

LAB - 0

Method 1:

Creating and initializing Data in dataframe directly.

code:

```
import pandas as pd.
```

```
data = {
```

```
    "Name": ["Preeti", "Prajval", "Prem"],  
    "USN": ["IBM22CS208", "IBM22CS220",  
            "IBM22CS341"],
```

```
    "Grade": ["O", "S", "A+"]
```

```
}
```

```
df = pd.DataFrame(data, index=[1, 2, 3])
```

```
df
```

	Name	USN	Grade
1	Preeti	IBM22CS208	O
2	Prajval	IBM22CS220	S
3	Prem	IBM22CS341	A+

Method 2:

Import dataset from sklearn.datasets

code:

```
from sklearn.datasets import load_iris
```

```
iris = load_iris()
```

```
df = pd.DataFrame(iris.data, columns=iris.feature_names)
```

`df['target'] = iris.target_names[iris.target]`
`df.head()`

output:

	sepal length(cm)	sepal width(cm)	petal length(cm)
0	5.1	3.5	1.4
1	4.9	3.0	1.4
2	4.7	3.2	1.3
3	4.6	3.1	1.5
4	5.0	3.6	1.4

petal width (cm)	target
0.2	setosa

Method 3:

Importing dataset in .csv form.

code:

`file_path = '/content/sample_data/california_housing_train.csv'`

`df = pd.read_csv(file_path)`
`df.head()`

output:

	longitude	lat	house-age	b_rooms	b_beds	populat ⁿ
0	-114.31	34.19	15.0	5612.0	1283.0	1015.0

Method 4:-

Importing data from platforms like UCI, Kaggle etc.

code:

```
df = pd.read_csv('Iceland / Dataset_of_Diseases')
print(df.head())
```

output

ID	No.	Patien	Gender	AGE	Urea	C _r	HbA1c
0	502	19975	F	50	4.7	46	4.9

#1 Yahoo Finance api

import yfinance as yf

import pandas as pd

```
import matplotlib.pyplot as plt
```

tickers = ["Reliance.NS", "TCS.NS", "TATA.NS"]

```
print(data.head())
print(data.shape)
print(data.columns)
reliance_data = data['RELIANCE.NS']
print(reliance_data.describe())
reliance_data['Daily Return'] = reliance_data['close'].pct_change()

reliance_data['Date'] = plt.figure(figsize=(12, 6))
plt.subplot(2, 1, 1)
plt.subplot(2, 1, 2)
plt.tight_layout()
plt.show()

reliance_data.to_csv('reliance-stock-data.csv')
```

Q8B
Ans: 3/3) 25

Data Pre-processing is all the transformations applied to data before feeding into algorithm.

1. Data cleaning : Handling Missing values,
Handling categorical data,
Handling Outliers.
2. Data Transformation : Min-max scalar,
Standard scalar.

Implementation :

```
import pandas as pd
```

loading.

```
df = pd.read_csv('housing.csv')
```

display information

```
df.head(), df.info()
```

display statistical information

```
df.describe()
```

display count of unique labels for 'Ocean Proximity'

```
df['Ocean Proximity'].unique()
df['Ocean-proximity'].value_count()
```

Display which attributes in dataset have missing value.

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

1. Adult Income Dataset (adult.csv)

- No columns had missing values.

Diabetes Dataset

- No columns had missing values.

Handling Method: Since no missing values were found, no imputation was performed.

2. Which categorical columns did you identify in the dataset? How did you encode them?

Adult Income Dataset (adult.csv)

- Categorical columns:

workclass, education, marital-status, occupation, relationship, race, gender, native-country, income

Encoding Method: One-hot encoding was applied using pd.get_dummies (df, drop-first = True), which converted these categorical columns into numerical format.

Diabetes Dataset

- Categorical columns:

Gender, class

Encoding: One-hot encoding was applied using pd.get_dummies (df, drop-first = True) to transform them into numeric format.

3. What is difference b/w Min-Max scaling and standard scaling? When would you use one over the other.

Feature

Definition: Scales values b/w a fixed range.

Effect

when to

use

in linear regression and logistic regression

MinMax Scaling:

Scales values b/w a fixed range

$$x_{\text{scaled}} = \frac{x - x_{\text{min}}}{x_{\text{max}} - x_{\text{min}}}$$

$$x_{\text{std}} = \frac{(x - \bar{x})}{(\sigma)}$$

$$x_{\text{scaled}} = x_{\text{std}} * (\text{max} - \text{min}) + \text{min}$$

Preserves original data distribution but scales it within a limited range.

when preserving relationships and original range of data is important.

