

1. classification

2. prediction

3. clustering

4. dimensionality reduction.

- pandas, Numpy, Matplotlib, Seaborn

- Python

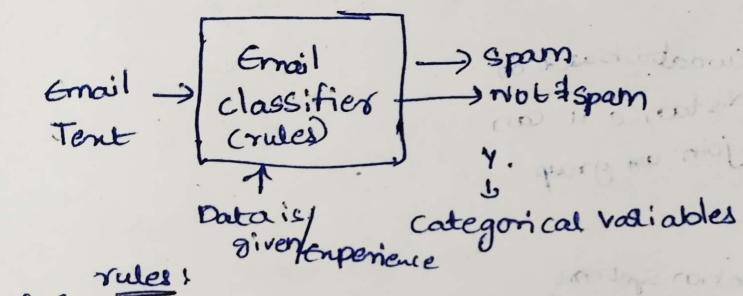
- Jupiter Notebook (Anaconda)

- Colab

classification

examples:

Email classification



if { rules }

- domain email id is not repeated

- contains words like

- discount, free, lottery, win,

then spam

}

else {

not spam.

→ Based on given data we need to write rules

Data:-

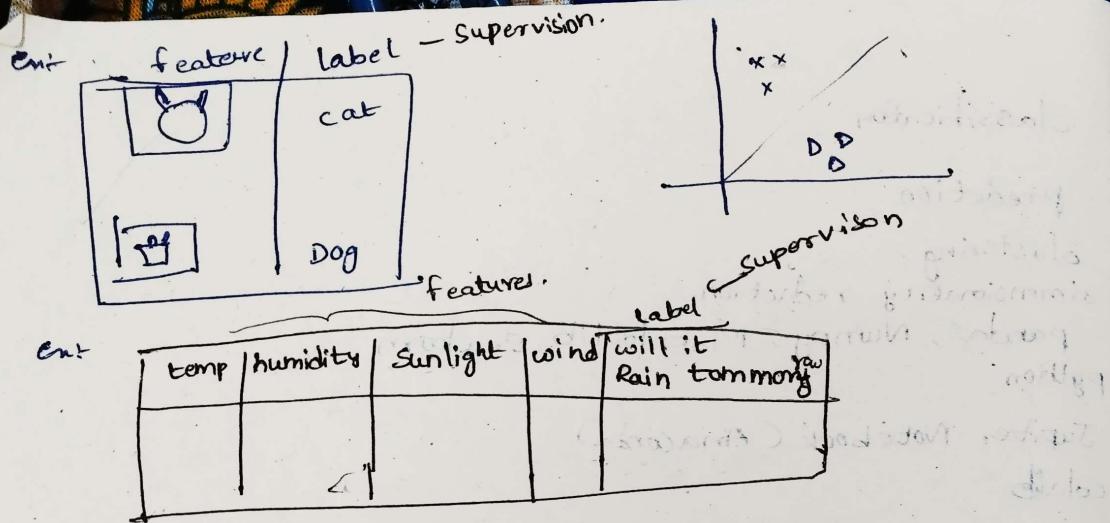
Data can be of - supervised approach

- unsupervised approach.

(a) Supervised : we have both X & Y. Based on given data

X (Email) (Feature)	Label
meeting at 10:00 am	not spam
won \$ 1000.00	spam

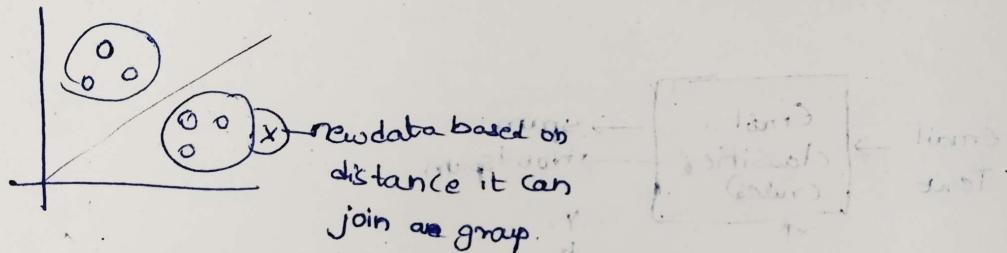
Learning something if we have X & Y & Label is Supervised



(b) UnSupervised Data:

→ have only x NO y (x)

Ex-1



Ex-2:

movie recommendation system:

Language	Actor	Director	Genre	Songs	Budget

(c) Reinforce learning:

try - fail - update

Learning by trying

(d) Semisupervised:

Dataset+

① labelled

1) Supervised label

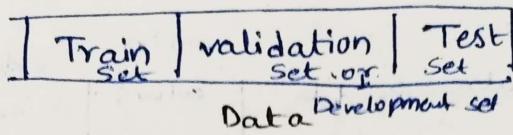
② unlabelled

1) unsupervised label

Supervised learning

- 1) Train
- 2) Test
- 3) Deployment

Splitting that data



Ex:- Train :- Solve examples

validation :- solve exercise

Test :- Solve Exam Questions

if total data is 100%
then train 50 or more than
50% data
validation :- 25 or 10 / 20
test :- 25 or 10 / 20.

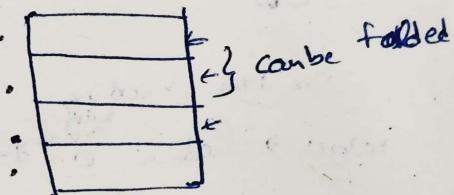
Cross validation

- if dataset size is small

model: (x, y) True label

$x \rightarrow [f(x)] \rightarrow \hat{y}$ predicted label

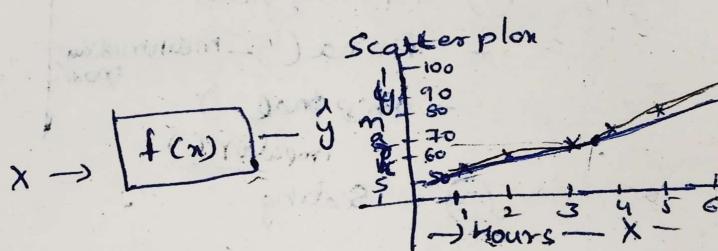
$$\hat{y} = f(x) \rightarrow \text{model.}$$



Regression:

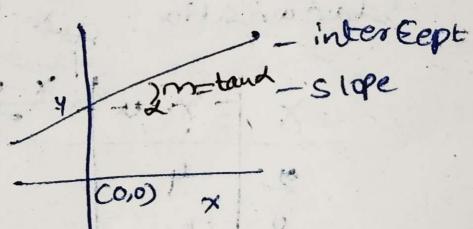
↳ simple linear Regression:

x	y
Hours	marks
1	50
2	55
3	65
4	70
5	80
6	?



$$y = mx + c$$

slope intercept



Data Preprocessing Steps

1) Read Dataset

2) Missing value

3) Encoding - Categorical data

4) Feature Scaling

5) Train - val - Test split

- Missing

One Hot Encoding: Machine may give more importance based on
label \rightarrow Red - 0 $\xrightarrow{\text{label encoding}}$ Red - 1 (Machine may give more importance based on number)

Green - 1
Blue - 2

So we use One Hot Encoder to not to prioritize a specific thing.

