# LAB-1

→ Importing and exporting data using Pandas library functions.

① Importing a CSV file using read-csv() function

```
import pandas as pd
data = pd.read_csv("C:\\Users\\Admin\\cleaned-iris-data.csv")
data.head()
```

output:-

|  | Unnamed: 0 | sepal-length-in-cm | sepal-width-in-cm | class |
| --- | --- | --- | --- | --- |
| 0 | | | | |
| 1 | | | | |
| 2 | | | | |
| 3 | | | | |
| 4 | | | | |

② Reading data from URL

```
import pandas as pd
url = "https://....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()
```

output:-

| sepal-length-in-cm | sepal-width-in-cm | petal-length-in-cm | petal-width-in-cm | cl |
| --- | --- | --- | --- | --- |
| 0 | | | | |
| 1 | | | | |
| 2 | | | | |

⑧ Exporting the Data frame to a CSV file

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

Output - Exports the file to the current working directory

Ava
7/5/24

# WEEK-2

## End to End machine Learning Project

Main steps

① Framing the problem & looking at the big picture
② Get the data
③ Explore the data
④ Data preparation
⑤ Shortlisting promising models
⑥ Fine tuning the model & combine them into a great soln
⑦ Present the solution

---

①

The data includes features such as
- Population
- median income
- Median housing price for each block grp in california

PROBLEM :- Prediction will serve as an input to model that attempts to increase the company's ROI

Performance measure : RMSE

② Get the Data

Downloading the data:
import os
import tarfile
import urllib

```
DOWNLOAD_ROOT = "https://raw....."
HOUSING_PATH = os.path.join ("data ", "01")
HOUSING_URL = DOWNLOAD_ROOT + "datasets/housing/housing.tga"
```

```python
def fetch_housing_data(housing_url = HOUSING_URL,
                       housing_path = HOUSING_PATH):
    os.makedirs(name = housing_path, exist_ok = True)
    tgz_path = os.path.join(housing_path, "housing.tgz")
    urllib.request.urlretrieve(url = housing_url, filename = tgz_path)
    housing_tgz = tarfile.open(name = tgz_path)
    housing_tgz.extractall(path = housing_path)
    housing_tgz.close()

fetch_housing_data()


import pandas as pd

def load_housing_data(housing_path = HOUSING_PATH):
    data_path = os.path.join(housing_path, "housing.csv")
    return pd.read_csv(data_path)
```

③ Discover & Visualize the Data

strat_train_set.shape, strat_test_set.shape
strat_test_set.reset_index(), to_feather(fname='data/01/
                                strat_test_set.f')
housing = strat_train_set.copy() ; housing.shape

Visualizing Geographical Data
housing.plot(kind = scatter', x = 'longitude', y = 'latitude')
plt.show

# 4. Prepare the Data for Machine Learning Algorithms

```python
housing = strat_train_set.drop("median_house_value", axis=1)
housing_labels = strat_train_set["median_house_value"].copy()
housing.shape, housing_labels.shape
```

## Data Cleaning

```python
from sklearn.impute import SimpleImputer
imputer = SimpleImputer(strategy='median')
housing_num = housing.drop("ocean_proximity", axis=1)
```

# 5. Select & Train Model

```python
from sklearn.linear_model import LinearRegression
lin_reg = LinearRegression()
lin_reg.fit(X=housing_prepared, y=housing_labels)
```

## Cross-Validation

```python
scores = cross_val_score(estimator=lin_reg,
                X=housing_prepared, y=housing_label,
                scoring='neg_mean_squared_error', cv=10)
```

# 6. Fine Tune your Model

```python
param_grid = [
    {'n_estimators': [3, 10, 30], 'max_features': [2, 4, 6, 8]},
    {'bootstrap': [False], 'n_estimators': [3, 10],
                'max_features': [2, 3, 4]}

]
forest_reg = RandomForestRegressor()
grid_search = GridSearchCV(estimator=forest_reg,
                param_grid, scoring='neg_mean_square_error', cv=5, return_train_score=True,
                n_jobs=-1)


grid_search.fit(X=housing_prepared, y=housing_labels)
```

# WEEK 3
## Linear Regression

① df_sal = pd.read_csv (r'C:\Users\STUDENT ... csv')
df_sal.head()

② plt.title ('Salary Distribution Plot')
sns.displot (df_sal ['salary'])
plt.show()

③ plt.scatter (df_sal ['Years Experience'], df_sal ['salary'],
color = 'lightcoral')

plt.xlabel ('Years of Exp')
plt.ylabel ('salary')
plt.box (false)
plt.show()

