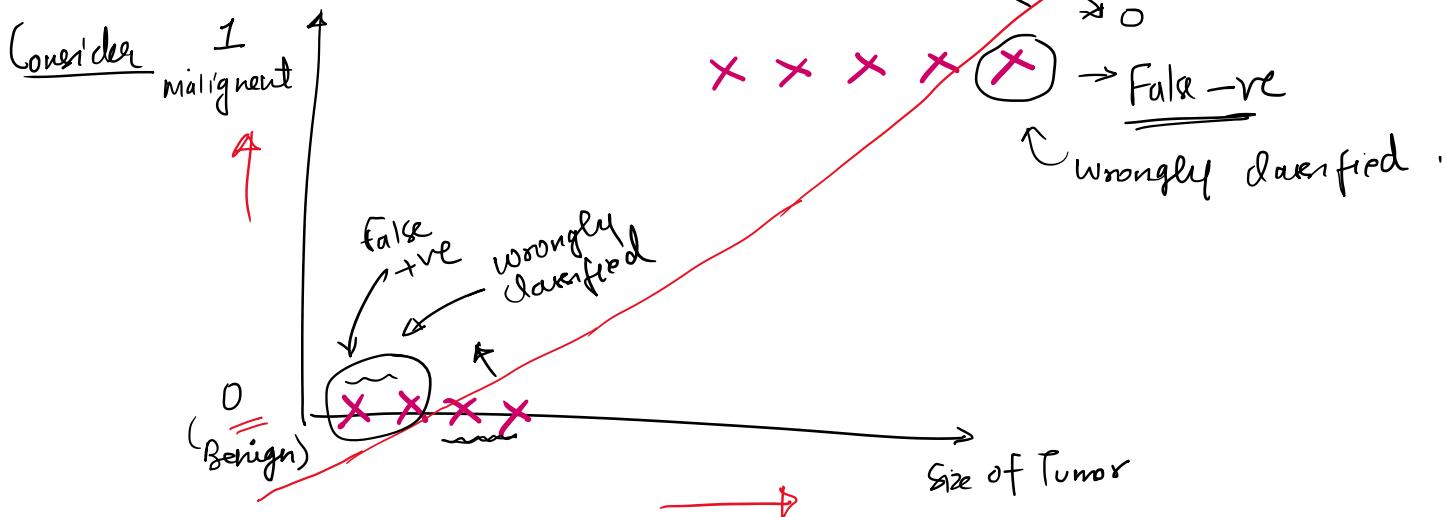
| <b>1</b> | <b>Solve any two (5 marks each)</b>  | <b>CO</b> |         |          |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
|----------|--|-----------|---------|----------|-------|----------|---|-------|-----|------|--------|---|-------|----|------|--------|---|-------|-----|------|--------|---|-------|-----|------|--------|---|-------|----|------|---------|---|-------|-----|------|--------|---|-------|----|------|--------|---|-------|----|------|----------|---|-------|-----|------|--------|----|-------|----|------|--------|--|
| A        | Justify why Linear Regression is not used for performing classification.   | CO1       |         |          |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
| B        | Elaborate on:<br>(i) Karl Pearson's Coefficient of Correlation.<br>(ii) R square method.   | CO2       |         |          |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
| C        | Discuss the importance of the Receiver Operating Characteristic (ROC) curve and the Area Under the Curve (AUC) in evaluating binary classifiers. How do these metrics help in comparing models?  | CO2       |         |          |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
| <b>2</b> | <b>Solve any two (5 marks each)</b>  |           |         |          |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
| A        | Outline the steps involved in developing a Machine Learning application for predicting housing prices.   | CO1       |         |          |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
| B        | Justify the significance of the F1 score in imbalanced classification problems. How does it address the limitations of using precision or recall alone? Provide a scenario where the F1 score might still be misleading.   | CO2       |         |          |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
| C        | Justify the bias-variance trade-off in the context of model complexity. How does this trade-off influence the choice of model in a real-world ML application? Provide examples to support your explanation.  | CO2       |         |          |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
| <b>3</b> | <b>Solve anyone (10 marks each)</b>  |           |         |          |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
| A        | For the dataset given below, construct a decision tree using Gini Index, and determine which attribute is a root attribute.  | CO2       |         |          |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
|          | <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th>Sr. No.</th> <th>Weather</th> <th>Parent</th> <th>Money</th> <th>Decision</th> </tr> </thead> <tbody> <tr><td>1</td><td>Sunny</td><td>Yes</td><td>Rich</td><td>Cinema</td></tr> <tr><td>2</td><td>Sunny</td><td>No</td><td>Rich</td><td>Tennis</td></tr> <tr><td>3</td><td>Windy</td><td>Yes</td><td>Rich</td><td>Cinema</td></tr> <tr><td>4</td><td>Rainy</td><td>Yes</td><td>Poor</td><td>Cinema</td></tr> <tr><td>5</td><td>Rainy</td><td>No</td><td>Rich</td><td>Stay In</td></tr> <tr><td>6</td><td>Rainy</td><td>Yes</td><td>Poor</td><td>Cinema</td></tr> <tr><td>7</td><td>Windy</td><td>No</td><td>Poor</td><td>Cinema</td></tr> <tr><td>8</td><td>Windy</td><td>No</td><td>Rich</td><td>Shopping</td></tr> <tr><td>9</td><td>Windy</td><td>Yes</td><td>Rich</td><td>Cinema</td></tr> <tr><td>10</td><td>Sunny</td><td>No</td><td>Rich</td><td>Tennis</td></tr> </tbody> </table> | Sr. No.   | Weather | Parent   | Money | Decision | 1 | Sunny | Yes | Rich | Cinema | 2 | Sunny | No | Rich | Tennis | 3 | Windy | Yes | Rich | Cinema | 4 | Rainy | Yes | Poor | Cinema | 5 | Rainy | No | Rich | Stay In | 6 | Rainy | Yes | Poor | Cinema | 7 | Windy | No | Poor | Cinema | 8 | Windy | No | Rich | Shopping | 9 | Windy | Yes | Rich | Cinema | 10 | Sunny | No | Rich | Tennis |  |
| Sr. No.  | Weather  | Parent    | Money   | Decision |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
| 1        | Sunny  | Yes       | Rich    | Cinema   |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
| 2        | Sunny  | No        | Rich    | Tennis   |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
| 3        | Windy  | Yes       | Rich    | Cinema   |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
| 4        | Rainy  | Yes       | Poor    | Cinema   |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
| 5        | Rainy  | No        | Rich    | Stay In  |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
| 6        | Rainy  | Yes       | Poor    | Cinema   |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
| 7        | Windy  | No        | Poor    | Cinema   |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
| 8        | Windy  | No        | Rich    | Shopping |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
| 9        | Windy  | Yes       | Rich    | Cinema   |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
| 10       | Sunny  | No        | Rich    | Tennis   |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |
| B        | Justify the significance of Cross Validation. List the type and discuss k-fold cross validation in detail?   | CO3       |         |          |       |          |   |       |     |      |        |   |       |    |      |        |   |       |     |      |        |   |       |     |      |        |   |       |    |      |         |   |       |     |      |        |   |       |    |      |        |   |       |    |      |          |   |       |     |      |        |    |       |    |      |        |  |