One Hot Encoder \rightarrow one way of encoding

	R	G	B
Red	1	0	0
Green	0	1	0
Blue	0	0	1

① missing values

- fillna() - mean, median, mode
- dropna()
- imputer()

② Scaling

- Standardization

- min max scaling

Standardization: $\frac{x - \mu}{\sigma}$ $\xrightarrow{\text{standard normal dist}}$ $\xrightarrow{\text{normal dist}}$

X	X - μ	$\frac{x - \mu}{\sigma}$
2	-2	-1
4	-1	-0.5
6	0	0
8	2	1
10	4	2

Mean = 6

Variance = 8

Scaling: - converting everything in a smaller value or a specific scale based on x. we calculate

Weight	Wage (0-1) / 100
50	0.5
60	0.6
70	0.7
80	0.8
100	1

z value & substitute all x with z so that all are in a range.

Min-Max normalization: $x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$

x	$\frac{x - x_{\min}}{x_{\max} - x_{\min}}$
2	$\frac{2-2}{8} = 0$
4	$\frac{4-2}{8} = 0.25$
8	$\frac{8-2}{8} = \frac{6}{8} = 0.75$
10	$\frac{10-2}{8} = 1$

min 2 max 10

Mean Normalization

$$x' = \frac{x - \mu}{\sigma}$$

based on number
noting.

④ Split

Data no. of study hours (x)	marks (y)	Model 1		Model 2		$\sum xy_i$	x^2
		\hat{y}_i	$e_i = y_i - \hat{y}_i$	\hat{y}_i	$e_i = y_i - \hat{y}_i$		
1	50	15	35	42	8	50	1
2	55	20	35	44	11	110	4
3	65	25	40	46	19	195	9
4	70	30	40	48	22	280	16
5	80	35	45	50	30	400	25
$\Sigma x = 15$		$\Sigma y = 320$		Total error = 195		Total error = 90	$\Sigma x^2 = 55$

1) model 1: $\hat{y} = 10 + 5x$

2) model 2: $\hat{y} = 40 + 2x$

Error = $y - \hat{y}$

Training - finding the parameters that will give the minimum error

making all years positive

$$\text{error} = |\hat{y}_i - y_i|$$

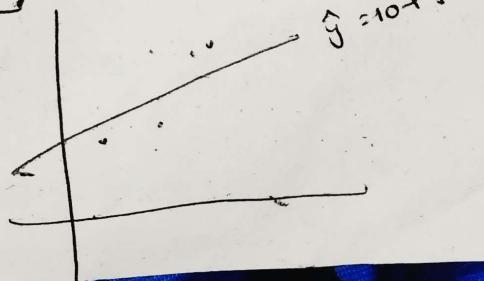
$$\text{error} = (\hat{y}_i - y_i)^2$$

amplifier the error.

$$\hat{y} = \theta_0 + \theta_1 x$$

Parameters: θ_0, θ_1

$SSE = \sum (y_i - \hat{y}_i)^2$
Sum of Squared errors



find θ_0 and θ_1 such that error is minimum.

① Using formul.

$$\theta = \bar{y} - \theta_0 \bar{x}$$

$$\theta_1 = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

$$\theta_1 = \frac{5 \cdot (1035) - (15)(320)}{5(55) - (15)^2} =$$

$$\theta = \frac{5175 - 4800}{225 - 225} = \frac{375}{50} = 7.5$$

$$\theta = 64 - (7.5)(5) = 41.5$$

$$\hat{y} = 41.5 + 7.5x$$

$$41.5 + 7.5t$$

actual model

y	ϵ_i
49	1
56.5	1.5
64	1
71.5	1.5
79	1

error=6

$$\begin{array}{r} 41.5 \\ 14 \\ \hline 55 \end{array}$$

①

model

(has)
parameters
has

② Error

③ Train. - find parameters such that error is minimum.

- What is Machine Learning
- ML focusses on developing algorithms that can
- Learn from historical data
 - Identify patterns and relationships with the data.
 - predict outcomes or take actions on new, unseen data.
- Key characteristics
- Data Driven: ML systems rely heavily on data to make decisions.
 - Iterative: The more data an ML model processes, the better it gets
 - Adaptive: ML Models can adapt/adjust as the environment or data changes

Why is ML Important?

- personalized recommendations
 - predictive maintenance in industries
 - spam detection in emails
 - Medical diagnosis and image analysis
 - Autonomous vehicles and robotics
- ML helps business in
- Automate decision-making
 - Reduce costs and improve efficiency.
 - Discover insights that humans may overlook.

Types of Machine Learning

based on nature of data and tasks 3 types of ML

(1) Supervised Learning:

- Learning from labeled data (l/p & o/p pairs).
- en:- predicting house prices based on size, bedrooms & location.

(2) Unsupervised learning:

- Learning from unlabeled data to identify patterns
- en:- Grouping customers into segments based on spending habits

(3) Reinforcement Learning:

- Learning through trial and error to maximize rewards in an environment.
- Eg:- Teaching a robot to navigate a room by rewarding correct moves.

Types of Machine Learning problems.

- ① → **Supervised**: In supervised learning, the model learns a mapping from input features and their corresponding target labels. It is called supervised because the learning process is guided by labeled data. Goal: Train a model to make predictions on unseen data based on labeled training data.

Key concepts:

Input Features (X): The independent variables

Target Labels (y): The dependent variables (Known outcome)

Training Data: Labeled data used to train the model.

Testing Data: Data used to evaluate the model's performance.

Types

- 1) **Classification**: Predicts a categorical label
ex: email - spam or not spam (binary classification)
handwritten digits into classes :- (0-9)

customer ID	Age	Income	purchased
1	25	50,000	Yes
2	32	60,000	No
3	47	80,000	Yes
4	52	45,000	No

- 2) **Regression**: Predicts a continuous value

ex: Predict the price of a house based on size and location
forecast the temp for next day.

house ID	Size	Bedrooms	Age	Price
1	1500 sqft	3	10	3,00,000
2	2000	4	5	4,50,000
3	1800	3	8	3,50,000

example Algo's

- Linear Regression
- Logistic Regression
- k-Nearest Neighbours
- Decision trees
- Random trees.

- ② Unsupervised Learning
 → the model learns patterns and structures from unlabeled data without any predefined labels or outcomes.
 Goal: Discover hidden structures or patterns in data.

Key concepts:

- Input Features (X): Independent variables with no target labels.
- Clusters: Groups of similar data points.
- Dimensionality Reduction: Simplifying data by reducing features.

Types:

1. Clustering: Group data points into similar clusters.
2. Segment customers into different groups based on their purchasing behaviour.