④ Splitting variables:
x = df_sal.iloc [:, :1]
y = df_sal.iloc [:, 1:]

⑤ x_train, x_test, y_train, y_test = train_test_split ( x, y,
test_size = 0.2, random_state = 0)

⑥ regressor.fit (x_train, y_train)

⑦ Predin:
y_pred_test = regressor.predict (x_test)
y_pred_train = regressor.predict (x_train)

Steps:  ① Import libraries    ② Import data
③ Analyze data    ④ split data
⑤ Predict results    ⑥ Visulize predictions

# Multiple Linear Regression

① df-start = pd.read_csv ('/content/so_startup.csv')
   df-start.head()

② plt.title ('Profit Distribution Plot')
   sns.displot (df-start ['Profit'])
   plt.show()

③ Splitting
   x = df-start.iloc [:, :-1].values
   y = df-start.iloc [:, -1].values

④ x-train, x-test, y-train, y-test = train-train-split (x, y,
                         test-size = 0.2, random-state = 0)

⑤ regressor = LinearRegression()
   regressor.fit (x-train, y-train)

⑥ Predict
   y-pred = regressor.predict (x-test)

   Steps :-
   → Import libraries
   → Import data
   → Analyze data
   → Split into Independent/Dependent variables
   → Predict Results
   → Compare predictions

Decision Tree                         sevitha-bms

① cols = df.columns[0: -1]
   for i in cols:
        sns.boxplot (y = df [i])
        plt.show ()


② dt = Decision Tree Classifier (max-depth = 3, min-samples-leaf =10,
                              random-state = 1)

   dt.fit (x ,y)

            ↘
            x = df.drop ("species", axis = 1)
            y = df ["species"]


③ features = X-columns
   dot-data = export-graphviz (dt, out-file = None,
                        feature-names = features)
   graph = pydotplus.graph-from-dot-data (dot-data)
   image (graph.create-png())


④ dt = Decision Tree Classifier (random-state =1)
   dt.fit (X-train, y-train)
   y-pred-train = dt.predict (x-train)
   y-pred = dt.predict (x-test)
   y-prob = dt.predict-proba (x-test)


⑤ print ("Accuracy of Decision Tree-Train", accuracy-score (y-pred-train,
                                                    y-train))

   print ("Accuracy of Decision Tree-Test", accuracy-score (y-pred,
                                                    y-test))

```
                    ┌─────────────────────────┐
                    │ petal width (cm) <= 0.8 │
                    │   gini = 0.666          │
                    │   samples = 146         │
                    │   value = [47, 49, 50]  │
                    └─────────────────────────┘
                       T /              \ F
         ┌──────────────────┐      ┌──────────────────────────┐
         │ gini = 0.0       │      │ petal width (cm) <= 1.75 │
         │ samples = 47     │      │   gini = 0.5             │
         │ value = [47,0,0] │      │   samples = 99           │
         └──────────────────┘      │   value = [0, 49, 50]    │
                                   └──────────────────────────┘
                                     T /              \ F
           ┌────────────────────────────┐   ┌──────────────────────────┐
           │ petal length (cm) <= 4.65  │   │ sepal length <= 6.25     │
           │   gini = 0.171             │   │   gini = 0.043           │
           │   samples = 53             │   │   samples = 46           │
           │   value = [0, 48, 5]       │   │   value = [0, 1, 45]     │
           └────────────────────────────┘   └──────────────────────────┘
              /             \                   /              \
  ┌──────────────────┐ ┌──────────────────┐ ┌──────────────────┐ ┌──────────────────┐
  │ gini = 0.05      │ │ gini = 0.408     │ │ gini = 0.165     │ │ gini = 0.0       │
  │ samples = 39     │ │ samples = 14     │ │ samples = 11     │ │ samples = 35     │
  │ value = [0,38,1] │ │ value = [0,10,4] │ │ value = [0,1,10] │ │ value = [0,0,35] │
  └──────────────────┘ └──────────────────┘ └──────────────────┘ └──────────────────┘
```

Logistic regression

```python
import pandas as pd
from matplotlib import pyplot as plt
%matplotlib inline

file path = d 'c \Naan ... csv'
df = pd. read_csv (file path)

df.head()

from sklearn.model_selection import train_test_split
plt.scatter (df.age, df.bought_insurance, marks = '', col

x_train, x_test, y_train, y_test = train_test_split (. df :: x::
                    df.bought_insurance

print (x_test)

model = LogisticRegression ()
model.fit (x_train, y_train)
print (x_test)

model = LinearRegression ()
model.fit (x_train, y_train)
print ("Coefficient (m):", model.coef_)
print ("Intercept (b):", model.intercept_)


def sigmoid(x):
    return 1/ (1+ math. exp(-x))
x=0
sigmoid_value = sigmoid (x)
print ("Sigmoid value at x:", x, ":", sigmoid_value)


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

