

To-Dodemodude

1) import pandas as pd  
 data = {  
 'Name': ['N', 'V', 'W', 'X', 'Y'],  
 'USN': [111, 222, 333, 444, 555],  
 'Marks': [64, 30, 80, 90, 50]  
 }  
 df = pd.DataFrame(data)

print("Sample data")  
 print(df.head())

2) from sklearn.datasets import load\_diabetes

diabetes = load\_diabetes()  
 df = pd.DataFrame(diabetes.data, columns=diabetes.feature\_names)

df['target'] = diabetes.target  
 print("Sample")  
 print('df.head()')

3) file\_path = '/content/Sample-data/california-housing.csv'

df = pd.read\_csv(file\_path)  
 print("Sample")  
 print('df.head()')  
 print("\n")

4) import pandas as pd  
df = pd.read\_csv('content/Dataset of  
Top 100 Banks India - 2021.csv')  
print(df.head(1))

df.to\_csv('output.csv', index=False)  
print("Data saved to output.csv")

5) import yfinance as yf

import pandas as pd

import matplotlib.pyplot as plt

tickers = ["HDFCBANK.NS", "ICICIBANK.NS",  
"KOTAKBANK.NS"]

data = yf.download(tickers, start="2020-01-01", end="2024-12-31",  
group\_by='ticker')

print("Print 5 rows of the dataset:")  
print(data.head())

6) plt.figure(figsize=(12, 6)).i  
 plt.subplot(2, 1, 1)  
 & HDFCBANK\_data['close'].plot  
 (title="HDFCBANK -  
 (1) Last 10 closing price")  
 plt.subplot(2, 1, 2)  
 HDFCBANK\_data['Daily Returns'].plot  
 (color='orange')  
 title="HDFCBANK - Daily  
 Returns",  
 color='orange'  
 plt.tight\_layout()  
 plt.show()

~~Save as pdf~~

~~File~~  
 "HDFCBANK.pdf" = saved.

~~File~~  
 "HDFCBANK.pptx" = saved.

~~File~~  
 "HDFCBANK.xls" = saved.

~~File~~  
 "HDFCBANK.csv" = saved.

~~File~~  
 "HDFCBANK.xlsx" = saved.

~~File~~  
 "HDFCBANK.ppt" = saved.

~~File~~  
 "HDFCBANK.pptx" = saved.

LAB-2ID3 Algorithm

ID3 algorithm is a popular decision tree algorithm used in machine learning. It aims to build decision tree by iteratively selecting the best attribute to split the data based on information gain.

Each node represents a test on an attribute and each branch represents a possible outcome of the test.

$$H(S) = - \sum (P_j * \log_2 (P_j)) \rightarrow \text{Entropy}$$

$$IG(A, D) = H(S) - \sum \frac{|S_v|}{|S|} \times H(S_v)$$

Pseudocode of ID3  $\rightarrow$  Information Gain

def ID3(D, A):

if D is pure or A is empty:  
return a leaf node with  
the majority class in D

else:

A-best = argmax (Information gain (D, A))

root = Node (A-best)

for v in values (A-best):

D-v = subset (D, A-best  
v)

child = ID3 (D-v, A-s  
A-best)

root.child = child (v, child)

return root

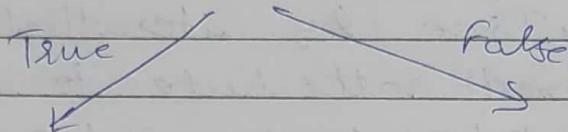
outlook  $\leq 10.5$

entropy = 0.94

samples = 14

value = [S, a]

class = 1



entropy = 0.0

samples = 4

value = [0, 4]

class = 1

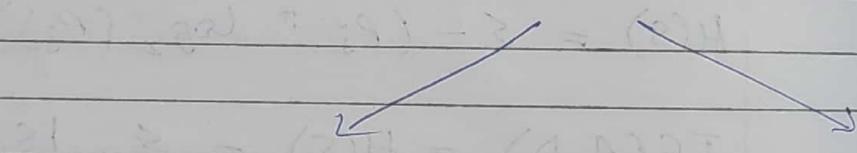
outlook  $\leq 10.5$

entropy = 1.0

samples = 10

value = [S, a]

class = 0



entropy = 0.971

samples = 5

value = [2, 3]

class = 1

entropy = 0.911

samples = 5

value = [3, 2]

