

VISVESVARAYA TECHNOLOGICAL UNIVERSITY

"JnanaSangama", Belgaum -590014, Karnataka.



ML Lab Report

Submitted by

PRAKHYATI BANSAL
(1BM21CS136)

Under the Guidance of

M Lakshmi Neelima
Assistant Professor
BMSCE

in partial fulfillment for the award of the degree of
BACHELOR OF ENGINEERING

In
COMPUTER SCIENCE AND ENGINEERING



B. M. S. COLLEGE OF ENGINEERING

(Autonomous Institution under VTU)

BENGALURU-560019

March 2024 to June 2024

Date : _____

Lab1

COMPLETE UPLOAD OF READING OF CSV FILE.

```
import pandas as pd  
url = "https://archive.ics.uci.edu/ml/  
machine-learning-databases/iris/iris.data"  
colnames = ["sepal-length-in-cm",  
            "sepal-width-in-cm", "petal-length-in-cm",  
            "petal-width-in-cm", "class"]  
iris_data = pd.read_csv(url, names =  
                        colnames)
```

print(iris_data.head(10))

Output is written in the file 20
Dropout prior step will do next

Output -

sepal-length-in-cm	sepal-width-in-cm	petal-width-in-cm	class
5.1	3.5	1.4	Iris-setosa

ND
2/3/24

Lab-2Housing DatasetEnd-to-End Machine Learning Project

1) We frame the problem, use median & RMSE regression model.

2) Getting the data-

Download-Root holds URL of housing data

10 Housing-Path is used to join "data" & "01" directories.

Housing-URL has it combined.

OS module for creating a directory.
Extract the files using tarfile.open()

• extractall()

3) imports required-

os for joining paths

tarfile for opening tarfiles

urllib for downloading data

matplotlib for plotting data

4) load-dataframe-

pd.read_csv(data-path) to create a dataframe.

Date: _____

5) Create Text-data:

np.random.permutations(0) → generates random numbers.

test_size → store the % of data in test_size

indices [: test_set_size] → first set.

indices [test_set_size:] → second set

use a seed to get the same
data always and store in a file.

15

NP

9/15/20

20

25

Date : _____

Python Implementation of Linear Regression

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

def estimate_coeff(x, y):
    n = np.size(x)
    m_x = np.mean(x)
    m_y = np.mean(y)
    ss_xy = np.sum(x * y) - n * m_x * m_y
    ss_xx = np.sum(x * x) - n * m_x * m_x
    b_1 = ss_xy / ss_xx
    b_0 = m_y - b_1 * m_x
    return (b_0, b_1)
```

```
def plot_regression_line(x, y, b):
    plt.scatter(x, y, color = "m", marker = "o")
    y_pred = b[0] + b[1] * x
    plt.plot(x, y_pred, color = "g")
    plt.xlabel('x')
    plt.ylabel('y')
```

```
def main():
    x = np.array([0, 1, 2, ..., 9])
    y = np.array([1, 3, 2, ..., 12])
```

Date : _____

b = estimate-coeff (x,y)
print(b)
plt - regression-line (x,y,b)

5 output:

$$(b_0, b_1) = (1.2363\ldots, 1.16969\ldots)$$

→ Multiple Linear Regression:

10 from sklearn.model_selection import
~~train-test split~~
import matplotlib.pyplot
import numpy as np
from sklearn import datasets, linear_model,
metrics
15 data-url = "http://lib.stat.cmu.edu/datasets
~~|| boston~~
raw_dt = pd.read_csv(data-url, sep="|",
skiprows=22, header=None)
20 x = np.hstack([raw_df.values[::2, :],
raw_df.values[1::2, :2]])
y = rawdf.values[1::2, 2]
x_train, x_test, y_train, y_test = train-test-
split(x, y, test_size=0.4, random_state=1)
25 reg = linear_model.LinearRegression()
reg.fit(x_train, y_train)

Date : _____

