**B.M.S. COLLEGE OF ENGINEERING BENGALURU**
Autonomous Institute, Affiliated to VTU



Lab Record

**Machine Learning**

*Submitted in partial fulfillment for the 6th Semester Laboratory*

Bachelor of Engineering
in
Computer Science and Engineering

*Submitted by:*

**R SHREYA**

1BM21CS152

Department of Computer Science and Engineering
B.M.S. College of Engineering
Bull Temple Road, Basavanagudi, Bangalore 560 019
Mar-June 2024

| Lab | Program |
| --- | --- |
| 1 | Write a python program to import and export data using Pandas library functions |
| 2 | Demonstrate various data pre-processing techniques for a given dataset |
| 3 | Demonstrate various data pre-processing techniques for a given dataset |
| 4 | Implement Linear and Multi-Linear Regression algorithm using appropriate dataset |
| 5 | Use an appropriate data set for building the decision tree (ID3) and apply this knowledge to classify a new sample. |
| 6 | Build Logistic Regression Model for a given dataset and 6 Build KNN Classification model for a given dataset. |
| 7 | Build Support vector machine model for a given dataset and Build k-Means algorithm to cluster a set of data stored in a .CSV file. And Implement Dimensionality reduction using Principle Component Analysis (PCA) |
| 8 | Implement Random forest ensemble method on a given dataset And Implement Boosting ensemble method on a given dataset. |
| 9 | Build Artificial Neural Network model with back propagation on a given dataset |

Write a python programming to import and export data using pandas library functions

Code:

```
import pandas as pd
x = pd. read_csv(' c: \\ Users \\ Admin \\ Downloads \\ iris \\
                                    iris.csv')
x.columne = col names
x.head()
```

Output:

| | Sepal length in cm | Sepal width in cm | petal length | petal width | class |
|---|---|---|---|---|---|
| 0 | 4.9 | 3.0 | 1.4 | 0.2 | Iris-Setosa |
| 1 | 4.7 | 3.2 | 1.3 | 0.2 | Iris-setosa |
| 2 | 4.6 | 3.1 | 1.5 | 0.2 | Iris-setosa |
| 3 | 5.0 | 3.6 | 1.4 | 0.2 | Iris-setosa |
| 4 | 5.4 | 3.9 | 1.7 | 0.4 | Iris-setosa |

Code:

```
url = " https: // archive.ics.uci.edu/ml/machine-learning-
                databases /iris/ iris-data"
col names = [" sepal-length in cm", "Sepal width in cm",
             " petal length-in-cm", "petal-width in cm",
             "class"]

iris-data = pd. read.csv (url, names = col-names)
```

output.

| | Sepal length in cm | Sepal width in cm | petal length in cm | petal width in cm | class |
|---|---|---|---|---|---|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | iris-setosa |

iris-data.to-csv(" cleaned. iris-data.csv")

1. Get the Data

Download the Data

import os
import tarfile
import urlib
Use these libraries to download /extract
web-based dataset that is housing dataset
Retrieve the date into "housing -tg2
load the data into housing.csv

housing.head()
housing.info()
housing.describe()

→ Using these commands, inspect the attributes
of the dataset

import matplotlib.pyplot as plt
import seaborn as sns

housing.hist (bins=50, figsize = (20,15))
plt. show()

→Using matplotlib & seaborn libraries detecting
the outliers

Create

Create Test Set

→ Splitting the dataset on test ratio = 0.2, i.e., training data is 80% of dataset and testing data is 20% of dataset

→ Stratified sampling is when random chosen data are representation of a whole target population. Each homogenous subgroup is called Strata.

Discover and Visualize the Data to gain insights

→ Visualize the data using matplotlib and seaborn libraries

→ Calculating the standard correlation coefficient of every pair of column.

Prepare the data for Machine learning algorithm

Data cleaning, handling text and categorical data, Custom Transformers, feature scaling, transformation pipelines etc, are done here

Select and Train model

→ At first linear regression model is used to train but the model is overfitting the data.

→ To Tackle this, DecisionTreeRegressor Model is used as it is capable of finding non-linear relationships within the data.

But, the decision tree model is also overfitting so badly that it performs worse than the linear regression model.

At last Random forest regress model is used. It is much better.

fine tune your model:

→ At last fine tuning of model is carried out, evaluating on the test set and then launch, monitor and maintaining the system.

→ Evaluate your system on the test set by using mean-squared error method.

Launch, Monitor & Maintain your system

We can automate this process by
- Collecting fresh data regularly and labeling it
- Writing script to train model and fine tune the hyper parameters
- Writing script to evaluate the model.