|     |  |
|-----|--|
| CO1 | Gain knowledge about basic concepts of Machine Learning and understand the difference between supervised and unsupervised techniques |
| CO2 | To select, apply and evaluate an appropriate machine learning model for the given  |
| CO3 | Ability to understand regression techniques.   |

Q1. A Justify why Linear Progression is not used for performing classification

Consider (Why linear regression is not used for classification)

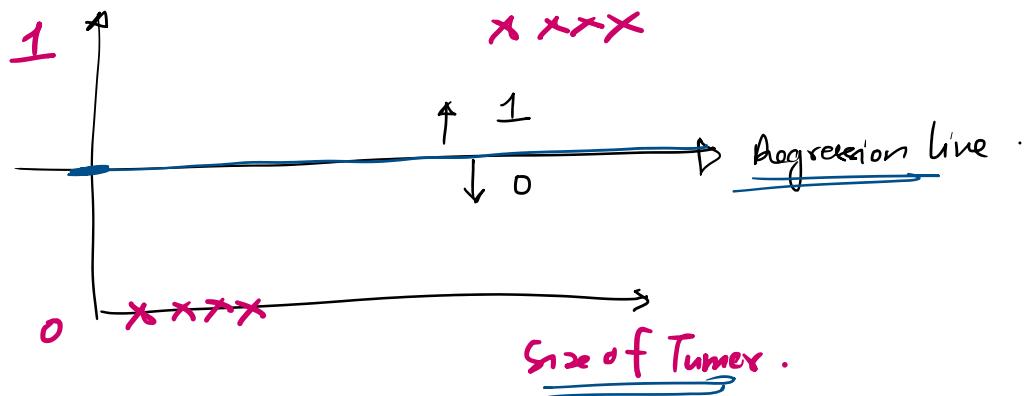
Consider



Suppose

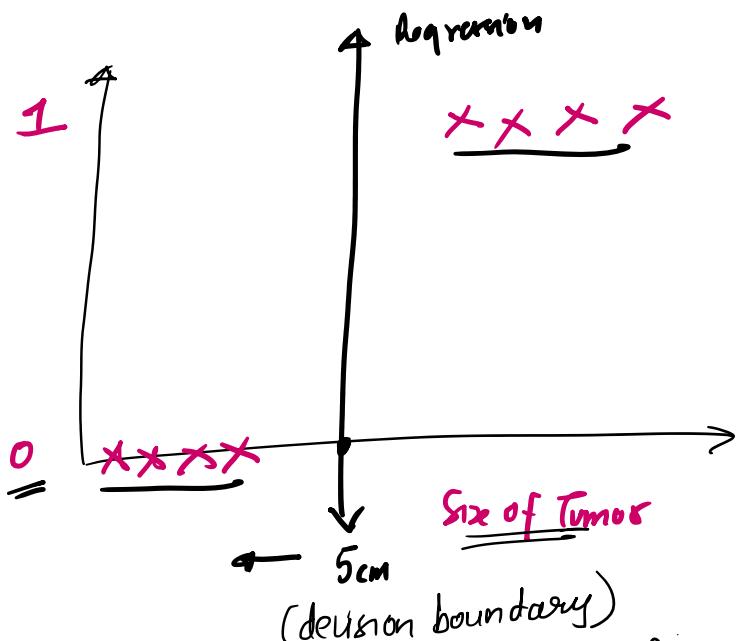
Malignancy

Not Correct



Consider

Malignancy



here we can observe for above regression line

If tumor size  $\leq 5\text{cm}$   $\rightarrow$  Yes (1) Malignant

If tumor size  $\leq 5\text{cm}$   $\rightarrow$  Yes (1) Malignant  
 tumor size  $> 5\text{cm}$   $\rightarrow$  No (0) Benign

