

## LAB - 1

### ① Reading data from URL

Code:

```
import pandas as pd
```

```
url = "https://archive.ics.uci.edu/ml/reddit-reviews.csv"
```

```
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
0	5.1	3.5	1.4
1	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

~~Skipped~~ Import pandas as pd

~~Skipped~~ data = pd.read\_csv("iris-data.csv")

```
data.head()
```

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

## LAB - 2

- look at the big picture.
- the features and columns of dataset must be examined.
- frame the problem statement and questions to which will be addressed in the next part. In this problem regression
- regression is considered
- Regress the Performance measure.
  - In regression problem root mean squared error is chosen as performance measure.
- There are many outliers mean absolute error is used.

### 3. Get the Data.

- It is preferable to create some util functions to facilitate the process of downloading / extracting web-based data sets.
- Links of URL are extracted into housing

Regression  
Dataset for all numerical data points are present.

Count

Data set is

Imaginary 80% sampled from the above (20%)  
you for training.

See  
Notes

## LAB - 3.

Create Data Set.

80% of the data set is taken as training set and 20% is taken as testing set.  
 Considering test data of program, we again, it generates different test set.  
 The test data are always preferred to be hidden. It is scalable and understandable.

## 3. Discover &amp; Visualize the Data.

For visualising training set is analyzed. It is preferred to make copy of training set.  
 Examining dataset since we have latitude and longitude information, map visualisation is done.  
 Various visualizations parameters like colour, alpha value, figsize, kind are used to make the visualisation look better.

Looking for correlation.

Various attributes will be interrelated or vary with each other in a pattern. This is given by correlation. Correlation coefficient lies between -1 to 1. When the coefficient is close to 1 it means there is strong positive correlation b/w 2 variables and vice versa.

## 2. Prepare Data for ML Algorithms.

- Data cleaning.

If it is the initial step where the missing values, null values etc are handled. Some

data which are not in format are also removed or transformed. Null values can be removed or filled using several methods like backward fill, forward fill or interpolation.

Handling Text & Categorical Attributes.  
Dataset has only 2 categorical data. Grouping can be used for this.

Feature Scaling.

ML steps don't perform well on data with different scales. The 2 ways of scaling are:

- i) Min Max Scaling.
- ii) Standardization

### 5. Select and Train Model.

The selected and cleaned data is used to train model. first linear Regression Model is used on the training dataset.

The overfitting and underfitting is handled during training.

Ridge model and decision tree model is applied.

RFECV model fits better. The model is saved after training so that it can be used anytime. joblib is used for storing.

### 6. Fine Tune Your Model.

Sch-Learning GridSearchCV.

Along with these RandomizedSearchCV is used for better purpose.

Q. Queso, Elevator & LAB-9

### Linear Regression.

i) Simple linear regression.

It involves predicting a dependent variable based on single independent variable.

ii) Multiple linear regression.

Involves predicting a dependent variable based on multiple independent variables.

$$y = m_x + c$$

↓  
dependent      independent

Code:

```
def estimate_coeff(x, y)
    n = np.size(x)
    m_x = np.mean(x)
    m_y = np.mean(y)
```

$$SS_{xy} = np.sum(y * x) - n * m_y * m_x$$

$$SS_{xx} = np.sum(x * x) - n * m_x * m_x$$

$$b_1 = SS_{xy} / SS_{xx}$$

$$b_0 = m_y - b_1 * m_x$$

return(b\_0, b\_1).

~~def plot\_sug\_line(x, y, b):~~

plt.scatter(x, y, color="m", marker="o", s=80)

$$y_{\text{pred}} = b[0] + b[1] \times x$$

plt.plot(x, y\_pred, color="g")

plt.xlabel('x')

plt.ylabel('y')

def main():

x = np.array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])

y = np.array([1, 3, 2, 5, 4, 3, 9, 10, 12, 14])

b = estimate\_coeff(x, y)

print("Estimated Coef: In B-0 = 2.1  
nb-1 = 3, format).

Lab - 5.

Decision Tree - ID3.

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

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

```
import os
```

```
for dirname, _, filenames in os.walk('kaggle\\input'):
```

```
    for filename in filenames:
        print(os.path.join(dirname, file))
```