3. Identify communities in a social network.

Customer ID	Age	Spending Score
1.	22	85
2.	40	150
3.	35	60
4.	28	77

2. Dimensionality Reduction:

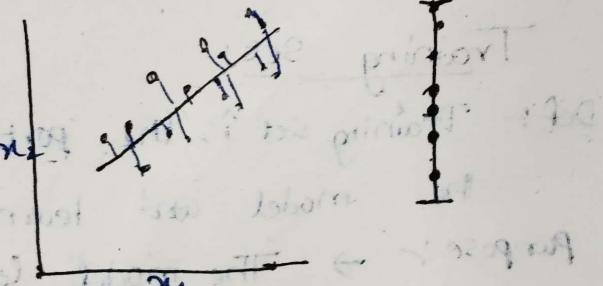
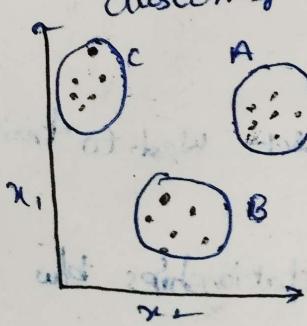
- Reduce the number of features in data while retaining important information.

e.g. Compressing data for feature preprocessing.

Sample ID	Petal Length	Petal width	Sepal length	Sepal width
1	1.4	0.2	4.7	3.2
2	1.5	0.4	5.1	3.5
3	4.5	1.5	6.4	2.9

Goal: use techniques like PCA to reduce these features into fewer dimensions while retaining important patterns.

Plant flower clustering



Example Algorithms

- 1) K-means clustering.
- 2) Hierarchical clustering
- 3) Principal Component Analysis (PCA)
- 4) t-SNE

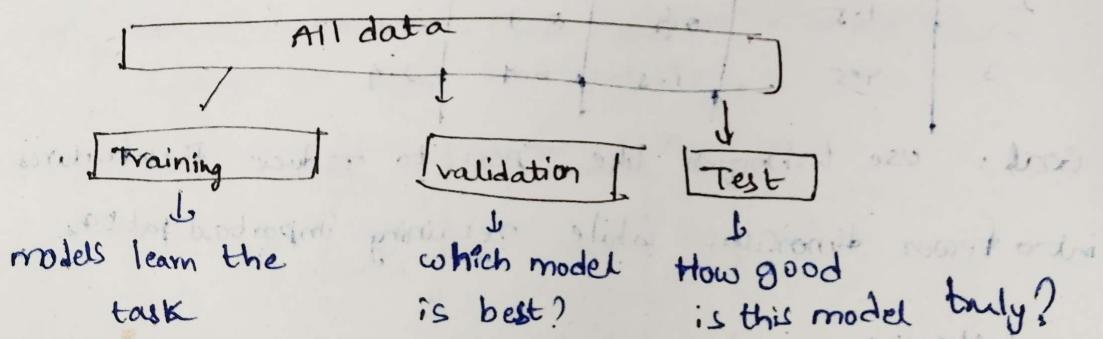
③ 3) Reinforcement Learning:

Definition: RL is a type of ML where an agent learns to make decisions by performing actions in an environment to maximize cumulative rewards.

Differences b/w Supervised, unsupervised, Reinforcement Algos

name	Definition	Example	Reinforcement Algos
1. Supervised learning	Learning from labeled data to make predictions	Spam detection, House price prediction	Linear Regression, Decision Trees.
2. Unsupervised learning	Finding hidden patterns in unlabeled data	Customer segmentation, Anomaly detection	K-means, PCA, Hierarchical clustering.
3. Reinforcement learning	Learning by interacting with an env to maximize rewards	Game playing, Robotics, Self-driving cars	Q-learning, DQN, policy gradient

Evaluation:



1) Training set:

Def → Training set is the portion of dataset used to train the model and learn patterns.

Key purpose: → The model learns the relationships b/w input features and target labels.

- The model adjusts its parameters to minimize the error on this data.

example: from a 1k samples, 800 samples are used to train the model.

2. validation set:

Def: The validation set is a portion of dataset used to tune the model's hyperparameters and assess performance during training.

Key Purpose:

- It helps in identifying the model's ability to generalize before testing
- It is used for hyperparameter tuning (e.g. adjusting learning rate, regularization strength, etc.)
- prevents overfitting by providing feedback on intermediate model performance.

example: out of 1k, 100 are for reserved.

3. Test Set:

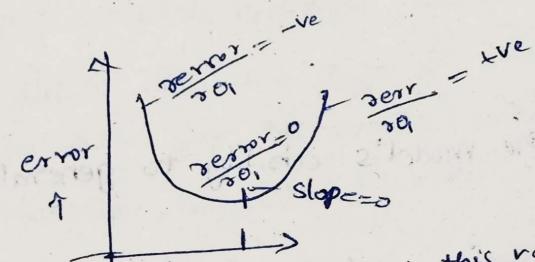
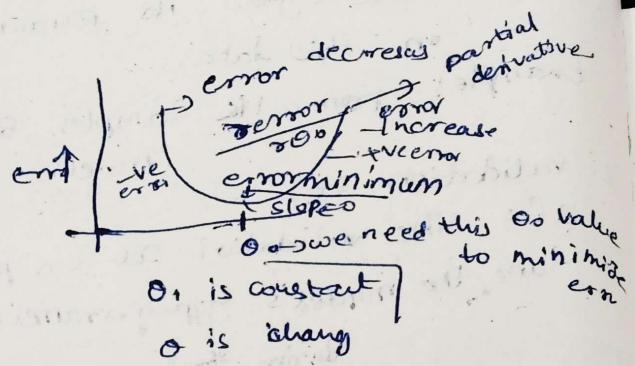
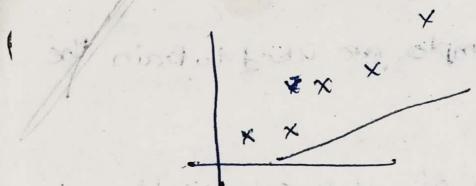
Def: portion of dataset reserved for final evaluation of model's performance on unseen data.

Key purpose:

- provides an unbiased evaluation of trained model.
- Helps estimate the model's generalization performance on new data

ex:- from 1k, 100 will be for testing.