```
print("coefficient:", reg.coef)
print("variance score: {}".format
      (reg.score(x-test, y-test)))
plt.style.use('fivethirtyeight')
5 plt.scatter(reg.predict(x-train),
              reg.predict(x-train)-y-train,
              color="green", s=10, label="Train Data")
plt.scatter(reg.predict(x-test),
              reg.predict(x-test)-y-test,
              color="blue", s=10, label="Test Data")
10 plt.hlines(y=0, xmin=0, xmax=50,
              linewidth=2)
plt.legend(loc="upper right")
plt.title("residual errors")
15 plt.show()
```

Decision Trees (ID3)

```

import numpy as np
import pandas as pd
df = pd.read_csv('url')
5 df.head()
df.info()
df.describe()
def find_entropy(df):
    target = df.keys()[-1]
    entropy = 0
    values = df[target].unique()
    for value in values:
        fraction = df[target].value_counts()[value] / len(df[target])
        entropy += fraction * np.log2(fraction)
    return entropy
15
def BuildTree(df, tree=None):
    target = df.keys()[-1]
    node = find_winner(df)
    att = np.unique(df[node])
    if tree is None:
        tree = {}
        20
        if att == {}:
            tree[target] = np.argmax(df[target])
        else:
            tree[target] = {}
            for value in att:
                sub = get_subtable(df, node, value)
                dValue, counts = np.unique(sub[target], return_counts=True)
                if len(dValue) == 1:
                    tree[target][value] = dValue[0]
                else:
                    tree[target][value] = BuildTree(sub, tree=None)
    25
    for value in att:
        sub = get_subtable(df, node, value)
        dValue, counts = np.unique(sub[target], return_counts=True)
        tree[target][value] = BuildTree(sub, tree=None)

```

Date : _____

```
if len(counts) == 1:  
    tree[Node3Value] = value_Co3  
else:  
    tree[Node3Value] = buildTree(subtable)  
return tree  
tree = buildTree(df)  
import pprint  
pprint.pprint(tree)
```

Decision tree (sklearn):-

```
import pandas as pd  
import numpy as np  
import sklearn.model_selection import  
DecisionTree  
from sklearn.tree import plot_tree.  
df = pd.read_csv('url')  
df.head()  
df.info()  
df.isnull().sum()  
cols = df.columns[0:-1]  
for i in cols:  
    sns.boxplot(y=df[i])  
    plt.show()  
x = df.drop('Species', axis=1)  
y = df['Species']
```

Date : _____

x-train, x-test, y-train =
train-test-split(x, y, test_size=0.3)

dt = DecisionTreeClassifier(max_depth=3)

dt.fit(x, y)

y-pred-train = dt.predict(x-train)

y-pred = dt.predict(x-test)

accuracy-score(y-pred, y-test).

① ~~minimum size of leaf node~~

② ~~maximum depth of tree~~

③ ~~number of splits~~

④ ~~number of nodes~~

⑤ ~~number of trees~~

⑥ ~~minimum entropy~~

⑦ ~~maximum information gain~~

20

25

Logistic Regression

```
(S.0325) import pandas as pd  
import matplotlib.pyplot as plt  
file = "insurance.csv"  
df = pd.read_csv(file)  
from sklearn.linear_model import LogisticRegression  
model = LogisticRegression()  
model.fit(x_train, y_train)  
print(model.score)
```

```
y-predicted = model.predict(x-test)  
print(y-predicted)
```

O/P: Age: 22
probability: 0.1059

Date : _____

KNN

import pandas as pd
from sklearn.datasets import load_iris
5 iris = load_iris()

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

10 df['target'] = iris.target
df['flower-name'] = df['target'].apply(lambda n: iris.target_name[n])

from sklearn.model_selection import
train_test_split

15 m = df.drop(['target', 'flower-name'],
axis='columns')

y = df['target']

20 x_train, x_test, y_train, y_test =
train_test_split(m, y, test_size=0.2)

from sklearn.neighbors import
KNeighborsClassifier

KNN = KNeighborsClassifier(
n_neighbors=10)

25 KNN.fit(x_train, y_train)

Date : _____

Knn. score(x-test, y-test)

0.966

Knn. predict([[4.8, 3.0, 1.5, 0.3]])

5

O/P's 0 (versa)

for i in range(0, 100):
 if i % 2 == 0:

print(i, "is even")
 else:

print(i, "is odd")

for i in range(0, 100):
 if i % 2 == 0:

print(i, "is even")
 else:

print(i, "is odd")

(0) = even

else:
 print(i, "is odd")

Date : _____

K-means Implementation

