## Lab-1

1) Reading data from URL:

Code:

```
import pandas as pd
url = "https://archive.ics.uci.edu/ml/ml-database/
        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)
iris_data.head()
```

O/p =

| | sepal_length_in_cm | sepal_width_in_cm | petal_length_in_cm |
|---|---|---|---|
| 0 | 5.1 | 3.5 | 1.4 |
| | 4.9 | 3.0 | 1.4 |
| 2 | 4.7 | | |

| | petal_width_in_cm | class |
|---|---|---|
| | 0.2 | iris_setosa |
| | 0.2 | iris_setosa |

2) Reading data from csv file

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

O/p =

| | sepal_length_in_cm | sepal_width_in_cm | petal_length_in_cm | petal_width_in cm |
|---|---|---|---|---|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 |

1

## Lab-2

1) Look at the big picture.
=> Features & columns of dataset must be maintained.
=> Frame problem and the questions to be addressed.
In the process, regression task is also considered.

-) Performance task
-) Regression problems need root mean squared error.
If there are many outliers mean absolute error
is used.

2) Get the data
It is preferable to create some util function to
automate the process of downloading/extracting
web-based data sets

Contents of URL are ~~extracted~~ extracted into housing path
Histogram for all numerical data points are plotted
against the count

Testing data set is sampled from the above (20%)
& remaining 80% goes for training

## Lab-3

→ Create data set

80% of the data set is taken as training set and 20% is taken as testing testing set.

Considering test if program is run again it generates different test set

The test data are always preferred to be hidden. It is scalable & extendable.

→ Discover & Visualize the data

For visualizing training set is analysed. It is preferred to make copy of training set.

Examining dataset since we have latitude & longitude info, map visualisations is done

Various visualisations parameters like colour, alpha value, figsize, kind are used to make the visualisation look better.

→ Looking for Correlation.

Various attributes will be inter-related or vary with each other in a pattern. This is given by correlation. Correlation efficient lies between -1 to 1 when coefficient is close to be 1 It means there is strong (+)ve correlation b/w variable.

→ Prepare data for Ml- Algorithms.

Data cleansing

It is initial step where the missing values, Null values etc are handled. Some

→ Select and train model

=) Linear regression model can be used to train but it can overfit the data.

-) Hence we can use decision tree regress or model is used cause it is capable of finding non-linear relationships within the data but it can be overfitting and can perform worse than linear regression model. To prevent this, Random Forest regression is used.

-) Tuning your Model
-) The model is fine tuned evaluated in test set and launch, monitor and maintaining system.

-) Using mean squared error method, evaluate your system on a test set.

Launch, Monitor & Maintain
-) Collect fresh data regularly and label it.
-) Write script to train model & tune its parameters
-) Write script to evaluate the model

Python's implementation of Linear Regression

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

def estimate_cost(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)
    plt.scatter(x,y, color="m", marks="0", 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_cost(x,y)
    print(b)
    plot_regression_line(x,y,b)


o/b = (b_0, b_1) = (1.236... 1.64...)
```

# Multiple Linear Regression.

```python
from sklearn.model.selection import train_test_split
import matplotlib.pyplot
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
    = 21, 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.linear_regression()
reg.fit(x_train, y_train)

print("coefficents=", reg.wilt)
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")
```

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

plt.legend (label = 'upper right')
plt.title ("suicided error")
plt.show()

# Lab 5

1) Implementation of ID3.

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

```
from google colab import drive.
drive.mount ('/content/drive')
Path = 'drive/My drive/ml datasets/Play Tennis.csv
```

O/p = Mounted at /content/drive.

|   | outlook | temp | humidity | windy | play |
|---|---------|------|----------|-------|------|
| 0 | sunny   | hot  | high     | windy | play |
| 1 | sunny   | hot  | high     | false | no   |
|   | sunny   | hot  | high     | true  | no   |

=) print (f'Rows : {df.shape[0]}, Columns: {df.shape[1]}')

O/p = Rows = 14, Columns = 5

=) print(df columns)

o/p = index (['outlook', 'temp', 'humidity', 'windy', 'play'],
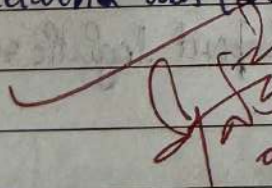
=) df.info() ; df.describe()

```
def 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]]
    return entropy
```

```python
def aug_info (df, attribute):
    target = df.keys()[-1]
    target_variables = df[target].unique()
    variables = df[attributes].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])
            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)
```

09.05.24

# Lab 6

## Logistic Regression

```
import pandas as pd
from matplotlib import pyplot as plt
from sklearn.metrics import accuracy
dt.pd.read_csv (path)

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

plt.scatter (data['Score'], data['Admitted'].
    marker: ':', colors: 'purple')
x_train, x_test, y_train, y_test = train_test_split
    (data['SCORE']], data['Admitted'], size = 0.8)

from sklearn linear_model import logistic regression.
    model = LogisticRegression ()
    model.fit (x_train, y_train)
    y_predicted = model.predict (x_test)
    model.score (x_test, ytest)
    print (y_predicted)
    print (x_test)

from sklearn linear model import Linearregression
    model = Linear Regression()
    model.fit (x_train, y_train)
    print ( "Co-ifficient (m)" : model.coit)
    print ("Intercept (b)", model.Intercept_)