Minimizing the error



$\rightarrow \theta_1$ (or θ) we need this value to minimize error.

slope is
-ve
slope is +ve
 $\frac{dy}{dx}$ derivative

if $\frac{dy}{dx}$ is -ve; -ve slope

if $\frac{dy}{dx}$ is +ve; +ve slope

$\frac{dy}{dx} = 0$; 0 slope

Deriving formulas from intuition

$$\hat{y} = \theta_0 + \theta_1 x$$

find θ_0 and θ_1 such that

$$\frac{\partial \text{error}}{\partial \theta_0} = 0 \quad \text{and} \quad \frac{\partial \text{error}}{\partial \theta_1} = 0$$

$$\rightarrow \text{SSE} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \hat{y} = \theta_0 + \theta_1 x$$

$$\sum (y_i - \theta_0 - \theta_1 x_i)^2$$

$$\frac{\partial \text{SSE}}{\partial \theta_0} = -2 \sum (y_i - \theta_0 - \theta_1 x_i)(-1)$$

$$= -2 \sum (y_i - \theta_0 - \theta_1 x_i)$$

$$= -2 \sum (y_i - \theta_0 - \theta_1 x_i) = 0$$

$$\sum (y_i - \theta_0 - \theta_1 x_i) = 0$$

$$\sum y_i - \sum \theta_0 - \sum \theta_1 x_i = 0$$

$$n\theta_0 = \sum y_i - \theta_1 \sum x_i$$

$$\Rightarrow \theta_0 = \frac{\sum y_i - \theta_1 \sum x_i}{n} \quad \text{--- } \textcircled{1} \text{ eq}$$

$$= \frac{\sum y_i}{n} - \theta_1 \frac{\sum x_i}{n}$$

$$\boxed{\theta_0 = \bar{y} - \theta_1 \bar{x}}$$

formula $\textcircled{1}$

$$\begin{aligned}
 \frac{\partial SSE}{\partial \theta_1} &= -\frac{\sum (y_i - \theta_0 - \theta_1 x_i)^2}{n} \\
 &= -2 \sum (y_i - \theta_0 - \theta_1 x_i) (-x_i) \\
 -2 \sum x_i (y_i - \theta_0 - \theta_1 x_i) &= 0 \\
 \sum x_i (y_i - \theta_0 - \theta_1 x_i) &= 0 \\
 \sum x_i y_i - \theta_0 \sum x_i - \theta_1 \sum x_i^2 &= 0 \\
 \sum x_i y_i - \left(\frac{\sum y_i - \theta_0 \sum x_i}{n} \right) \sum x_i - \theta_1 \sum x_i^2 &= 0 \quad \rightarrow \text{from eq(1)} \\
 \sum x_i y_i - \frac{\sum x_i y_i - \theta_0 \sum x_i}{n} - \theta_1 \sum x_i^2 &= 0 \\
 n \sum x_i y_i - \sum x_i y_i + \theta_0 \sum x_i - n \theta_1 \sum x_i^2 &= 0 \\
 \Rightarrow n \sum x_i y_i - \sum x_i y_i = -\theta_1 (\sum x_i)^2 + \theta_0 \sum x_i \\
 n \sum x_i y_i - \sum x_i y_i &= \theta_1 (n \sum x_i^2 - (\sum x_i)^2) \\
 \theta_1 &= \frac{n \sum x_i y_i - \sum x_i y_i}{n \sum x_i^2 - (\sum x_i)^2} \quad \text{(2) formula}
 \end{aligned}$$

Data Processing

- Data preprocessing is a data mining technique that involves transforming raw data into an understandable format.
- Real world data is often incomplete, inconsistent, and/or lacking in certain behaviours or trends, and is likely to contain many errors. Data processing is a proven method of resolving such issues. Data preprocessing prepares raw for further processing.

Common steps involved in Data Processing

1. Read the Dataset
2. Dealing with missing values
3. Dealing with Categorical data
 - ordinal feature mapping
 - label encoding

- one-hot Encoding

- Feature scaling

- splitting dataset into training and testing sets.

1) Reading Dataset:

Pandas is a software library written for the Python programming language for easy and high performance data manipulation and analysis.

```
df = pd.read_csv("weather.csv")
```

2) Dealing with Missing values

→ many real world datasets contains missing values, often encoded as blanks, NaNs or other placeholders.

Various ways to handle missing values of categorical ways.

1. Ignore observations of missing values if we are dealing with large data set and less number of record has missing value.

2. Ignore variable, if it is not significant

3. Replace with means, median, mode (for numeric data)

4. Replace by most frequent value (for categorical data)

5. Missing values can be treated as a separate category by itself (for categorical data)

6. Develop model to predict missing values.

We need to check whether dataset contains missing values or not.

to check any columns contains missing values.

```
df.isnull().any()
```

to check count of missing values column wise

```
df.isnull().sum()
```

to check count of missing values row wise

```
df.isnull().sum(axis=1)
```

Deleting rows from data frame.

```
x = x.drop([1, 3], axis=0)
```

```
x = x.drop(1, axis=0)
```

Deleting columns from Dataframe

To delete rows and columns from DataFrame, pandas uses the

drop() function.

The drop() function in Pandas is used to delete rows from a DataFrame, with axis set to 0.

Imputation of missing values

Computing the overall mean, median or mode is a very basic imputation method and it is very fast.

from sklearn.preprocessing import Imputer

```
imputer = Imputer(missing_values = "NaN", strategy = "mean")  
numeric_values = x.iloc[:, 1:3]
```

```
x.iloc[:, 1:3] = imputer.fit_transform(numeric_values)
```

Dealing with categorical data

→ As we know scikit-learn models only numerical data. If any categorical values present in dataset we need to convert them into numerics.

① nominal:

② Ordinal: - T-shirt size (XL > L > M)

Ordinal feature Mapping

We should convert the categorical string values into integers.

However, since there is no convenient function that can automatically derive the correct order of the labels of our size feature, we have to define the mapping manually.

```
color = ['green', 'red', 'blue']
```

```
size = ['M', 'L', 'XL']
```

```
a = list(zip(color, size))
```

```
df1 = pd.DataFrame(data=a, columns=["color", "size"])
```

```
size_mapping = {'M': 1, 'L': 2, 'XL': 3}
```

```
df1['size'] = df1['size'].map(size_mapping)
```

```
df1
```

```
size_mapping.items()
```

Label Encoding (Nominal feature encoding)

Label encoding is a utility class to help normalize labels such that they contain only values 0 and n-classes-1. Because scikit-learn's estimators treat class labels without any order, we can use LabelEncoder class to encode the string labels into integers.

will not accept nan's

```
from sklearn.preprocessing import LabelEncoder
```

```
labelencoder_x = LabelEncoder()
```

```
x.iloc[:, 0] = labelencoder_x.fit_transform(x.iloc[:, 0])
```

```
labelencoder_y = LabelEncoder()
```

```
y = labelencoder_y.fit_transform(y)
```

OneHotEncoding (Nominal feature encoding)

OneHotEncoding transforms categorical feature to one new feature of integers (0 to n-categories-1).

one-hot encoding can be obtained with OneHotEncoder, which transforms each categorical feature with n-categories possible values into n-categories.

```
from sklearn.preprocessing import OneHotEncoder
```

```
onehot = OneHotEncoder(categorical_features=[0])
```

make sure x has only numeric values

```
x = onehot.fit_transform(x).toarray()
```

Imputation with categorical data

we can handle in two ways

1. Replace by most frequent value (mode)

2. Missing values can be treated as a separate category

by itself.

getting mode

```
md = x['outlook'].mode()
```

filling missing values with mode

```
x['outlook'] = x['outlook'].fillna(md[0])
```

(or)

as a separate category
 $x[\text{outlook}] = x[\text{outlook}].\text{final unknown}]$.

Feature Scaling

→ Used to standardize the range of independent variables, also known as data normalization.

→ If one of the features has a broad range of values, the distance will be governed by particular features. Therefore, the range of all features should be normalized so that each feature contributes approximately proportionately to final distance.