We can say let 'P' denote Probability that  $y=1$  when  $\underline{\underline{X=x}}$ .

$$P = \underline{\underline{P}}(y=1 | \underline{\underline{X=x}}) = \frac{\beta_0 + \beta_1 x}{\text{In linear Reg}}$$

$P$  = probability lies bet<sup>n</sup>  $\underline{\underline{0 \text{ to } 1}}$

But linear function are unbounded.

and Expected o/p here is 0 or 1

So we cannot use regression to build classifier

$\therefore$  linear regression is not suitable for classification.

Q1. B

Elaborate on:

- (i) Karl Pearson's Coefficient of Correlation.
- (ii) R square method.

- \* To evaluate performance of m.l. models
- \* To find how good the model fits on given data

### (1) Karl Pearson's Coefficient of Correlation ( $\gamma$ )

→ To calculate relationship bet two variable

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

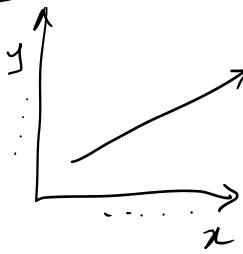
$\gamma$  = quantifies the strength of relationship bet<sup>n</sup> two variables.

The value of  $\gamma$  be between  $\pm 1$  and  $-1$

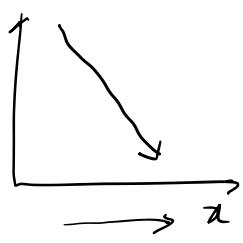
If  $\gamma = 1$  ⇒ Total +ve correlation ⇒ If  $x \uparrow$  then  $y \uparrow$

If  $\gamma = -1$  ⇒ Total -ve correlation ⇒ If  $x \downarrow$  then  $y \uparrow$

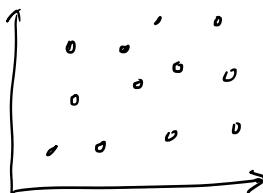
> gives strength (degree) & direction of correlation. or  $x \uparrow$  then  $y \downarrow$



+ve  
correlation



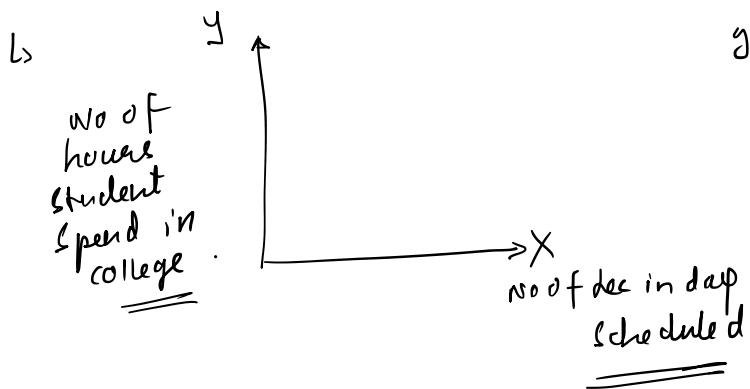
-ve  
correlation



Zero  
Correlation

## ② $R^2$ method (R Square)

↳ It gives information about good of fit feature of the model.



If  $r^2 = 0.85$

Variation in no of hours that students spend in college is 85% dependent on no of lec scheduled.

↳ Indicate percentage of variance in dependent and independent variable pair

> Value varies from 0 to 1

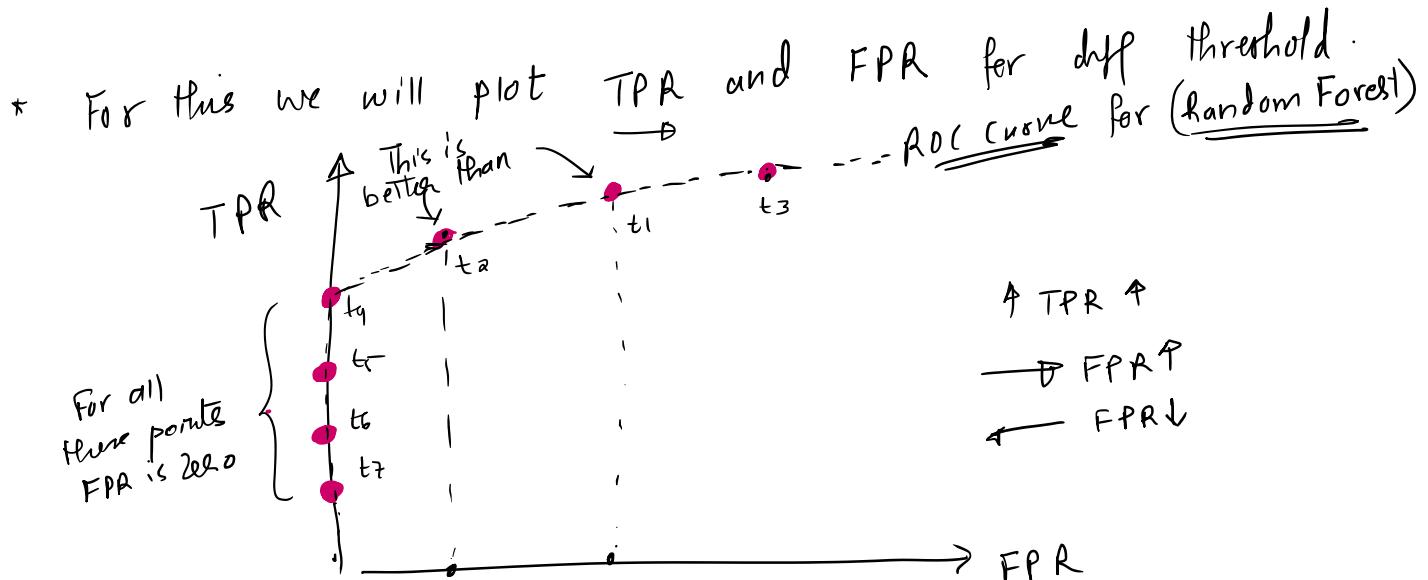
If  $r^2 = 1 \Rightarrow$  no diff bet<sup>n</sup> actual & predicted value.