```
If > df = pd.read_csv('path')
```

```
df > df.drop(['day', axis=1])
```

```
df.head()
```

```
print("Rows: ", {df.shape[0]}), columns: {df.shape[1]})
```

```
print(df.columns)
```

```
df.info()
```

```
df.describe()
```

Handling missing data

$$\text{Entropy} = -\frac{p_m}{p_m + p_n} \log_2 \left( \frac{p_m}{p_m + p_n} \right) - \frac{n}{p_m + p_n} \log_2 \left( \frac{n}{p_m + p_n} \right)$$

(i) We compute entropy for dataset.

(ii) Take average info entropy for current attribute.

(iii) Calculate gain.

- (iv) Pick highest gain attribute.  
 (v) Repeat until we get desired tree.

Average information:

$$I(\text{Attribute}) = \sum_{p+n} p_i n_i \cdot \text{Entropy}(\text{Attribute})$$

Information gain:

$$\text{Gain} \rightarrow \text{Entropy}(S) - I(\text{Attribute})$$

Building Decision Tree

def buildTree (df, tree = None):

target = df.keys() [-1]

if tree is None:

tree = {}

tree[node] = {}

for value in att\_value:

subtable = get\_subtable (df, node\_value)

if len(counts) == 1:

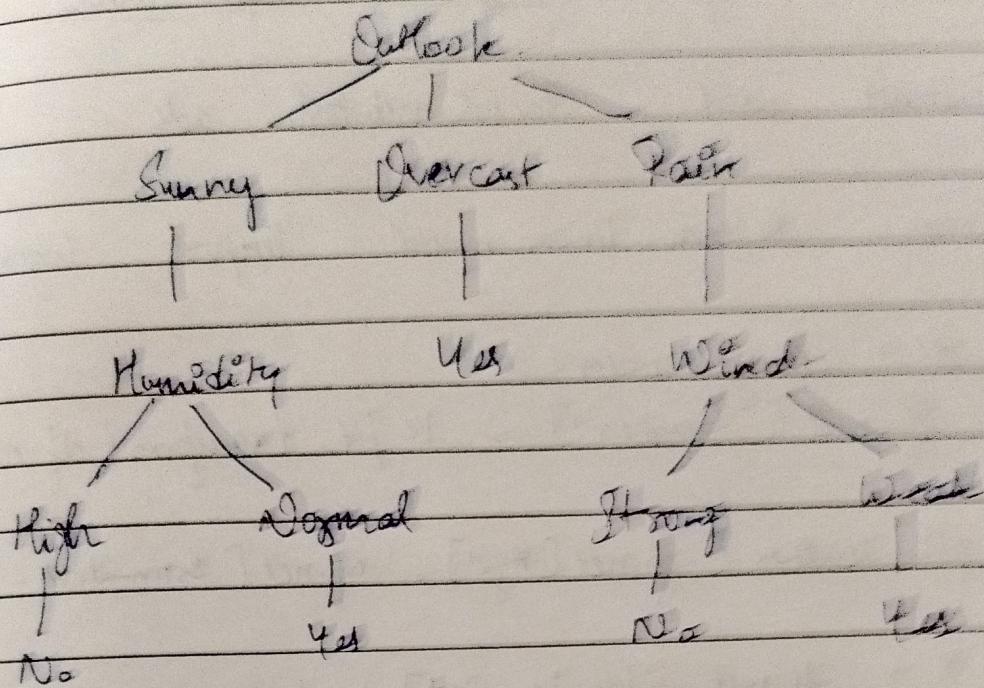
tree[node][value] = cl\_value[0]

else  
 return tree

tree = buildTree (df)

pprint pprint (tree)

## OUTPUT



## Code Output

```
I 'outlook': { 'Overcast': 'Yes', 'Rain': 'False',  
  'Strong': 'No', 'Weak': 'Yes' },  
  'Sunny': { 'Humidity': { 'High': 'No',  
    'Normal': 'Yes' } }
```

949 -

## Lab - 6.

### Logistic Regression and KNN

Dataset used : Social-network-ads