→ Another reason is that gradient descent converges much faster with feature scaling than without it.

feature scaling methods are

1. Standardization

2. Rescaling (min-max normalization)

3. Mean normalization

4. scaling to unit length.

Standardization

→ Standardization is widely used for normalization in many ML algorithms (e.g., Support vector machine, logistic regression).

→ Feature Standardization makes the values of each feature in the data have zero-mean (when subtracting the mean in numerator) and unit-variance.

X	$\frac{X - \mu}{\sigma}$
2	-4
4	-2
6	0
8	2
10	4

$$\text{mean} = 0$$

$$\text{variance} = 1$$

$$\frac{\sum x}{n} = 6 \quad \mu = \frac{\sum x}{n} = 6$$

Standardization:

Standardization is widely used for normalization

$$x' = \frac{x - \mu}{\sigma}$$

2. Min-Max Normalization

Scales the values b/w a given range (usually 0 to 1)

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

Given data : [10, 20, 30, 40, 50]

$$\text{Min}(x_{\min}) = 10, \text{Max}(x_{\max}) = 50$$

Original data (x) | Min-Max scaled (x')

10	0.00	$\frac{10 - 10}{50 - 10} = 0$
20	0.25	
30	0.50	
40	0.75	
50	1.00	

3) Mean Normalization:

Scales the values b/w -1 and 1 using mean and range of data

$$\text{Formula } x' = \frac{x - \mu}{x_{\max} - x_{\min}}$$

Given data : [10, 20, 30, 40, 50]

$$\text{Mean } (\mu) = 30, \text{Min} = 10, \text{Max} = 50$$

Original data (x) | Mean Normalize (x')

10	-0.67
20	-0.33
30	0.00
40	0.33
50	0.67

Implementing using scikit-learn

standardization :

```
from sklearn.preprocessing import StandardScaler  
scaler = StandardScaler()  
x_scaled = scaler.fit_transform(X)
```

Min-Max Normalization :

```
from sklearn.preprocessing import MinMaxScaler  
scaler = MinMaxScaler()  
x_scaled = scaler.fit_transform(X)
```

Mean Normalization

$$x_{\text{scaled}} = \frac{(x - \bar{x}_{\text{mean}}) / (\bar{x}_{\text{max}} - \bar{x}_{\text{min}})}{(\bar{x}_{\text{max}} - \bar{x}_{\text{min}})}$$

Regression

Regression is a fundamental technique in ML used to model relationships b/w variables. It is widely applied to predict outcomes based on input data.

Regression can be broadly categorized into

1. Linear Regression:
2. Polynomial Regression

1. Linear Regression:-

Linear Regression assumes a straight-line relationship b/w the independent variable(s) and the target variable.

The general equation for a linear regression model is

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

where \hat{y} : predicted value n : number of features

x_i is i^{th} feature value θ_j : The $j^{\text{-th}}$ model parameter
(including the bias term θ_0 and the feature weights
 $\theta_1, \theta_2, \dots, \theta_n$).

Simple Linear Regression

→ when there is only one feature ($n=1$), the eq. simplifies to

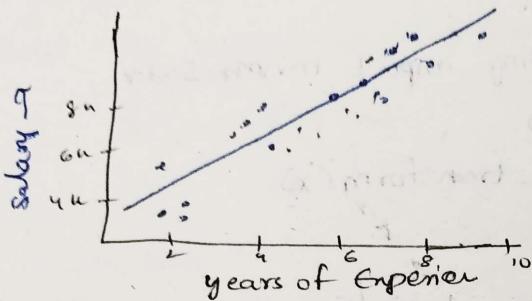
$$\boxed{\hat{y} = \theta_0 + \theta_1 x_1}$$

ex: a company uses experience as independent variable to predict salary.

x_1 : Experience \hat{y} : predicted salary

- goal: Understand how experience impacts salary.

Salary vs Experience



Derivation of SLR

SLR aims to model the relationship b/w one independent variable x and one dependent variable y . The rls is

$$\hat{y} = \theta_0 + \theta_1 x$$

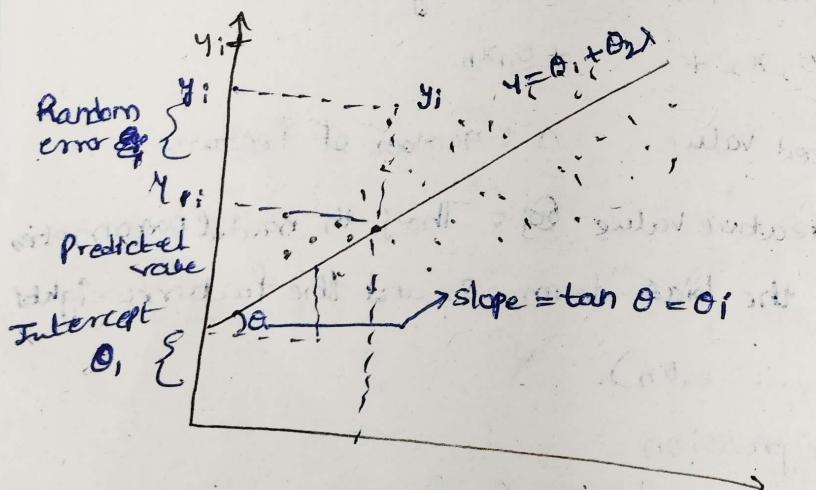
\hat{y} : Predicted value x : independent variable

θ_0 : Intercept θ_1 : Slope (rate of change value of y when $x=0$)

→ to determine values of slope θ_1 and intercept θ_0 that minimize the sum of squared residuals (errors).

1. Sum of Squared Residuals

The residual (error) for each datapoint is the difference b/w the observed value y_i and the predicted value \hat{y}_i .

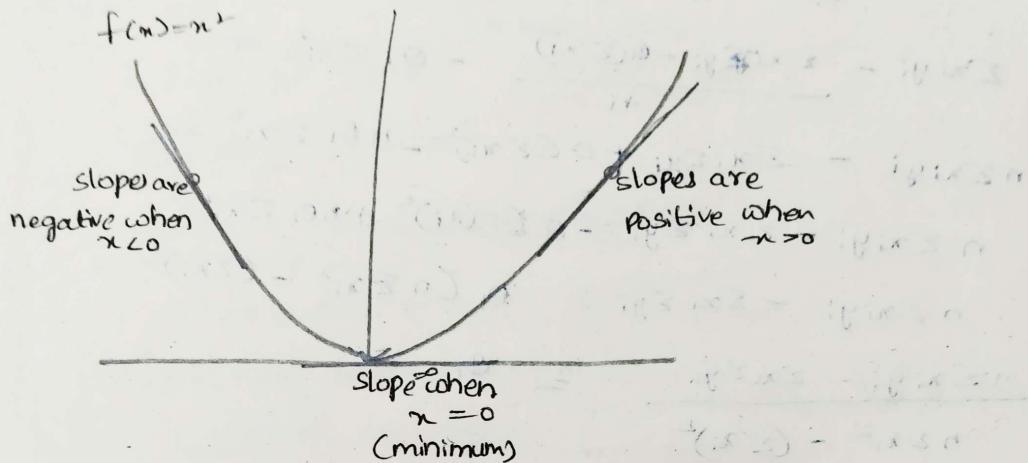


Sum of squared Errors (SSE)

$$SSE = \sum_{i=1}^n (y_i - (\theta_0 + \theta_1 x_i))^2$$

2. Deriving the slope (θ_1)

To minimize SSE, we differentiate SSE with respect to θ_0 and θ_1 , and set the derivative to zero.