$r^2 = 0$  means the model does not learn any relationship bet<sup>n</sup> variables -

Discuss the importance of the Receiver Operating Characteristic (ROC) curve and the Area Under the Curve (AUC) in evaluating binary classifiers. How do these metrics help in comparing models?

### ROC [Receiver Operator characteristic]

- \* Consider for logistic regression, where we identify a threshold point and prepare confusion matrix and calculate 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 analyze the performance at diff threshold and want to identify the best of it.

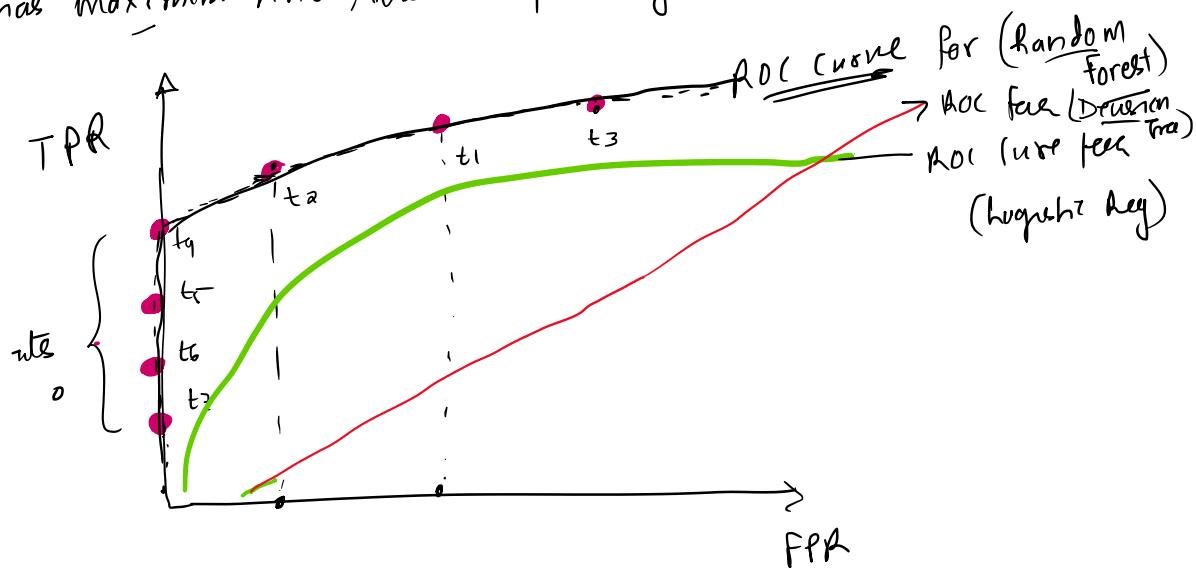


→ from the above ROC curve for Random Forest model, it will help us to find the best threshold.

AUC [Area Under Curve] → It is a method to compare ROC for more than one method and will help to judge which one is better.

\* → Here will find Area Under ROC for each method.

\* → The one will find maximum Area, the corresponding method will be best.



From above ROC Curves, we find that there is maximum area under ROC Curve for Random forest, hence the method Random forest will be the best method.

Outline the steps involved in developing a Machine Learning application for predicting housing prices. Include data collection, preprocessing, model selection, training, evaluation, and deployment.

Solution:-

Steps

① Data Collection → Identify source to collect data  
i.e from real estate website, govt databases  
or real estate agencies-

- \* Scrape or download data containing  
relevant features and historical prices

② Data Preprocessing:

- \* Data Cleaning → Handling missing values  
removing duplicates  
correcting inconsistencies

- + Feature Engineering → Create new features
  - + Normalize or standardize data
  - + Encode categorical variables

- + EDA → Visualize data distributions  
Correlations  
Identify patterns and anomalies

③ Model Selection

- + Algorithm choice
- + Justification

- (4) Data Splitting
- \* Train Test Split (approx 80% train, 20% test)
  - \* Validation set (approx 70% train, 10% validation, 20% test)

- (5) Model Training
- \* Model Training Algo
  - \* Hyperparameter tuning

- (6) Model Evaluation
- \* Performance metrics
  - \* Validation

- (7) Model Improvement
- \* Feature Selection
  - \* Algo optimization
  - \* Regularization

- (8) Deployment
- \* Model Serializ'
  - \* API Development
  - \* Integration

Justify the significance of the F1 score in imbalanced classification problems. How does it address the limitations of using precision or recall alone? Provide a scenario where the F1 score might still be misleading.