class = 0

Entropy = 0.94

Gini coefficient = 0.45

Information value = 0.94

Information gain = 0.49

Information loss = 0.49

Information entropy = 0.94

Information uncertainty = 0.94

Information disorder = 0.94

ID 3 - Code Implementation  
 without using scikit-learn  
 library

```
import pandas as pd
import numpy as np

data = pd.read_csv("weather.csv")
df = pd.DataFrame(data)

def entropy(data):
    class_counts = data.value_counts()
    prob = class_counts / len(data)
    return -np.sum(prob * np.log2(prob))
```

```
def information_gain(df, feature, target):
    total_entropy = entropy(df[target])
    feature_value = df[feature].unique()
    weighted_entropy = 0
    for value in feature_values:
        subset = df[df[feature] == value]
        weighted_entropy += (len(subset) / len(df)) * entropy(subset[target])
    return total_entropy - weighted_entropy
```

```
def id3(df, features, target):
    if len(df[target].unique()) == 1:
        return df[target].iloc[0]
```

if not features:

return df[target].mode()[0]

gains = {feature: information\_gain}

(df, feature, target) for  
feature in features}

best\_feature = max(gains, key=gains)

tree = {best\_feature: {}}

for value in df[best\_feature].unique():

subset = df[df[best\_feature] == value]

tree[best\_feature][value] = id3(subset,

{f for f in features if

F != best\_feature},

target)

return subtree

features = ['outlook', 'temperature',  
'humidity', 'wind']

target = 'PlayTennis'

desc\_tree = id3(df, features, target)

print(desc\_tree)

### OUTPUT

{'Outlook': {'Sunny':

{'Humidity': {'High': 'No', 'Low': 'Yes'}}}

'Overcast': {'Yes': 'Rain', 'No': {'Wind':

{'Weak': 'Yes',

'Strong': 'No'}}}}

 CamScanner

## End to End ML Project

### ① Get Data

```
from sklearn.datasets import fetch_california_housing  
cal_house = fetch_california_housing()  
(as frame + tree)  
print(cal_house.bDESCR)  
cal_house.head()
```

### ② Discover & Visualize the data to gain insights

```
import matplotlib.pyplot as plt  
import seaborn as sns  
  
print(cal_house.frame.isnull().sum())  
  
plt.figure(figsize=(8, 6))  
sns.histplot(cal_house.frame['MedHouseval']  
            Kde = True)  
plt.title('Dist')  
plt.show()
```

### ③ prepare the data for ML Algo

```
from sklearn.model_selection import train_test_split  
from sklearn.preprocessing import StandardScaler
```

```
X = cal_house.frame.drop(['MedHouseval']  
                           , axis=1)
```

$y = \text{cal\_house\_flame} [ \text{'MedHouseVal'} ]$

$X\text{-train}, X\text{-test}, Y\text{-train}, Y\text{-test} =$   
 $\text{train-test-split}(X, y, \text{test\_size}=$

$\text{val\_size=0.2, train\_size=0.8}$ ,  
 $\text{random\_state=42})$

(223d. house) ->

(0.000000000000000)

at what we get result is correct  
 at least one

the no. the different from  
 and no. about 1000

((0.000000000000000)) -> (0.000000000000000)

((0.000000000000000)) -> (0.000000000000000)  
 ((0.000000000000000)) -> (0.000000000000000)  
 ((0.000000000000000)) -> (0.000000000000000)

(0.000000000000000)

(0.000000000000000)

0.000000000000000

at approximately taken with me!  
 different

- house -> house price average month me!  
 Standard -

house -> each month me! -> X  
 (0.000000000000000)

Scales = StandardScaler()

X-train-Scale = scales.fit\_transform(  
X-train)

X-test-Scale = scales.transform(X-test)

#### ④ Select a model and train it

from sklearn.linear\_model import Linear  
Regression

from sklearn.metrics import mean -

squared\_error

linear\_model = LinearRegression()

linear\_model = fit(X-train-Scale,  
y-train)

y-pred-linear = linear\_model.predict

(X-test-mse-linear =  
mean-square-error

(y-test, y-pred,  
linear)

sd-linear = sd-score(y-test, y-pred)

(linear)

print(f'mse-linear), {sd-linear})

#### ⑤ Fine Tune Your Model

from sklearn.metrics import root  
mean\_squared

import numpy as np

y-pred = model.predict(X-test)

rmse = root\_mean\_squared\_error(y-test,

print(f'rmse) {rmse})

LAB-4

## Simple Linear Regression

function linear reg (x, y, z, num)

$$m = 0$$

$$b = 0$$

$$n = \text{length}(x)$$

for  $i \leftarrow 1$  to  $\text{num}$ 

pred = []

for  $i \leftarrow 0$  to  $n-1$ pred =  $m[i] + b$ 

append pred to pred

$$\text{error} = [ ]$$

for  $j \leftarrow 0$  to  $n-1$ 

$$\text{error}[j] = \text{pred}[j] - y[j]$$

error.append(error)

$$\text{cost} = (\frac{1}{n}) \times \sum (\text{error}[j]^2) \quad \text{for } j \in \text{error}$$

$$m\text{-gradient} = \frac{2}{n} \times \sum (\text{error}[j] \times x[j]) \quad \text{for } j \in \text{error}$$