Standard scaling

centers data around mean with a standard deviation of 1

$$z = \frac{(x - \mu)}{\sigma}$$

changes data distribution to have zero mean and unit variance.

when data follows a normal distribution or contains outliers

- Min - max

- you need data to be within fixed range
- the dataset does not contain extreme outliers

- standard scale

- the dataset has varying units and a gaussian
- There are significant outliers, as standardization is less sensitive to them.

~~checked B
10/3/25~~

LAB 2

ID3 algorithm implementation

```
import pandas as pd
```

```
import numpy as np
```

```
import math
```

```
from graphviz import Digraph
```

```
def cal_entropy(data, target):
```

```
total_samples = len(data)
```

```
class_counts = data[target].value_counts()
```

```
entropy = 0
```

```
for count in class_counts:
```

```
probability = count / total_samples
```

```
entropy -= probability * math.log2(probability)
```

```
return entropy
```

```
def cal_info_gain(data, feature, target):
```

```
total_samples = len(data)
```

```
weighted_entropy = 0
```

```
for value in data[feature].unique():
```

```
subset = data[data[feature] == value]
```

```
subset_entropy = cal_entropy(subset, target)
```

```
weighted_entropy += (len(subset) /
```

```
total_samples) * subset_entropy
```

```
return cal_entropy(data, target) - weight_entropy
```

```
def build_tree(data, target, features, parent_node_class='None'):
```

```
if len(data[target].unique()) <= 1:
```

```
return data[target].unique()[0]
```

if len(features) == 0:

return parent_node_class

parent_node_class = data[target].mode[0]

best_feature = max(features, key = lambda

feature: calculate_info_gain(data, feature
forget)

tree = {best_feature: {}}

features.remove(best_feature)

for value in data[best_feature].unique():

subset = data[data[best_feature] == value]

subtree = build_tree(subset, target, features,

copy(), parent_node_class)

return tree

def visualize_tree(tree, dot=None, node_name=
=None)

if not dot:

dot = Digraph(comment='Decision
tree')

if isinstance(tree, dict):

feature = list(tree.keys())[0]

dot.node(node_name, label=feature)

else:

dot.node(node_name, label=tree)

return dot.

APL

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

```
target = 'PlayTennis'
```

```
features = list(data.columns)
```

```
features.remove(target - column)
```

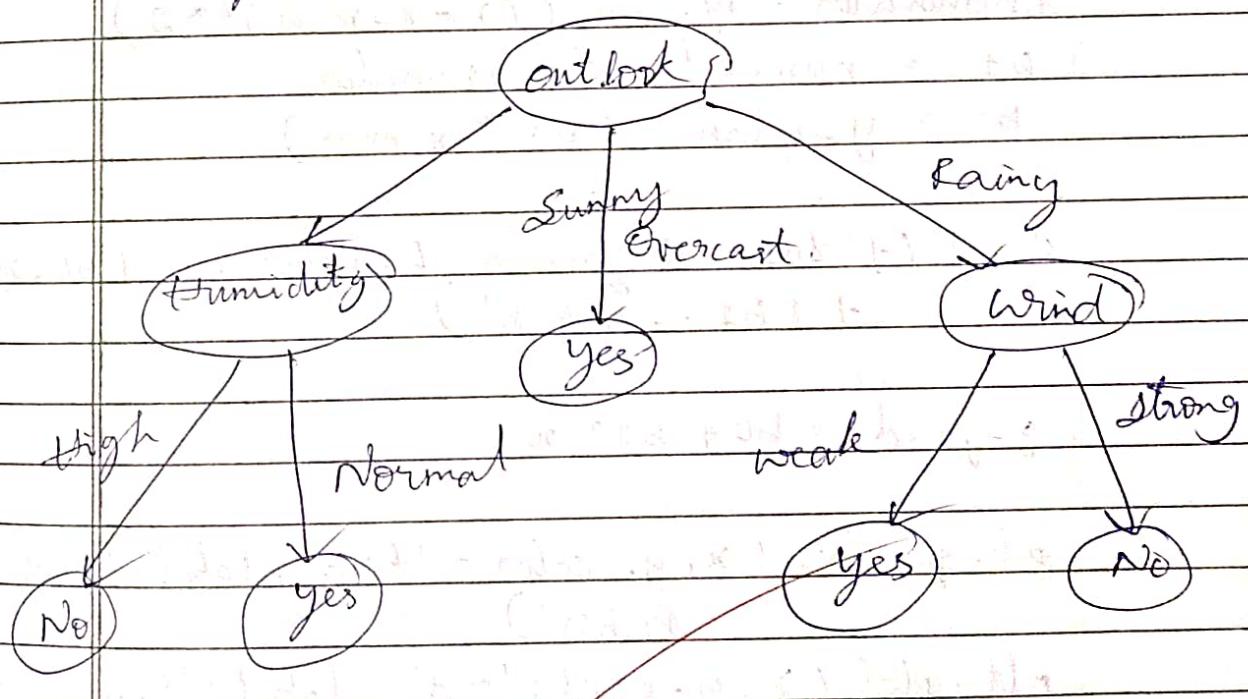
```
tree = build_tree(data, target, features)
```

```
dot = visualize(tree)
```

```
dot.render('decision-tree')
```

```
dot
```