### F1-Score

- \* In reality we need a metric that takes into account both precision and recall.
- \* F1 score is a metric that takes into account both precision & recall.
- \* F1 score is harmonic mean of Precision & Recall.

$$\text{F1-Score} = \frac{\frac{2}{\text{Precision}} + \frac{1}{\text{Recall}}}{\text{Precision} + \text{Recall}} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

If F1-Score = 1  $\Rightarrow$  when Precision = 1  
Recall = 1

- \* When Precision and Recall both are high

then F1-Score is high

### When to Use F1 Score $\rightarrow$

$\rightarrow$  Accuracy is not a good metric to use when we have class imbalance.

Ex  $\rightarrow$  let say 99% of people visiting site are onlookers and

Harmonic mean of two variables  $\frac{2ab}{a+b}$   
 $H = \frac{2}{\frac{1}{a} + \frac{1}{b}} = \frac{2ab}{a+b}$ .  
 for n variables

$$H = \frac{n}{\frac{1}{n_1} + \frac{1}{n_2} + \dots + \frac{1}{n_n}} = \frac{n}{\sum_{i=1}^n \frac{1}{n_i}}$$

Ex → let say 99% of people visiting site are onlookers and not purchasing anything.

→ Suppose we have a model that predicts that 99% people visiting site are onlookers.

→ The model is 1% wrong, Generally 1% Error is Acceptable

→ But such model in this case is useless

→ In such case instead of accuracy, we will prefer

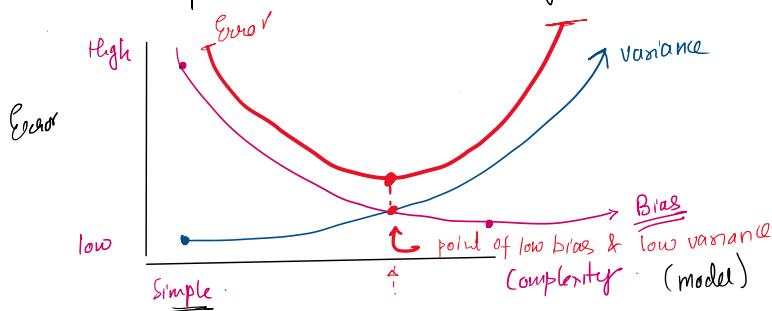
F1-Score.

Justify the bias-variance trade-off in the context of model complexity. How does this trade-off influence the choice of model in a real-world ML application?

### Bias - Variance Tradeoff

① If model is too simple and  
Very few parameters → Underfitting  
→ High bias  
→ Low variance

② If model is complex and  
has large no of parameters → Overfitting  
→ Low bias  
→ High variance.



### 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 and low variance.

The ods to achieve this :-

\* Regularization (Penalize 'Q' parameter)

\* Boosting }

\* Bagging }

\* Impact on choice of models -

\* Simple model (high bias, low variance)

Ex → Logistic regression

Linear regression

\* Complex model (low bias, high variance)

Ex Decision Tree

## Polynomial Regression

- \* Ensemble method > Random Forest  
Gradient Boost machines } can achieve low bias & variance as compared to individual models.
- \* Regularization Techniques
  - \* Helps preventing overfitting
  - \* Need careful tuning of regularization parameters

Ex: Ridge regression  
Lasso regression.

1 meth

For a dataset given below, construct a decision tree using Gini Index, and determine which attribute is a root attribute.

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

Solution  $\Rightarrow$  Independent features:

Weather  
Parent  
Money

Decision/Outcome

Cinema  
Tennis  
Stay In  
Shopping

→ shopping.

Step 1

We will calculate Gini Index for Overall collection of Outcomes of Training Examples.

These are 4 possible outcomes for devision

Cinema → 6 instances  
 Tennis → 2 instances  
 StayIn → 1 instance  
 Shopping → 1 instance.

$$G_{ini} = 1 - \left( \underline{\underline{(6/10)^2}} + \underline{\underline{(2/10)^2}} + \underline{\underline{(1/10)^2}} + \underline{\underline{(1/10)^2}} \right)$$

$$(devision) = 1 - \left( \frac{42}{100} \right) = \underline{\underline{0.58}}$$

Step 2 → find Gini Index for money

Gini (money) → a possible value

Rich  
Poor

(1) Gini (money = Rich) → f instance

Cinema → 3 time  
 Tennis → 2 time  
 StayIn → 1 time  
 Shopping → 1 time.  
 possible devision

$$G_{ini} = 1 - \left( \underline{\underline{(3/7)^2}} + \underline{\underline{(2/7)^2}} + \underline{\underline{(1/7)^2}} + \underline{\underline{(1/7)^2}} \right)$$

$$(money = Rich) = \underline{\underline{0.694}}$$

(2) Gini → n class → Cinema

(2)  $Gini_{(money=poor)}$   $\rightarrow$  3 instance  $\xrightarrow{\text{Decision}}$  Cinema.