```
def predicted_function(age):
    m = 0.042
    b = -1.53

    z = m * age + b

    y = sigmoid(z)

    return y


predicted_probability = predicted_function(35)
print("Predicted probability of buying insurance for
    age 35: ", predicted_probability)
```

Quas
30-5-2024

## KNN Classifier

```
P1) import os
    for dirname, _, filenames in os.walk ('/kaggle/input'):
        for filename in filenames:
            print (os.path.join (dirname, filename))


df = pd.read.csv ('play-tennis.csv')

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.log 2 (fraction)
    return entropy


def average_information (df, attribute):
    target = df.keys () [-1]
    target_variables = df [target].unique ()
    variables = df [attribute].unique ()
    entropy 2 = 0

    for variable in variables:
        entropy = 0
        for target_variable in target_variables:
            num = len (df [attribute] [df [attribute] == variable)
                    [df [target] == target_variable)
            fraction = num/(den + eps)
            entropy += - fraction * log (fraction + eps)
        $
    return abs (entropy 2)
```

```python
def build_tree(df, tree=None):
    target = df.keys()[-1]

    node = find_winner(df)

    attValue = np.unique(df[node])

    if tree is None:
        tree = {}
        tree[node] = {}

    for value in attValue:

        subtable = get_subtable(df, node, value)
        clValue, counts = np.unique(subtable[target], return_counts=None)


        if len(counts) == 1:
            tree[node][value] = clValue[0]
        else:
            tree[node][value] = build_tree(subtable)

    return tree
```

p2)  __SVM__

```python
plt.scatter(df.age, df.bought_insurance, marker='+', color='red')
x_train, x_test, y_train, y_test = train_test_split(df[['age']],
    df.bought_insurance, train=)


model = LogisticRegression()
model.fit(x_train, y_train)


y_predicted = model.predict(x_test)
model.predict_proba(x_test)
model.score(x_test, y_test)


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


def prediction_function(age):
    z = 0.042 * age - 1.53
    y = sigmoid(z)
    return y
```

ANN model with back propagation

```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

wh = np.random.uniform(size=(input_layer_neurons, hiddenlayer_
bh = np.random.uniform(size=(1, hiddenlayer_neurons))
wout = np.random.uniform(size=(hiddenlayer_neurons, output_neu
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
    hlayer_act = sigmoid(hinp)
    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
```

p2) Random Forest Algorithm

```
iris = load_iris()
x = iris.data
y = iris.target
x_train, x_test, y_train, y_test = train_test_split (x, y,
                        test_size = 0.5, random_state = 32)

rf_classifier = RandomForest classifier()
rf_classifier.fit (x_train, y_train)
y_pred = rf_classifier.predict (x_test)
accuracy = accuracy_score (y_test, y_pred)

classification_rep = classification_report (y_test, y_pred)
```

p3) Adaboost Algorithm

```
iris = load_iris()
x = iris.data
y = iris.target
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_test)
accuracy = accuracy_score (y_test, y_pred)
print (" Accuracy :", accuracy)
```

## K-Means algorithm to cluster

```python
import numpy as np
import pandas as pd
from google.colab import drive
drive.mount('/content/drive')
Path = '/mydrive/ML-dataset/ans.cv'
df = pd.read-cv(Path)
# df

X = df.iloc(:, :-1).values

class kmeans:
    def __init__(self, n-clusters=4):
        self.k = n-clusters

    def fit(self, x):
        self.centroids = x [np.random.choice(len(x), self.k,
                                             replace=False)]
        self.initial-centroids = self.centroids
        self.prev-label, self.labels = None, np.zeros(len(x))
        while not np.all(self.labels == self.prev-label):
            self.prev-label = self.labels = None, np.zeros(len(x))
            self.labels = self.predict(x)
            self.update-centroid(x)
        return self

    def predict(self, x):
        return np.apply-along-axis(self.compute-label, 1, x)

    def compute-label(self, x):
        return np.argmin(np.sqrt(np.sum((self.centroids-x)**2,
                                        axis=1)))

fig = plt.figure(figsize=(8,6))
plt.scatter(x[:,0], x[:,1], c=y)
```

# PCA

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

data = load_breast_cancer()
data.keys()

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

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

plt.figure(figsize =(10,10))
plt.scatter(x[:0], x[:1], c= data = ['target'], cmap = 'plasma')

plt.xlabel('pc1')
plt.ylabel('pc2')
```