```
from sklearn.linear_model import LogisticRegression
```

```
al = LabelEncoder()
```

```
df['Gender'] = le.fit_transform(df['Gender'])
```

```
plt.scatter(df['Age'], df['Estimated Salary'])
```

```
X = df.iloc[:, :-1].values
```

```
Y = df.iloc[:, -1].values
```

```
X_train, X_test, Y_train, Y_test = train_test_split
```

```
(X, Y, test_size=0.25, random_state=1)
```

```
sc = StandardScaler()
```

```
X_train = sc.fit_transform(X_train)
```

```
X_test = sc.transform(X_test)
```

```
classifier = LogisticRegression(random_state=0)
```

```
classifier.fit(X_train, Y_train)
```

```
y_pred = classifier.predict(X_test)
```

```
accuracy_score(Y_test, y_pred)  
→ 0.82.
```

```
print(f"f1 score: {f1_score(Y_test, y_pred)}")  
→ f1 score: 0.78
```

```
cf_matrix = confusion_matrix(Y_test, y_pred)
```

ans. heatmap (cf-matrix)

# for KNN, till feature scaling procedure is same.

from sklearn.neighbors import KNeighborsClassifier  
as KNN

classifier = KNN(n\_neighbours = 5, metric = "minkowski",  
 $p=2$ )

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

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

accuracy = score(y-test, y-pred)  
→ 0.88

- point (+ "f1 score = f1-score(y-test, y-pred)")  
- f1 score = 0.86

cf-matrix = confusion\_matrix(y-test, y-pred)

ans. heatmap (cf-matrix)

for sample data,

logistic regression, age = 45, salary = 97000  
predicted result = 1 (purchased ads)

KNN, Age = 30, salary = 87000.

predicted result = 0 (not purchased)

Correlation matrix.

for LR:

$y^{\text{pred}}$	0	52	6
	1	11	31
		0	1

$y^{\text{true}}$

for KNN:

$y^{\text{pred}}$	0	50	8
	1	4	38
		0	1

$y^{\text{true}}$

0, 0 → True negative

0, 1 → False positive

1, 0 → False negative

1, 1 → True negative.

Lab - 2

## 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[:, :7]
```

```
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",
```

```
map = plt.cm.Spectral,
```

```
alpha=0.8,
```

```
xlabel=cancer.feature_name[0],
```

```
ylabel=cancer.feature_name[1],
```

```
)
```

```
plt.scatter(x[:, 0], x[:, 1], c=y, s=20, edgecolor="r")
```

```
plt.show()
```

PLA:

```
import pandas as pd
```

```
import numpy as np
```

```
from sklearn.decomposition import PCA
```

```
from sklearn.preprocessing import StandardScaler
```

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

```
scaling = StandardScaler()
```

```
scaling.fit(data)
```

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')

## K-Means.

```
import pandas as pd
```

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

```
data.head()
```

```
import numpy as np
```

```
import seaborn as sns
```

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

```
from sklearn.cluster import KMeans
```

```
X, y = load_iris(return_X_y=True)
```

```
kmeans = KMeans(n_clusters=3, random_state=8)
```

```
kmeans.fit(X)
```

```
y_pred = kmeans.fit_predict(X)
```

```
y_pred
```

Week - 8.

## Ensemble Random forest and ADA Boost

```
from sklearn.datasets import make_moons
from sklearn.metrics import accuracy_score
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier,
                           AdaBoostClassifier, GradientBoostingClassifier
```

$x, y = \text{make_moons}(n\_samples=10000, \text{noise}=.5, \text{random\_state}=0)$

$x\text{-train}, x\text{-test}, y\text{-train}, y\text{-test} = \text{train-test-split}(x, y, \text{test-size}=0.2, \text{random-state}=42)$

```
clf = DecisionTreeClassifier()
clf.fit(x_train, y_train)
y_pred = clf.predict(x_test)
accuracy_score(y_test, y_pred)
→ 0.753.
```

$\text{clf} = \text{RandomForestClassifier}(n\_estimators=100, \text{max_features='auto'}, \text{random-state}=0)$

```
clf.fit(x_train, y_train)
y_pred = clf.predict(x_test)
accuracy_score(y_test, y_pred)
→ 0.7965
```

$\text{clf} = \text{AdaBoostClassifier}(n\_estimators=100)$

```
clf.fit(x_train, y_train)
y_pred = clf.predict(x_test)
accuracy_score(y_test, y_pred)
→ 0.7965.
```

clf = AdaBoostClassifier(n\_estimators=100)

clf.fit(x\_train, y\_train)

y\_pred = clf.predict(x\_test)

accuracy\_score(y\_test, y\_pred)

$\Rightarrow 0.833$ .

Sp. I  
of 16/24