$$= 1 - ((3/3)^2) = 0_{//}$$

Weighted Average  $Gini_{(money)} = (Gini_{(money=Rich)} * \text{proportion of Rich}) + (Gini_{(money=poor)} * \text{proportion of poor})$

$$= (0.694 * 7/10) + (0 * 3/10)$$

$$\boxed{Gini_{(money)}} = \underline{0.485}$$

### Step 3 Gini Index on Parent

For Parent feature  $\xrightarrow[\text{values}]{\text{possible values}}$  Yes No

$Gini_{(parent=Yes)} = 5 \text{ instances} \xrightarrow{\text{possible decision}} \text{Cinema}$

$$= 1 - ((5/5)^2) = 0_{//}$$

$Gini_{(parent=No)} = 5 \text{ instances} \xrightarrow{\text{possible decision}} \begin{array}{l} \text{Tennis} \rightarrow 2 \text{ times} \\ \text{StayIn} \rightarrow 1 \text{ times} \\ \text{Shopping} \rightarrow 1 \text{ time} \\ \text{Cinema} \rightarrow 1 \text{ time} \end{array}$

$$= 1 - ((2/5)^2 + (1/5)^2 + (1/5)^2 + (1/5)^2)$$

$$= 1 - (7/25) = 0.72$$

Weighted Average of  $Gini_{(parent)} = (0 * 5/10) + (0.72 * 5/10) = \underline{0.36}$

$$\boxed{Gini_{(parent)}} = 0.36$$

$$\boxed{G_{\text{ini}}^{\prime}(\text{parent}) = 0.36}$$

Step 4) Gini Index for Weather.

Weather       $\xrightarrow[\text{values}]{\text{Possible}}$

- Sunny = 3 instances
- Windy = 4 instances
- Rainy = 3 instances

$G_{\text{ini}}^{\prime}$  ( $\text{weather} = \text{sunny}$ )  $\Rightarrow$  3 instances

possible outcomes

$$= 1 - \left( \left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 \right) = \underline{0.444}$$

$G_{\text{ini}}^{\prime}$  ( $\text{weather} = \text{windy}$ )  $\Rightarrow$  4 instances

possible outcomes

$$= 1 - \left( \left(\frac{3}{4}\right)^2 + \left(\frac{1}{4}\right)^2 \right) = \underline{0.375}$$

$G_{\text{ini}}^{\prime}$  ( $\text{weather} = \text{Rainy}$ )  $\Rightarrow$  3 instances

possible outcomes

$$= 1 - \left( \left(\frac{2}{3}\right)^2 + \left(\frac{1}{3}\right)^2 \right) = \underline{0.444}$$

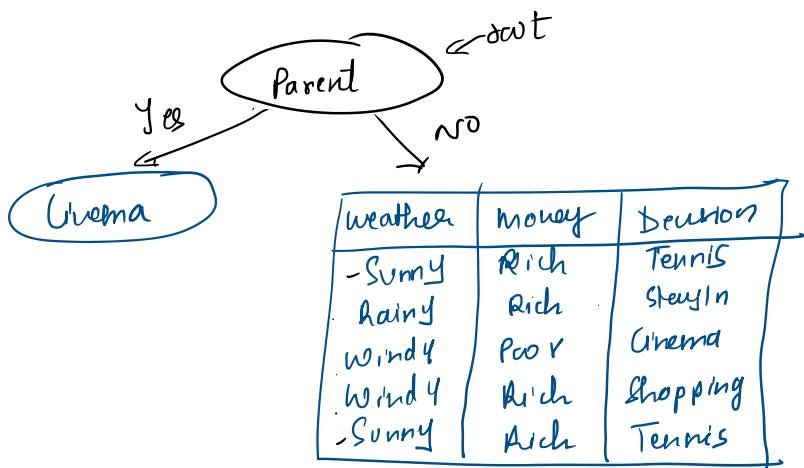
Weighted Average  $G_{\text{ini}}^{\prime}(\text{Weather}) = \frac{(0.444 \times 3/10) + (0.375 \times 4/10) + (0.444 \times 3/10)}{3} = \underline{0.414}$

$$\boxed{1 - r_{\text{ini}} = 0.486}$$

$$\boxed{\begin{aligned} Gini(\text{money}) &= 0.486 \\ Gini(\text{parent}) &= 0.36 \\ Gini(\text{weather}) &= 0.416 \end{aligned}}$$

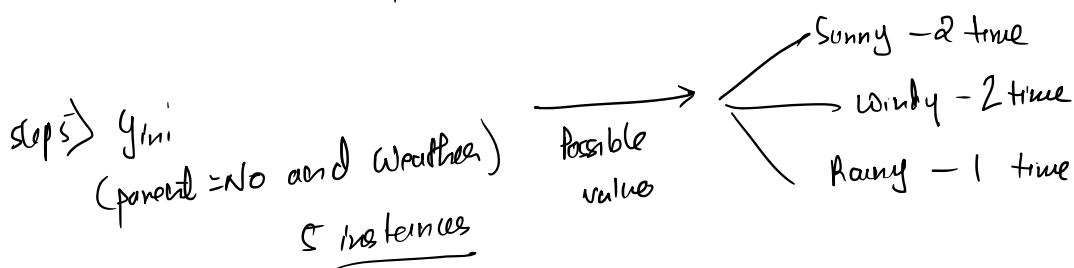
Minimum Gini  $\rightarrow$  Minimum Impurity in Decision.  
Here Minimum Gini Value =  $Gini(\text{parent}) = 0.36$