$$b\text{-gradient} = \frac{2}{n} \times \sum (\text{error}) \quad \text{for } j \in \text{error}$$

$$m = m - m\text{-gradient} \times \lambda$$

$$b = b - b\text{-gradient} \times \lambda$$

return m, b

## Logistic linear Regression

function logistic - regression ( $x_{ij}$ ,  $y_{ij}$ , num)

$$m = \text{length}(x[0])$$

$$n = \text{length}(x)$$

weights = array of 0 to length in  
bias = 0

for i = 1 to num:

$$\text{pred} = []$$

for i < 0 to n - 1:

$$\text{lin} = \text{dot\_product}(\text{weights}, x[i]) + \text{bias}$$

$$\text{pred} = \text{sigmoid}(\text{lin})$$

$$\text{predappend}(\text{pred})$$

$$\text{error} = []$$

for j < 0 to n - 1:

$$\text{error} = \text{pred}[j] - y[j]$$

$$\text{error.append}(\text{error})$$

$$\text{cost} = -\frac{1}{n} * \sum(y[i] * \log(\text{pred}[i]))$$

$$+ (1 - y[i]) * \log(1 - \text{pred}[i])$$

for j from 0 to n - 1

$$\text{weight\_grad} = \frac{1}{n} * \text{dot\_product}(\text{error}, x)$$

$$\text{bias\_grad} = \frac{1}{n} * \sum(\text{error})$$

$$\text{weights} = \text{weights} - \text{weight\_grad}$$

$$\text{bias} = \text{bias} - \epsilon * \text{bias\_grad}$$

return weight, bias

## Algo Linear Regression

1. Initialise  $m$  &  $b = 0$
2. make prediction  $mx + b$
3. Find error by  $\hat{y} - y$  - pred
4. update  $m$  &  $b$  using gradient descent
5. Repeat 4 steps for minimization

## Algo logistics Regression

1.  $m$  be the length of features initialized to 0.
2.  $n = \text{no. of data points}$
3. Initialize weight & bias to 0
4. Make predictions weight  $x + bias$
5. Use sigmoid function for pred.
6. Find errors using pred -  $y[j]$
7. Update weight & bias using gradient descent
8. Repeat steps 4 Return weight & bias

## Algo Multilinear Regression

Same as linear but with multiple features.

## Multivariate Multilinear Regression

function multilinear - e.g.  $f(x, y, z, m)$

$$m = \text{length}(x[0])$$

weights - array of  $\phi$  of length  $n$

bias -  $\theta_0$  with starting value

$$n = \text{length}(x)$$

for  $i \leftarrow 1$  to  $m$ :

$$\text{pred} = [ ]$$

for  $j \leftarrow 0$  to  $n-1$ :

$$\text{pred} = \text{dot-product}(\text{weight}, x[j]) + \text{bias}$$

$$\text{pred} = \text{append}(\text{pred})$$

$$\text{error} = [ ]$$

for  $j \leftarrow 0$  to  $n-1$ :

$$\text{error} = \text{pred}[j] - y[j]$$

$$\text{error} = \text{append}(\text{error})$$

$$\text{cost} = \frac{1}{n} \times \text{sum}(1 - y[j]) \cdot \text{error}^2 \text{ in } \text{error} = \text{error} \cdot \text{error}$$

$$\text{weight\_grad} = \frac{1}{n} \times \text{dot\_prod}(\text{error}, x)$$

$$\text{bias\_grad} = \frac{1}{n} \times \text{dot\_sum}(\text{error})$$

return weight & bias

## KNN algorithm

- choose the number of neighbors ( $k$ )

$$\text{eg } k = 3$$

- calculate distance from the test point to all training points

use a distance metric like:

Euclidean distance:

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

- Sort the distances in ascending order of the Euclidean's distance

- Select the  ~~$k$  nearest~~ neighbors  
Pick the top  $k$  points with the smallest distance.

- Do a majority vote (for classification)  
 → Count the labels of the  $k$  neighbors.  
 → The label with the most votes is the predicted class.

(for regression): average of the  $k$  nearest neighbors

## SVM algorithm

1. Data loading and ~~Image processing~~
2. Split the Dataset
3. Initialize the SVM classifier
4. One-vs-Rest Training Strategy
5. Binary SVM Training
6. Prediction Phase
7. Evaluation

- • Store the score  
• Predict the class with maximum decision score

→ ~~Compare predicted labels with true labels~~

~~Calculate accuracy :-~~

~~Accuracy = Total test samples Number of correct predictions × 100~~

Inputs :