output



Subha
17/3/25

LAB - 3

LINEAR REGRESSION

import numpy as np

import matplotlib.pyplot as plt

$x = np.array([1, 2, 3, 4, 5])$

$y = np.array([2, 4, 5, 4, 3])$

$x_mean = np.mean(x)$

$y_mean = np.mean(y)$

numerator = np.sum((x - x_mean) * (y - y_mean))

denominator = np.sum((x - x_mean) ** 2)

$b_1 = \text{numerator} / \text{denominator}$

$b_0 = y_mean - (b_1 * x_mean)$

print ("Linear Regression Equation : $y = {} b_0 + {} b_1 * x$ "
 $+ " + {} b_1 : . 2 f y x")$

$y_pred = b_0 + b_1 * x$

plt.scatter(x, y, color = 'blue', label = 'Data Points')

plt.plot(x, y_pred, 'red', label = 'Regression Line')

plt.xlabel('x')

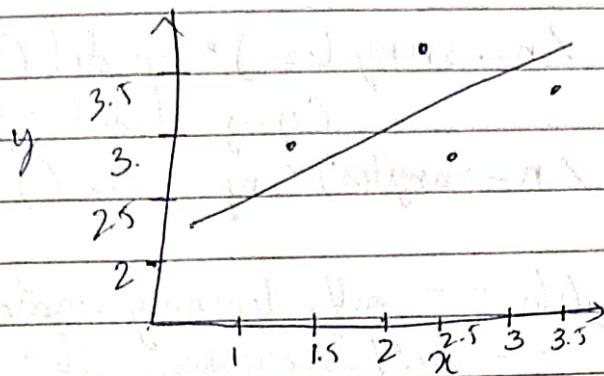
plt.ylabel('y')

plt.title('Linear Regression')

plt.legend()

plt.show()

output: $y = 2.20 + 0.60x$



LOGISTIC REGRESSION

import numpy as np.

class LogisticRegression:

def __init__(self, learning_rate=0.01,
 $n=1000$):

self.learning_rate = learning_rate

self.num_iterations = num_iterations

self.weights = None

self.bias = None

def sigmoid(self, z):

return 1 / (1 + np.exp(-z))

def fit(self, x, y):

n_samples, n_features = x.shape

self.weights = np.zeros(n_features)

self.bias = 0

for _ in range(self.num_n):

linear_model = np.dot(x, self.weights)
 $+ \text{self.bias}$

$y_{predicted} = \text{self_sigmoid}(\text{linear_model})$

$$dw = (1/n_samples) * \text{np.dot}(x.T, (y_{predicted} - y))$$

$$db = (1/n_samples) * \text{np.sum}(y_{predicted} - y)$$

self.weights = self.learning_rate * dw

self.bias = self.learning_rate * db

def predict(self, x):

linear_model = np.dot(x, self.weights) + self.bias

$y_p = \text{self_sigmoid}(\text{linear_model})$

$y_{p_cls} = [1 \text{ if } i > 0.5 \text{ else } 0 \text{ for } i \text{ in } y_p]$

return np.array(y_{p_cls})

if name == "main":

~~x = np.array([[1, 2], [2, 3], [3, 4], [4, 5], [5, 6]])~~

~~y = np.array([0, 0, 1, 1, 1])~~

~~model = LogisticRegression()~~

~~model.fit(x, y)~~

$y_{pred} = \text{model.predict}(x)$

print("Predictions : ", y_pred)

output: [0 1 1 1 1]