So the root of Decision is Parent



We need to find  $Gini(\text{parent} = \text{No and weather})$

Also  $Gini(\text{parent} = \text{No and money})$



$$\begin{aligned} Gini(\text{parent} = \text{No and weather} = \text{Sunny}) &\xrightarrow{1 instance} \text{Tennis} = 1 - \left( \left(\frac{1}{2}\right)^2 \right) = 0 // \\ Gini(\text{parent} = \text{No and weather} = \text{Windy}) &\xrightarrow{2 instances} \text{Possible Outcome: Cinema} = 1 - \left( \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right) = 0.5 \\ &\quad \text{Possible Outcome: Shopping} = 0.5 \end{aligned}$$

$$Gini'_{(parent=No \text{ and } Weather=Rainy)} \xrightarrow[1 \text{ instance}]{\text{possible outcome}} StayIn = 1 - ((Y_1)^2) = 0$$

$$\boxed{\text{Weighted Average } Gini'_{(parent=No \text{ and } Weather)} = 0.5 * \frac{4}{5} = \underline{\underline{0.2}}}.$$

Step 6)  $Gini'_{(parent=No \text{ and Money})}$

5 instance

Rich = 4 time  
 poor = 1 time  
 possible values.

$$Gini'_{(parent=No \text{ and Money}=Rich)} \xrightarrow[\text{(4P)}]{\text{possible outcome}} \begin{array}{l} \text{Tennis} - 2 \\ \text{StayIn} - 1 \\ \text{Shopping} - 1 \end{array}$$

$$= 1 - \left( \left(\frac{2}{4}\right)^2 + \left(\frac{1}{4}\right)^2 + \left(\frac{1}{4}\right)^2 \right) = \underline{\underline{0.625}}$$

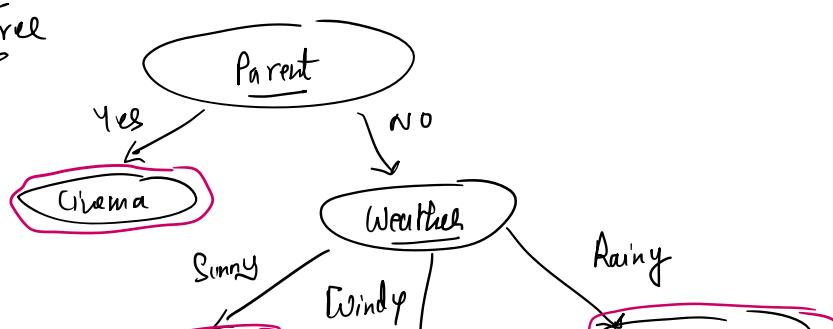
$$Gini'_{(parent=No \text{ and Money}=Poor)} \xrightarrow{\text{possible outcome}} Cinema = 1 - ((Y_1)^2) = 0$$

$$\boxed{\text{Weighted Average } Gini'_{(parent=No \text{ and Money})} = 0.625 * \frac{4}{5} = 0.5}$$

Here  $Gini'_{(parent=No \text{ and Weather})}$  has smallest value

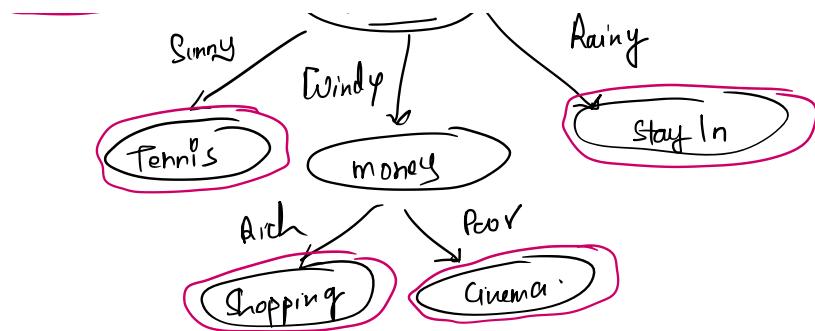
So now Next Node = Weather

\* Updated Decision Tree



Aho

Ans.



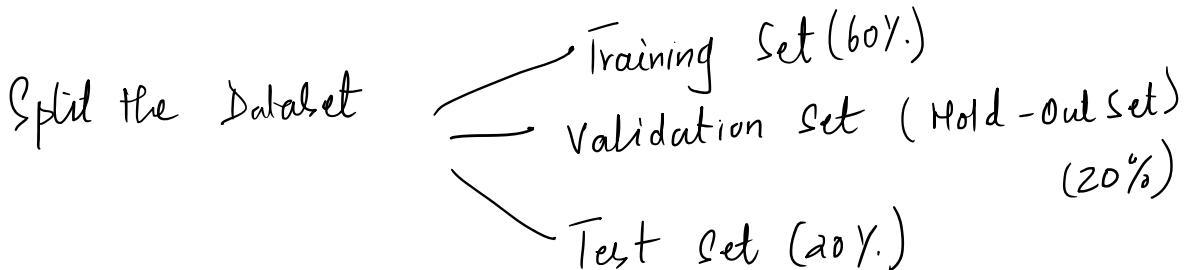
Q3.B

Can you delve into the significance of Cross Validation .List the type and discuss k-fold cross validation in detail?