O/p : Prediction : 0.999
```

# K-NN implementation:

```
import numpy as np
import pandas as pd

from google.colab import drive
drive.mount('drive')
ds = pd.read_csv('iris.csv')
ds.head()

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

feature.columns = ['Sepal length', 'sepal width', 'petal length',
                   'petal width']

X = ds[feature_columns].values
Y = ds['Species'].values

from sklearn.preprocessing import label Encoder
    le = label Encoder()
    Y = le.fit.transform(Y)
from sklearn.model_selection import train.test_split

x_train, x_test, y_train, y_test = train.test(X, Y, test_size=0.2,
                  random_state=0)

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

from sklearn.neighbours import kNeighbourClassify
```

```python
from sklearn.metrics import confusion_matrix,
    accuracy_score
from sklearn.model_selection import cross_val_score

classifier = kNeighbours.Classifier (n_neighbours = 3)
classifier.fit (X_train, Y_train)

Y_pred = classifier.predict (x_test)
accuracy = accuracy_score (y_test, y_pred)*100
print ("Accuracy", str (round (accuracy, 2)) + '%')
```

O/p = Accuracy = 96.67 %

Lab-7

1) SVM
-)

```python
import numpy as np
import matplotlib.pyplot as plt
from sklearn import datasets
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.svm from SVC
from sklearn metrics import accuracy_score


iris = datasets.load_iris()
X = iris.data (:, :2)
y = iris.target


X_train, X_test, Y_train, Y_test = train_test_split.
                        (X, Y, test_size = 0.3, state=42)


scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
svm = SVC(kernel = 'linear', random_state=42)
svm.fit(X_train, y_train)


y_pred = svm.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
```

o/p = ~~0.9733~~ 0.977 / 97.7%.

2) PCA

```python
→ import numpy as np
import matplotlib.pyplot as plt
from sklearn import datasets
from sklearn.model selection import StandardScaler
from sklearn.preprocessing import PCA

iris= dataset.load_iris()
X= iris.data()
V= iris.target()

X_train, Y_test, X_test, Y_train= tt.split (test_size = 0.3,
                  stati = 42)

scalar = StandardScalar()
X_train_std = scalar.fit_transform (X_train)
X_test_std = scalar.transform (X_test)

pca= PCA (n components = 2)

X_train_pca = pca.fit_transform (X_train_std)

X_test_pca = pca.transform (X_test_std)

random_forst = randomForestClassifier (n_estimators = 100,
              random_stati = 42)

random_forst.fit (X_train_pca, y_train)

y_pred= random_forest.predict (X_test_pca)
accuracy= accuracy_score (y_test, y_pred)
```

```
print ("Accuracy", accuracy)
```

o/p = A = 95.5

3) K-mean Cluster

```
:) import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib
import sklearn datasets
import local Ris
import sklearn. cluster
import Ate kMeans.
```

```
X, Y = load_iris (return X_Y = True)
KMean = kMeans (n.clusters = 3, random stat = 2)
kMean = fit (x)
KMean = cluster_center
```

1) Implement Random forest method

-)

```
import pandas as pd
from sklearn.model_selection import train test split
from sklearn.ensemble import Random Forest Classifier
from sklearn.metrics import accuracy_score

titanic_data = pd.read csv('url')

print (titanic_data ())

titanic_data.drop (['Name', 'Ticket', 'Cabin', 'Embarked'], axis=1,
    inplace=True)
titanic_data ['Sex'] = titanic_data ['Sex'].map ({'male': 0, 'female':
titanic_data.fillna (titanic_data.mean(), inplace=True)
X = titanic_data.drop ('Survived', axis = 1)
Y = titanic_data ('Survived')

X_train, X_test, Y_train, y_test = train test split (X, y, test size=0.2
    random_state = 42)

random_forest = Random Forest Classifier (n_estimators =100, random_state=
    42)
random_forest.fit (X_train, Y_train)

y_pred = random_forest.predict (X_test)

accuracy = accuracy_score (y_test, y_pred)
print ("Accuracy:", accuracy)
```

o/p: 80%

2] Implement boosting Method.

c) 
```
import numpy as np
import pandas as pd
from sklearn.datasets import load wine
from sklearn.model selection import train_test split
from sklearn.ensemble import AdaBoost Classifier
from sklearn.metrics import accuracy score.


data = load wine()
df = pd.Data frame (data.data, columns = data.feature_names)
df ['target'] = data.target
print (df.head())
X = df.drop ('target', axis = 1)
y = df ['target']
x train, x test, y train, y test = traintest_split (X, Y, test_size 0.2,
random state: 42)


adaBoost = AdaBoost Classifier (n estimates =50, random_state = 42)
adaBoost. predict (x+b)


accuracy = accuracy score (x+iytest, y pred)
print ("Accuracy", accuracy)
op = 92%
```

30.05.2024