LAB-01

KNN algorithm:

```
import numpy as np
import pandas as pd
from collections import Counter

data = {'Feature 1': [1, 2, 3, 4, 5, 6, 7, 8, 9, 10],
        'Feature 2': [2, 4, 1, 3, 5, 7, 6, 8, 10, 9],
        'class': ['A', 'A', 'B', 'B', 'A', 'B', 'A',
                  'B', 'B', 'A']}
```

df = pd.DataFrame(data)

def euclidean_distance(x1, x2):

return np.sqrt(np.sum((x1 - x2) ** 2))

def knn(x_train, y_train, test_point, k=3):

distances = [euclidean_dist(test_point, x) for x in x_train]

k_indices = np.argsort(distances)[:k]

k_n_l = [y_train[i] for i in k_indices]

most_common = Counter(k_n_l).most_common(1)
return most_common[0][0]

x = df[['Feature 1', 'Feature 2']].values

y = df['class'].values

test_point = np.array([6, 6])

predicted_class = knn(x, y, test_point, k=3)

test points = $\{(6, 6), (7, 3), (9, 10)\}$

output:

Predicted class for $(6, 6)$: A

Predicted class for $(9, 10)$: B

Predicted class for $(7, 3)$: A

SVM

```
import numpy as np  
import pandas as pd  
import matplotlib.pyplot as plt
```

```
data = {'feature1': [1, 2, 3, 4, 3, 6, 7, 8, 7, 10],  
        'feature2': [2, 4, 1, 3, 5, 7, 6, 8, 10, 9],  
        }  
        }
```

~~df = pd.DataFrame(data)~~

~~x = df[['feature1', 'feature2']].values~~

~~y = df['class'].values~~

class SVM:

```
def __init__(self, learning = 0.001, l = 0.01,  
             n = 1000):
```

self.ln = learning

self.l = l

self.n = n

self.w = None

self.b = None

```
def fit(self, x, y):
    n_s, n_f = x.shape
    y = np.where(y <= 0, -1, 1)
```

```
self.w = np.zeros(n_f)
self.b = 0
```

for i in range(self.n):

for idx, x_i in enumerate(x):

condition = y[idx] * (np.dot(x_i, self.w) + self.b) >= 1

if condition:

self.w = self.w + x_i * self.lmbd * self.w

else:

self.w = self.w + self.lmbd * self.lmbd *
self.w - np.dot(x_i,
y[idx]))

self.b = -self.w - x_i * y[idx]

def predict(self, x):

app = np.dot(x, self.w) + self.b

return np.sign(app)

clf = SVC()

clf.fit(x, y)

test_points = [[6, 6], [2, 3], [9, 10]]

for point in test_points:

predictn = clf.predict(np.array(point))

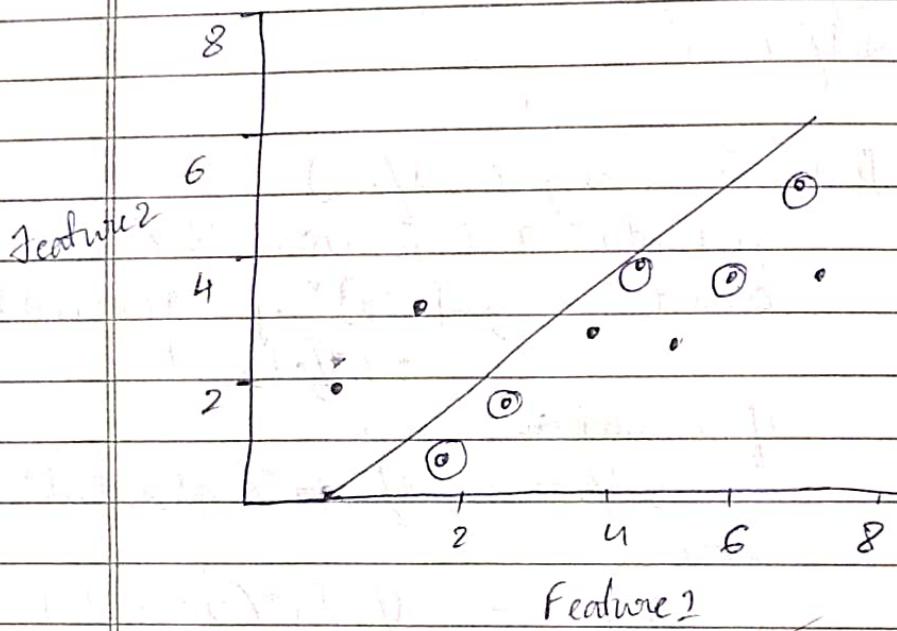
print(f"Predicted class for {point} : {predictn}")

output:

$$(6, 6) : 1$$

$$(2, 3) : -1$$

$$(9, 10) : 1$$



Feature 2

SB
7/4/25

Random Forest

Function Random_Forest (Data D, Integer

N-Trees, Integer m-features):