Step 1:- Partial Derivative with respect to θ_0 .

→ Taking the derivative of SSE with respect to θ_0

$$\frac{\partial SSE}{\partial \theta_0} = -2 \sum_{i=1}^n (y_i - (\theta_0 + \theta_1 x_i))^2$$
$$= -2 \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)$$

Set this derivative to zero.

$$\sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i) = 0$$

after Rearranging we get

$$\sum y_i = n\theta_0 + \theta_1 \sum x_i$$

$$\theta_0 = \frac{\sum y_i - \theta_1 \sum x_i}{n} \quad \text{--- eq(1)}$$

$$\boxed{\theta_0 = \bar{y} - \theta_1 \bar{x}}$$

Step 2:- Partial Derivative with respect to θ_1 .

Taking the derivative of SSE with respect to θ_1 ,

$$\frac{\partial SSE}{\partial \theta_1} = -2 \sum_{i=1}^n x_i (y_i - \theta_0 - \theta_1 x_i)$$

Set this derivative to zero.

$$\sum_{i=1}^n x_i(y_i - \theta_0 - \theta_1 x_i) = 0$$

expanding & rearranging

$$\theta = \left(\frac{\sum y_i - \theta \sum x_i}{n} \right) \text{ after eq ①}$$

$$\sum x_i y_i = \theta_0 \sum x_i + \theta_1 \sum x_i^2$$

$$\sum x_i y_i - \left(\frac{\sum y_i - \theta \sum x_i}{n} \right) \sum x_i + \theta_1 \sum x_i^2 = 0$$

$$\sum x_i y_i - \frac{\sum x_i \sum y_i - \theta (\sum x_i)^2}{n} - \theta_1 \sum x_i^2 = 0$$

$$n \sum x_i y_i - \sum x_i \sum y_i + \theta (\sum x_i)^2 - n \theta_1 \sum x_i^2 = 0$$

$$n \sum x_i y_i - \sum x_i \sum y_i = -\theta (\sum x_i)^2 + n \theta_1 \sum x_i^2$$

$$n \sum x_i y_i - \sum x_i \sum y_i = \theta_1 (n \sum x_i^2 - (\sum x_i)^2)$$

$$\frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2} = \theta_1$$

$$\boxed{\theta_1 = \frac{\sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}}$$

Regression's

$$\hat{y} = \theta_0 + \theta_1 x$$

$$\theta_0 = \bar{y} - \theta_1 \bar{x}$$

$$\theta_1 = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

$$\theta_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

Independent	
X	Y
1	50
2	55
3	65
4	70
5	80

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$$

$$y = \begin{bmatrix} 50 \\ 55 \\ 65 \\ 70 \\ 80 \end{bmatrix}$$

$$\hat{y} = x \theta = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \\ 1 & 5 \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$$

$$x = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \\ 1 & 5 \end{bmatrix}$$

$$\theta = (x^T y^{-1}) x^T y$$

$$= \frac{2 \times 5 \cdot 5 \times 2}{2 \times 2} \cdot \frac{2 \times 5 \cdot 5 \times 1}{2 \times 1}$$

2x1.

```
import numpy as np
```

```
x = np.array([[1, 0, 1, 2], [1, 1, 0, 4], [1, 5]])
```

```
y = np.array([50, 55, 65, 70, 80])
```

```
theta = np.linalg.inv(x.T @ x) @ x.T @ y
```

```
print(theta)
```

$$\text{sse} = \sum_{i=1}^n (y_i - \hat{y})^2 \quad \hat{y} = x\theta$$

$$= (y - \hat{y})^T (y - \hat{y}) \quad \{ \text{both are same.} \}$$

$$= (y - x\theta)^T (y - x\theta)$$

$$= y^T y - y^T x\theta - \theta^T x^T y + \theta^T x^T x\theta$$

$$y^T x\theta = \theta^T x^T y$$

$$1 \times 5 \times 5 \times 2 - 2x1 = 132 - 2 \times 5 - 5x1$$

$$1 \times 1$$

$$= y^T y - 2\theta^T x^T y + \theta^T x^T x\theta$$

$$2 - x^T x(\theta)^2$$

$$\frac{\partial \text{sse}}{\partial \theta} = \cancel{y^T y} - \cancel{2\theta^T x^T y} + \cancel{\theta^T x^T x\theta} \quad \cancel{2 - x^T x(\theta)^2}$$

$$\theta = \theta - 2 \cdot x^T y + 2 \cdot x^T x \theta$$

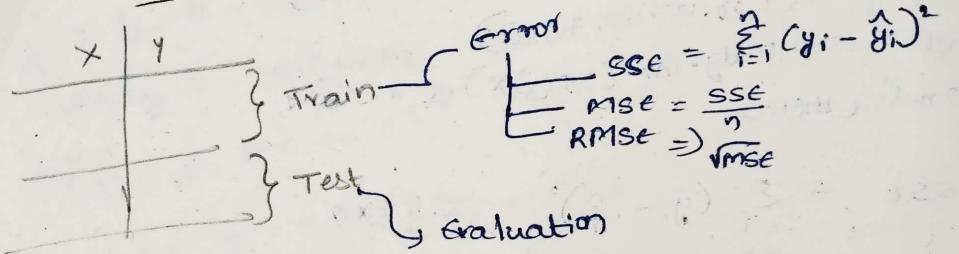
$$2x^T x \theta = 2x^T y$$

$$x^T x \theta = x^T y$$

$$\theta = \frac{x^T y}{x^T x}$$

$$\boxed{\theta = x^T y (x^T x)^{-1}}$$

Evaluation Metrics



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

$$\text{MSE} = \frac{\text{SSE}}{n}$$

$$\text{RMSE} \Rightarrow \sqrt{\text{MSE}}$$

Evaluation

Evaluation metrics

$$R^2$$

$$\text{Adjusted } R^2$$

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} \rightarrow \frac{\text{Variance explained by model}}{\text{Variance in data}}$$

If $R^2 \rightarrow$ good model / perfect fit
 $\rightarrow 0 \rightarrow$ bad model / not fitting.