```
import pandas as pd
data = pd.read_csv("iris.csv")
data.head()
```



```
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.datasets import load_iris
from sklearn.cluster import KMeans
```

`obj.y = load_iris (return_x_y = True)`

kmeans = kmeans(n_clusters=3,
random_state=2)

Want to implement K-means. It can find clusters

`pred = kMeans. fit_predict(cw)`

predator behaviour, 36 pages

1976-392 .193. dk - 9001

(B. C. D. E. F. G.)

John. W. G. T. (1860-1910)

~~(6) Roman Empire - Vikings = Vikings~~

$$25. \quad \text{If } f(x) = 2x^2 + 3x - 1, \text{ find } f(8) \text{ and } f(-2).$$

C'f' = 3.115 phs.

(Cawdor, May)

Date : _____

SVM

```
from sklearn.datasets import load_breast_cancer
import matplotlib.pyplot as plt
from sklearn.inspection import DecisionBoundaryDisplay
from sklearn.svm import SVC
cancer = load_breast_cancer()
X = cancer.data[:, :2]
y = cancer.target
SVM = SVC(kernel="rbf", gamma=0.5, C=1.0)
SVM.fit(X, y)

DecisionBoundaryDisplay.from_estimator(
    SVM,
    X,
    response_method="predict",
    cmap=plt.cm.Spectral,
    alpha=0.8,
    xlabel=cancer.feature_names[0],
    ylabel=cancer.feature_names[1])
plt.scatter(X[:, 0], X[:, 1], c=y, s=20,
            edgecolors='k')
plt.show()
```

PCA

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 %matplotlib inline
5 from sklearn.decomposition import PCA
6 from sklearn.preprocessing import StandardScaler
7
8 from sklearn.datasets import load_breast_cancer
9 data = load_breast_cancer()
10 data.keys()
11 print(data['feature_names'])
12 df1 = pd.DataFrame(data['data'],
13 columns = data['feature_names'])
14 scaling = StandardScaler()
15 scaling.fit(df1)
16 scaled_data = scaling.transform(df1)
17 principal = PCA(n_components=3)
18 principal.fit(scaled_data)
19 X = principal.transform(scaled_data)
20 plt.figure(figsize=(10,10))
21 plt.scatter(X[:, 0], X[:, 1], c=['target'],
22 cmap='plasma')
23 plt.xlabel('pc1')
24 plt.ylabel('pc2')

```

Date : _____

Random Forest Ensemble & Ada Boost

from sklearn.datasets import make_moons
ada, boost

5 from sklearn.metrics import accuracy_score

10 from sklearn.model_selection import train_test_split

15 from sklearn.tree import DecisionTreeClassifier

20 from sklearn.ensemble import RandomForestClassifier
AdaBoostClassifier

(n=10000, n_estimators=100,
noise=0.5, random_state=0)

x-train, y-train, x-test, y-test =

train-test-split

(x-train, y-train, x-test, y-test =
train-test-split)

clf = RandomForestClassifier(n_estimators=100)

max_features="auto", random_state=0)

clf.fit(x-train, y-train)

y-pred = clf.predict(x-test)

accuracy = accuracy_score(y-test, y-pred)

o/p:
0.7965

Date : _____

AdaBoost:

$\text{clf} = \text{AdaBoostClassifier(n_estimates=100)}$

$\text{clf}\cdot\text{fit}(\text{x_train}, \text{y_train})$

$\text{y_pred} = \text{clf}\cdot\text{predict}(\text{x_test})$

$\text{accuracy_score}(\text{y_test}, \text{y_pred})$

op! 0.8393