~~D : Training dataset  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$~~

~~where  $y \in \{-1, 1\}$~~

~~C : Regularization parameter  
(for soft margin SVM)~~

CAB - 6

## i) Random forest Ensemble Method

1. Load dataset D with features X and target Y
2. preprocess the data
  - Handle missing values
  - Encode categorical variables
  - split dataset into train set D-train and test set D-test
3. Initialise Randomforest with parameters
  - n\_estimators (no. of trees)
  - max\_depth (max depth of trees)
  - randomstate (seed of for replicability)
4. For each tree  $t = 1$  to  $n_{\text{estimator}}$ 
  - Randomly select a bootstrap sample of data points from D-train
  - fit
5. After all trees are built, make predictions:
  - a) for each test instance  $x_{\text{test}}$ , aggregate predictions from all trees using majority voting
  - b) prediction = mode / prediction from all trees
6. Evaluate the model:
  - calculate accuracy, precision, recall
  - F1 score
  - between Model Performance

LAB-7Boosting Ensemble Method

1. Read dataset D with features X and target Y

2. preprocess the data

a. handle missing values

b. Encode ~~categorical~~ categorical variables

c. Split dataset into training

set D-train and test set

D-test

3. Initialise the weights of all training samples equally

$$w_1 = w_2 = \dots = w_N = y_N$$

where ~~N~~ N is the no. of training samples

4. For each iteration t = 1 to n-iteration:

a. Train a weak classifier  $\delta_t$  using the weighted dataset D-train with sample weights  $w_i$

b. calculate the error of the classifier

$$\epsilon_t = \left( \sum_{i=1}^N w_i \cdot I(\delta_t(x_i) \neq y_i) \right) / \sum_{i=1}^N w_i$$

where  $I$  is the indicator function

c) Compute the classifier weight  $(\delta_t)$ :

$$\alpha_t = 0.5^{-1} \log((1 - \epsilon_t) / \epsilon_t)$$

d) Update the weights of unclassified points

e) normalize the weight so that they sum to 1:

$$w_{-i} = w_{-i}/\sum w_i$$

5. After all iterations, combine weak classifiers:

a) Final predictions for test point

$$x = \text{sign}(\sum \alpha_i b_i \mathbf{x}_i^T \mathbf{w}_i)$$

b) Evaluate the model:

c) calculate accuracy, precision, recall, F1-score

7. Return model performance

LAB-8

## K-means clustering Algorithm

1. Load dataset  $D$  with datapoints  $x_i$ , where  $x_i \in \mathbb{R}^d$  ( $d$ -dimensional space)
2. Preprocess the data:
  - a. Handling the missing values
  - b. Standardize the data
3. Initialize  $k$  centroids randomly from the dataset:  
 $c_1, c_2, \dots, c_k$
4. Repeat until convergence (and centroids no longer change):
  - a. For each data point  $x_i =$ 
    - i) Calculate the distance to each centroid  $c_j$ :
$$\text{distance}(x_i, c_j) = \|x_i - c_j\|_2$$
    - ii) Assign point  $x_i$  to the closest centroid
$$\text{cluster}(x_i) = \text{argmin}_j \text{distance}(x_i, c_j)$$
  - b. Update centroids by calculating the mean of points in each cluster:
5. convergence : if centroid do not change significantly  
 $\Rightarrow$  iterations, stop

6. Output:

- a. Find centroids  $C_1, C_2, \dots, C_k$
- b. cluster assignments for each data point

7. Optionally calculate clustering metrics  
(e.g. Silhouettes, etc. method)

LAB-9Principle Component Analysis

1. Load data D with data points  
 $\mathbf{x}_i$  where  $\mathbf{x}_i \in \mathbb{R}^{n \times d}$  ( $d$ -dimensional space)

2. Pre process the data:
  - a. Handle missing values
  - b. Standardize the data by subtracting mean of each feature

$$x_j = (x_{ij} - \bar{x}_{-j}) / \sigma_{-j}$$

where  $\bar{x}_{-j}$  is the mean of feature  $j$   
 $\sigma_{-j}$  is the standard deviation

3. Compute the covariance matrix  $\Sigma$  of standardized data:
- $$\Sigma = (\mathbf{Y}^T) * \mathbf{X}^T * \mathbf{X}$$
- where  $\mathbf{X}$  is the data matrix (each row is a data point)

4. Perform eigen value decomposition on the covariance matrix  $\Sigma$ :
- $$\Sigma * V = \Lambda * V$$
- where  $\Lambda$  are the eigenvalues, and  $V$  are the eigenvectors.

5. Sort the eigenvalues in descending order and select the top  $R$  eigenvalues.

6. project the data onto the top eigenvectors (PCA):  
 $X_{\text{new}} = X \cdot V$

where  $V$  is the matrix of the top  $k$  eigenvectors

output:

Transformed dataset  $X_{\text{new}}$   
(reduced dimensional)  
optionally, visualize the  
reduced data ( $\text{if } K=2$  or  
 $K=3$ )

~~Q~~  
~~16/5~~