Python Implementation of Linear Regression

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


def estimate_coef(x,y):
    n = np.size(x)
    mx = np.mean(x)
    my = np.mean(y)
    ss_xy = np.sum(y*x) = n*my*mx
    ss_xx = np.sum(x*x) - n*mx*mx
    b_1 = ss_xy / ss_xx
    b_0 = my - b_1*mx
    return (b_0, b_1)


def plot_regression_line(x,y,b):
    plot.scatter(x,y, color = "m", marker = "o", s=30)
    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])
    b = estimate_coef(x,y)
    print(b)
    plot_regression_line(x,y,b)
```

Output:

$(b_0, b_1) = (1.2363, ..., 1.16969, ...)$

# Multiple Linear regression

```
from sklearn.model-selection import train.test split
import matplotlib.pyplot as plt
import numpy as np
from sklearn import datasets, linear model, metrics


data_url = " http://lib.stat.cmu.edu/datasets/boston"
raw-df = pd.read.csv (data_url, sep = "\s+", skiprows=22
                      header = None)


x = np.hstack((raw_df.values [::2, 1], raw_df.values[1::2,:2])
y = raw_df.values [1::2, 2]


x_train, x_test, y_train, y_test = train_test split (x,y.
                test_size =0.4, random_state =1)


reg = linear_Model.LinearRegresion()
reg.fit( x_train, y_train)


print("Coefficients :", reg.coeff)
print (" Variance Score :{}", format( reg_score (x_test,
                                                  y_test))

plt_style.use ("fivethirtyeight")


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")
```

```python
plt.lines (y=0, xmin=0, xmax=50, linewidth=2)

plt.legend (loc= 'upper right')
plt.title ("Residual error")
plt.show ()
```

Implementation of ID3

```
import numpy as np
import pandas as pd
eps = np.info(float).eps
from numpy import log2 as log


from google.colab import drive
drive.mount('/content/drive')
Path = 'drive / My Drive / ml datasets / PlayTennis.csv'
df = pd.read_csv(Path)


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


def average_information(df, attribute):
    target = df.keys()[-1]
    target_variables = df[target].unique()
    variables = df[attribute].unique()
    entropy2 = 0
    for variable in variables:
        entropy = 0
        for target_variable in target_variables:
            num = len(df[attribute][df[attribute] == variable]
                      [df[target] == target_variable])
```

```python
den = len(df[attribute][df[attribute] == variable])
fraction = num / (den + eps)
entropy += - fraction * log(fraction + eps)
fraction2 = den / len(df)
entropy2 += - fraction2 * entropy
return abs(entropy2)


tree = buildTree(df)


import pprint
pprint . pprint(tree)
```
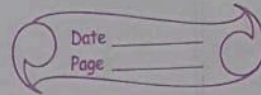
Output: {'outlook': {'overcast': 'yes',
                     'rainy': {'windy': {False: 'yes',
                                         True: 'no'}},
          'sunny': {'humidity': {'high': 'no',
                                 'normal': 'yes'}}}}

09.05.·Y

## Logistic Regression

```
import pandas as pd
from matplotlib import pyplot as plt
from sklearn.metrics import accuracy-score
Path = 'drive/MyDrive/mldatasets/insurance_data-csv's
df = pd.read-csv (Path)


from sklearn.model selection import train_test_split
from matplotlib import pyplot as plt
% matplotlib inline


plt.Scatter (data ['CET score'] vs data ['Admitted'].
     marker = '.', color = 'purple')
x_train, X_test, y_train, y_test = train.test_split (data.
     [['CETSCORE']], data ['Admitted'], train_size = 0.8)


from sklearn.linear-model import LogisticRegression
model = LogisticRegression ()
model.fit (x train, y_train)
y_predicted = model.predict (x test)
model.Score (x test, y_test)
print (y_predicted)
print (Xtest)


from sklearn.linear model import LinearRegression
model = LinearRegression ()
model.fit (x train, y train)
print(" coefficient cm) : model.coef.)
print (" Intercept (b): ", model.intercept_)
```

output:

# KNN Implementation

```python
import numpy as np
import pandas as pd

from google.colab import drive
drive.mount('/content/drive')
dataset = pd.read_csv('/content/drive/MyDrive/iris.csv')
dataset.head()


dataset.groupby('species').size()


feature_columns = ['sepal length', 'sepal width', 'petal length',
                   'petal width']


X = dataset[feature_columns].values
y = dataset['species'].values


from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
y = le.fit_transform(y)
from sklearn.model_selection import train_test_split


X_train, X_test, y_train, y_test = train_test_split(X, y,
                              test_size=0.2, random_state=0)


import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline


from sklearn.neighbours import KNeighbourClassifier
```

```python
from sklearn.metrics import confusion-matrix, accuracy
                                                      score
from sklearn.model_selection import cross-val-score


classifier = KNeighbors Classifier ( n-neighbors = 3)
classifier.fit (X_train, y-train)


y_pred = classifier.predict (x_test)
accuracy = accuracy-score (y_test, y-pred) * 100
print ('Accuracy of our model is equal +
          str (round (accuracy, 2)) + ' %')
```

Output.
Accuracy of our model is equal 96.67%

## SVM