Variance: how data distributed around mean.

err. & R^2

x	y	\hat{y}	s
50	49	1	
55	56.5	1.5	
65	64.1	1	
70	71.5	1.5	
80	79	1	
		$\bar{y} = 64$	

$$\begin{aligned} & (1 + (0.5)^2 + 1^2 + (-0.5)^2 + 1) \\ & = (1 + 0.25 + 1 + 0.25 + 1) \\ & = 7.5 \end{aligned}$$

$$\begin{aligned} & (50 - 64)^2 + (55 - 64)^2 + (65 - 64)^2 + (70 - 64)^2 \\ & + (80 - 64)^2 \end{aligned}$$

$$(49 - 64)^2 + (56.5 - 64)^2 + (64.1 - 64)^2 + (71.5 - 64)^2 + (79 - 64)^2 = 570$$

$$225 + 56.25 + 0 + 56.25 + 225 = 588$$

$$\frac{588}{570} = \frac{562.5}{570}$$

$$= 0.986$$

It is a good model.

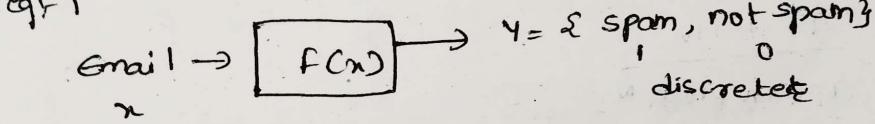
$$\text{Adjusted } R^2 = 1 - \frac{(1-R^2)(n-1)}{n-k-1}$$

$n \rightarrow \text{no. of data points} = 5$
 $k \rightarrow \text{no. of features(inputs)} = 1$

$$= 1 - \frac{(1-(0.98)^2)(4)}{5-1-1} = 0.94$$

Classification:

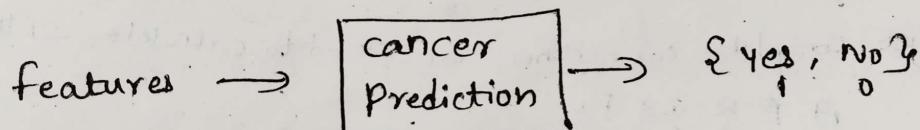
ex: 1



	x	y	\hat{y}
TN		0	0
TP		1	1
FN		1	0
FP		0	1
TN		0	0
TP		1	1
FP		0	1
FN		1	0
TN		0	0

\rightarrow

ex: 2



		Predicted			
		TP	FN		
Actual	1	•	•		
	0	•	•		
		TP	FN		
		FP	TN		

Type error

have cancer & Predicted as cancer \rightarrow TP

no cancer but Predicted as cancer \rightarrow FP

has cancer but, not predicted as cancer \rightarrow FN

no cancer and Predicted No cancer \rightarrow TN

Evaluation Metrics

- Accuracy
 - Precision
 - Recall
 - F1 Score.
- Near to 1 \rightarrow Good model
Near to 0 \rightarrow Bad model (not fitting model)

$$(A) \text{Accuracy} = \frac{\text{Total correct prediction}}{\text{Total prediction}} = \frac{6}{10} = 0.6 = 60\%$$

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$(P) \text{Precision} = \frac{\text{no. of correct positive predictions}}{\text{no. of positive predictions}} = \frac{TP}{TP + FP} = \frac{3}{5} = 0.6 = 60\%$$

$$(R) \text{Recall} = \frac{\text{no. of correct Positive Prediction}}{\text{no. of positive actual Positive cases}} = \frac{TP}{TP + FN} = \frac{3}{5} = 0.6 = 60\%$$

$$(F_1) \text{Score} = \frac{2 * P * R}{P + R} = \frac{2 * 0.6 * 0.6}{0.6 + 0.6} = 0.6$$

\rightarrow when data is unbalance the accuracy is high and info

will be wrong.

Example: 40% patient are left
 $A = 0.80\%$

$$P = 0$$

$$R = 0$$

$$F_1 = 0$$

TP	FN
0	2
FP	TN
0	8

\rightarrow to get complete correctness we need to calculate all the A P R and F1.

cancer

X	0	1
0	0	0
1	0	0
0	0	0
0	0	0
0	0	0
0	0	0
0	0	0
0	0	0
0	0	0
0	0	0

TP	FN
1	2

FP	TN
3	14

$$\text{Acc} = \frac{1+14}{1+2+3+14} = \frac{15}{20} = 0.75$$

$$\text{Pre} = \frac{1}{4} = 0.25$$

$$R = \frac{1}{3} = 0.3$$

$$F_1 = \frac{2 \times 0.25 \times 0.3}{0.25 + 0.3} = 0.27$$

Both A & F should be near values.

Multiclass classification

email	y	\hat{y}
	0	0
	0	0
	0	2
	0	0
	0	2
	1	1
	1	1
	1	1
	1	0
	1	1
	2	1
	2	1
	2	2
	2	2
	2	2

0 - Regular Email
1 - Social Email
2 - Spam Email } more than 2 classes, so multiclass.

0	1	2
0	3	0
1	1	4
2	0	2

$$\text{Accuracy} = \frac{10}{15} = \frac{2}{3} = 0.66$$

class 0: TP = 3, TN = 9, FP = 1, FN = 2

$$P = 0.75 \quad R = 0.6 \quad F_1 = \frac{2 \times 0.75 \times 0.6}{0.75 + 0.6} = \frac{0.9}{1.35} = 0.67$$

class 1: TP = 4, TN = 8, FP = 2, FN = 1

$$P = \frac{4}{6} = 0.66 \quad R = \frac{4}{5} = 0.8 \quad F_1 = \frac{2 \times 0.66 \times 0.8}{0.66 + 0.8} = \frac{0.96}{1.46} = 0.68$$

$$TP = 3 \quad TN = 8 \quad FP = 2 \quad FN = 2$$

$$P = \quad R = \quad F =$$

$$\text{Macro P} = \frac{1}{3}$$

$$= 0.65$$

$$\text{Macro R} = \frac{3}{3}$$

$$= 0.66$$

$$\text{Macro F}_1 = \frac{3}{3}$$

$$= 0.6$$

	En 2 ^c	Y	1
1	0	2	
	0	1	
	0	0	→ Regular Email
	0	0	→ Social Email
	0	0	→ Spam Email
	0	1	
	1	1	
	1	1	
	1	2	
	1	0	
	2	0	
	2	0	
	2	0	
	2	2	
	2	2	

0 → Regular Email
1 → Social Email
2 → Spam Email

$$\text{Accuracy} = \frac{7}{15}$$

0	1	2	
0	2	1	2
1	1	3	1
2	3	0	2

Class 0 ~~not marked~~ $TP = \frac{2}{6}$

not 0 not marked $TN = 6$ $FP = 0$ $FN = 3$

$$P = \frac{2}{6} = \frac{1}{3} \quad R = \frac{2}{5} \quad F_1 = \\ = 0.3 \quad =$$

Class 1

marked $TP = 3$ $not 1, not marked$ $TN = 9$

not 1 marked $FP = 1$ $FN = 2$

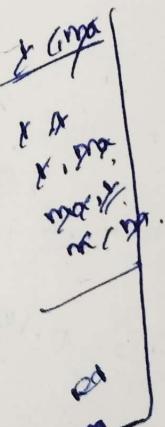
$$P = \frac{3}{4} \quad R = \frac{3}{5} \quad F_1 =$$

Class 2:

marked $TP = 2$ $not 2, not marked$ $TN = 7$

not 2 marked $FP = 3$ $FN = 3$

$$P = \frac{2}{5} \quad R = \frac{2}{5} \quad F_1 =$$



Classification Algorithms

① Decision Tree :-

$\{x, y\}$ can't be spam, Not spam.
 ↓
 discrete.