Forest $\leftarrow \{ \}$

For i from 1 to N-trees do:

Sample $D_i \leftarrow$ Bootstrap-sample(D)

Tree $T_i \leftarrow$ Build-dec-tree (D_i , m-features)

add T_i to forest

return Forest

Function Build_Dcision_Tree (Data D_i , Integer m):

If stopping-condition-mel(D_i):

Return Leaf-node (Most-common-label(D_i))

Features \leftarrow Random-Subset (All-features, m)

Best-feature, threshold \leftarrow Find-Best-Split
(D_i , Features)

Left-data, Right-data \leftarrow split-Data (D_i ,
Best-feature, threshold)

Left-node \leftarrow Build-Dcision (Left-data, m)

Right-node \leftarrow Build-Dcision (Right-data, m)

Return Node (Best-feature, threshold, left-node,
right-node)

Function Predict-Random (Forest, Instance x):

Predictions \leftarrow [Predict-tree (T_i , x) for each
 T_i in Forest]

Return Majority-vote (Predictions)

LAB-07

AdaBoost (Boosting)

Function AdaBoost (Dataset D, Integer T):

Initialize weights $w_i = 1/n$ for each training sample (x_i, y_i)

classifiers $\leftarrow []$ alphas $\leftarrow []$

for t from 1 to T do:

classifier $h_t \leftarrow \text{Train-weak-learner}(D, \text{weights}_t)$ Error $\epsilon_t \leftarrow \sum [w_i * \mathbb{I}(h_t(x_i) \neq y_i)]$ alpha $\alpha_t \leftarrow 0.5 * \log((1 - \epsilon_t) / \epsilon_t)$

for each i from 1 to n do:

 $w_i \leftarrow w_i * \exp(-\alpha_t * y_i * h_t(x_i))$ Normalize weights: $w_i \leftarrow w_i / \sum w_i$ Append h_t to classifiersAppend α_t to Alphas

Return classifiers, alphas

Function Predict - adaBoost (classifiers, alphas, Instance x):

Total $\leftarrow 0$

for t from 1 to T do:

Total $\leftarrow \text{Total} + \alpha_t * h_t(x)$

Return Sign(Total)

K-Means Clustering.

Function k-Means (DataSet D, Integer K,

Integer max-iter):

centroids \leftarrow Initialize-random (D, K)

For it from 1 to max-iter do:

clusters \leftarrow assign-points-to-centroids(D,
centroids)

New-centroids \leftarrow compute-mean-of-clusters
(clusters)

If converged (centroids, New-centroids)
then:

Break

centroids \leftarrow New-centroids.

Return clusters, centroids

Function assign-points-to-centroids (Data D,
centroids):

clusters \leftarrow empty dictionary

For each point x_i in D do:

closest \leftarrow find-nearest-centroid(x_i ,
centroids)

add x_i to clusters[closest]

Return clusters

Function compute-mean-of-clusters (clusters):

For each cluster:

compute mean of all points return update
centroids

LAB-09

Principal Component Analysis (PCA)

Function PCA (Data D, Integer K):

standardized \leftarrow Standardize - Data (D)

CovMatrix \leftarrow Compute - covariance - matrix
(standardized)

Eigenvalues, Eigenvectors \leftarrow Eigen - Decomposit.
(CovMatrix)

Sorted \leftarrow Sort - Eigenvectors - By - Eigenvalues
(Eigenvalues, Eigenvectors)

Top-k-vectors \leftarrow Select Top - k - eigenvectors
(Sorted, k)

Reduced Data \leftarrow Project - Data (standardized,
Top-k-vectors)

Return ReducedData

Function Standardize - Data (Data D):

For each feature :

subtract mean and divide by standard
deviation

Return standardized D

Function Project - Data (Data, Eigenvectors):

Return Matrix - Multiplication (Data,
Eigenvectors).