```
-from    sklearn.datasets import load-breast cancer
import   matplotlib.pyplot  as plt
from     sklearn.inspection  import  Decision Boundary Display
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)


Decision Boundary Display . from_estimator (
        svm,
        x,
        response_method = " predict",
        cmap = plt.cm. Spectral,
        alpha = 0.6,
        xlabel = cancer. feature-names[0],
        ylabel = cancer. feature_names[1],
    )


plt. scatter (x [:, 0], x[:, 1],
        c=y,
        s = 20, edgecolors="k")
    plt-show
```

### K-Means Clustering

```
import pandas as pd
import numpy as np
import seaborn as sns
import motplotlib
import sklearn . datasets import load.iris
import sklearn. cluster import KMeans


x, y = load. iris (return. x.y = True)


kmeans =  KMeans (n.cluster = 3 , random.state = 2)
kmeans. fit (x)


kmeans. cluster. centers


pred = kmeans - fit. predict (x)
pred
```

# Dimensionality Reduction using PCA.

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler

from sklearn.datasets import load_breast_cancer
data = load_breast_cancer()

print(data['target_names'])
print(data['feature_names'])

df1 = pd.DataFrame(data['data'],
                   columns = data['feature_names'])

scaling = StandardScaler()
scaling.fit(df1)

scaled_data = scaling.transform(df1)

principal = PCA(n_components = 3)
principal.fit(scaled_data)

x = principal.transform(scaled_data)
```
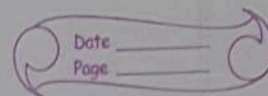
Implement Random Forest ensemble method.

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import sklearn
from sklearn import datasets
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score,
confusion_matrix, classification_report


iris = datasets.load_iris()
iris_data = pd.DataFrame({
    'sepal length' : iris.data[:,0],
    'sepal width' : iris.data[:,1],
    'petal length' : iris.data[:2],
    'petal width' : iris.data[:3],
    'Species' : iris.target
})


X = iris_data.iloc[:, :-1].values
y = iris_data.iloc[:, -1].values


X_train, X_test, y_train, y_test = train_test_split(X, y,
                        test_size = 0.33, random_state = 0)

model = RandomForestClassifier()
model.fit(X_train, y_train)
```

accuracy-score (y-tut, y-pred)

Output : 0.56

## Ada Boost

```
from sklearn.datasets import load_iris
from sklearn.model_selection import train-tut-split
from sklearn.ensemble import AdaBoostClassifier
from sklearn.metrics import accuracy-score

iris = load.iris()

X = iris.data
y = iris.target

x-train, x-tut, y-train, y-tut = train-tut-split (x, y,
              tut-size = 0.4, random-state = 32)

adaboost-clf = AdaBoostClassifier (n-estimators-30,
              learning-rate = 1.0, random-state = 42)

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

y-pred = adaboost-clf.predict (x-tut)

accuracy = accuracy-score (y-tut, y-pred)
print("Accuracy:", accuracy)
```

Output:
   Accuracy : 0.96667

## ANN Implementation

```python
import numpy as np
x = np.array(([2,9], [1,5], [3,6]), dtype = float)
y = np.array(([92], [86], [89]), dtype = float)
x = x/np.amax(x, axis=0)
y = y/100


epoch = 5000
lr = 0.1
input layer neurons = 2
hidden layer neurons = 3
output neurons = 1


wh = np.random.uniform(size = (input layer neurons,
                    hidden layer neurons))
bh = np.random.uniform(size = (1, hidden layer neurons))
wout = np.random.uniform(size = (hidden layer neurons,
                    output neurons))
bout = np.random.uniform(size = (1, output neurons))


def sigmoid(x):
    return 1/(1+np.exp(-x))


def derivatives_sigmoid(x):
    return x*(1-x)


for i in range(epoch):
    hinp1 = np.dot(x, wh)
    hinp = hinp1 + bh
```

```
outinp1 = np.dot (hlayer_act, wout)
outinp  = outinp1 + bout
output  = sigmoid (outinp)


EO      = y - output
outgrad = derivatives_sigmoid (output)
d_output = EO * outgrad
EH      = d_output.dot (wout.T)


hiddengrad = derivatives_sigmoid (hlayer_act)
d_hiddenlayer = EH * hiddengrad


wout += hlayer_act.T.dot (d_output) * lr
wh   += x.T.dot (d_hiddenlayer) * lr


print ("Input : \n" + str(x))
print ("Actual Output : \n" + str(y))
print ("Predicted output : \n" + output)
```

Input
```
[[0.66667   1. 1. ...]
 [0.33333   0.55556]
 [1. ...    , 0.666667]]
```

| Actual Output | Predicted output |
| --- | --- |
| [[0.92] | [[0.86729245] |
| [0.86] | [0.8451565] |
| [0.89]] | [0.8640413]] |