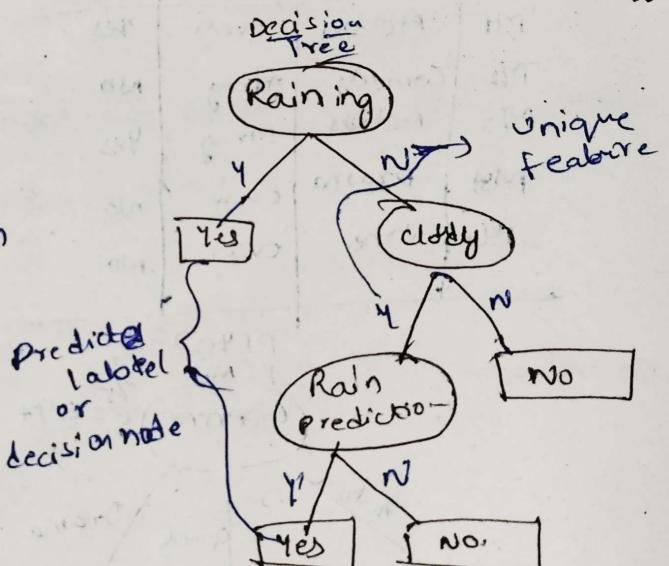
Ans -

cloudy,
 weather prediction
 raining. → will you
 take Umbrella → Yes/No

cloudy	weather prediction	raining	Umbrella ?	Pred
N	N	N	N	N
Y	N	N	N	N
Y	N	Y	Y	Y
Y	Y	Y	Y	Y

Acc = 100%.

→ If raining
 then Yes
 elseif cloudy
 if rain prediction
 then Yes
 else
 'No'
 else
 'No'



② Uncertainty and Entropy:-

→ How much certain an event is → measured by Entropy.

Cn: more certainty less entropy.
 will the students wear uniform tomorrow.

Today	C1	C2	C3	C4
Y	Y	Y	Y	Y
Y	Y	Y	N	Y
Y	Y	W	N	N

more certainty → Y
 entropy is less

less certainty
 entropy more.

$$c_1: P(\text{Uniform}) = \frac{4}{4} = 1$$

$$P(\text{Nonuniform}) = \frac{0}{4} = 0$$

$$c_2: P(U) = \frac{3}{4}$$

$$P(N_U) = \frac{1}{4}$$

$$c_3: P(U) = \frac{1}{4}$$

$$P(N_U) = \frac{3}{4}$$

$$c_4: P(U) = \frac{2}{4}$$

$$P(N_U) = \frac{2}{4}$$

$$\text{Entropy} = -p_1 \log_2 p_1 - p_2 \log_2 p_2$$

$$\text{Entropy}(c_1) = -1 \log_2 1 - 0 \log_2 0 = -1 \log_2 \frac{1}{2} = \frac{-\log_2 1}{2} = \frac{\log_2 1}{2} = 0$$

$$\text{Entro}(c_4) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = -\frac{1}{2} \log_2 \frac{1}{2} = \frac{-\log_2 1}{2} = \frac{\log_2 1}{2} = 0$$

$$E(c_2) = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = 0.81$$

$$E(c_3) = -\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4} = 0.81$$

Uncertainty \uparrow Entropy \uparrow

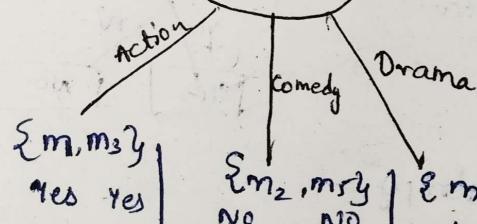
Certainty \uparrow Entropy \downarrow

Dear Watch Movie?

Movie	Genre	Time	watch?
M ₁	Action	Even	Yes
M ₂	Comedy	Morning	No
M ₃	Action	Morning	Yes
M ₄	Drama	Even	No
M ₅	Comedy	Even	No

$$\begin{aligned} P(\text{Yes}) &= \frac{2}{5} \\ P(\text{No}) &= \frac{3}{5} \\ E &= 0.97 \end{aligned}$$

$$\text{Genre}$$



$$\begin{aligned} P(\text{Yes}) &= \frac{2}{3} \\ P(\text{No}) &= 0 \\ E &\Rightarrow 0 \end{aligned}$$

$$\begin{aligned} P(\text{Yes}) &= \frac{2}{5} \\ P(\text{No}) &= \frac{3}{5} \\ E &= \end{aligned}$$

$$\text{Time}$$

$$\begin{aligned} \text{morning} & \quad \text{Evening} \\ \{M_2, M_5\} & \quad \{M_1, M_4, M_5\} \\ \text{No} & \quad \text{Yes} \\ P(\text{Yes}) &= \frac{1}{2} \\ P(\text{No}) &= \frac{1}{2} \\ E &\Rightarrow 1 \end{aligned}$$

$$\begin{aligned} P(\text{Yes}) &= \frac{1}{3} \\ P(\text{No}) &= \frac{2}{3} \\ E &\Rightarrow \end{aligned}$$

$$\text{avg} \\ E = 0$$

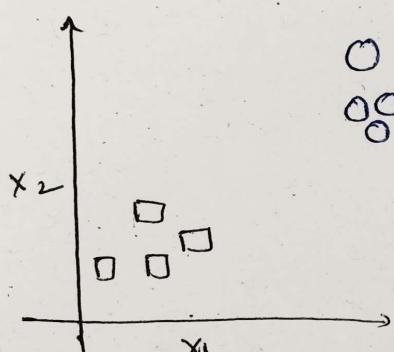
Introduction to classification

→ classification is the process of categorizing data or objects into predefined class or categories based on their features or attributes.

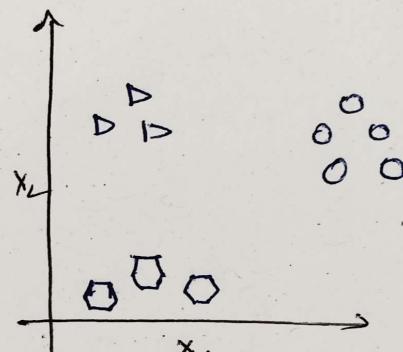
objective of classification (C)

→ main goal of C is to build a model that can accurately assign a label or category to a new observation based on its features.

ex:- problem :- classify images as either "dogs" or "cats".



Binary



Multi

I. Binary classification

→ The task is to classify input into one of two possible classes

ex:- predict if a person has a disease based on health data

positive & Negative

II. Multiclass classification

→ The task is to classify the input into one of several possible classes.

shape

- circle, triangle, rectangle, sphere etc.