### Cross Validation →

- In Supervised ML
- Train a model on a Dataset
- Trained model is used to predict the target given new sample.
- How to know if model we have trained will produce effective and accurate result on new input

Cross Validation → It is process that ensures the model will perform well on new Data.



Training Set = part of data on which model is trained.  
(This dataset will help to \* build model)

\* Validation Set ⇒ \* Evaluate the Model

- will help to chk if model overfits or Underfits -
- Update the parameters and again train the model.

- will help to improve the model
- update the parameters and again train the model
- repeat this until the model performs best on Validation set

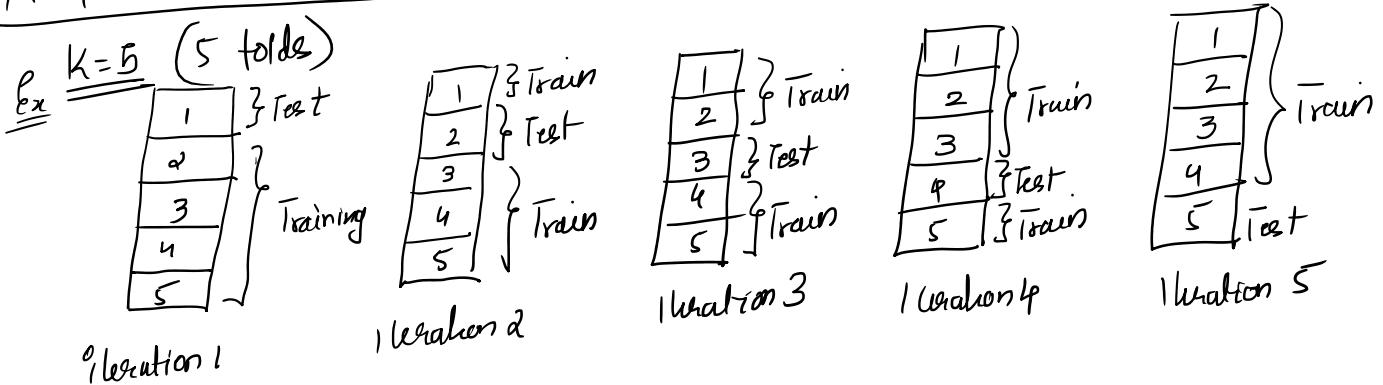
Test Set  $\Rightarrow$  \* Prediction

- The fully trained model after being evaluated on validation set can be used on test set to generate Estimation.

Q) Types of Cross Validation  $\Rightarrow$

- ① The standard Validation Set Approach
- ② Leave one out cross validation
- ③ K-fold Cross Validation

K-fold Cross Validation  $\Rightarrow$



K fold  $\Rightarrow$  K fold helps us to build the model in generalized form.

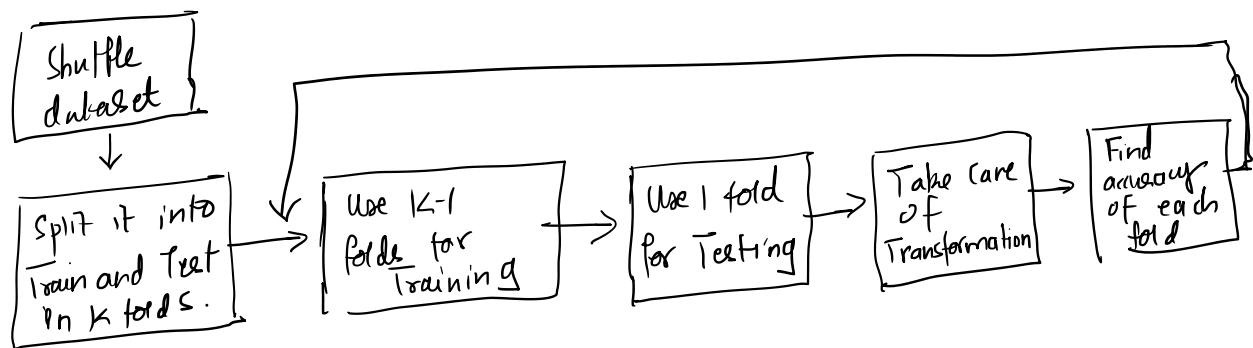
→ To achieve this K fold Cross Validation splits the dataset into Training, Testing & Validation.

→ 10 common ways

dataset into Training Testing & Validation.

→ Here Test and Train data will support building the model.

→ life cycle of K-fold Cross Validation.



- \* The No of iterations ideally is  $K$  time
- \* Finding mean of accuracy score of each iteration will give the consistency of the Trained model.

### Rules

①  $\underline{K \geq 2}$

if  $K=2 \rightarrow$  just 2 iterations.

if  $K=n \underline{\geq 2} \Rightarrow n-1$  for Training  
1 for Testing

② most commonly used value of  $\underline{K=10}$

③ If  $K$  is very large then the running time of process will increase.

④ The value of  $K$  is inversely proportional to size of data i.e. if dataset size is small then number of

∴ true number of folds = 10  
data i.e if dataset size is small then number of folds can increase.