

# **I N D E X**

NAME: Sagar (M) STD.: \_\_\_\_\_ SEC.: \_\_\_\_\_ ROLL NO.: \_\_\_\_\_ SUB.: ML - lab

- 1) Write a python program to import and export data using pandas library function

- 1) Reading data from CSV file

```
import pandas as pd
```

```
airbnb_data = pd.read_csv("/content/sample_data/california_housing_test.csv")
```

```
airbnb_data.head()
```

Output :-

	longitude	latitude	housing_median_age	total_rooms
0	-122.05	37.37	27.0	3885.0

total_bedrooms	population	households	median_income	median_house_value
661.0	1537.0	606.0	6.6085	

- 2) Reading ~~existing~~ data from URL

```
wrl = 'https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data'
```

```
col_names = ['sepal-length', 'sepal-width', 'petal-length',  
            'petal-width', 'class']
```

```
iris_data = pd.read_csv(wrl, names=col_names)
```

```
iris_data.head()
```

- 3) Exporting dataframes to a CSV file

```
iris_data.to_csv("cleaned_iris_data.csv")
```

## Output

	sepal length	sepal width	petal length	petal width	class
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3	1.4	0.2	Iris-setosa

S.P.R.  
21/3/24

Week - 2

## 2 End to End Machine learning Project

## Get the Data

```
import os
import tarfile
import urllib
```

```
download_root = "https://raw.githubusercontent.com/  
ageron/handson-ml2/master/"
```

```
house_path = os.path.join("data", "01")
```

```
housing_url = download_root + "datasets/housing/housing.tgz"
```

```
def fetch_housing_data(housing_data=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)
```

```
housing = load_housing_data()  
housing.head()
```

```
housing.info()
```

```
housing['ocean_proximity'].value_counts()
```

```
housing.describe()
```

```
import matplotlib.pyplot as plt  
import seaborn as sns
```

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

Create a Test set:

```
import numpy as np
```

```
def split_train_test(data, test_ratio=0.2):
```

```
shuffled_indices = np.random.permutation(len(data))  
test_set_size = int(len(data) * test_ratio)  
return data.iloc[:train_indices], data.iloc[test_indices:]
```

~~Train set, test set = split\_train\_test(data=housing)~~  
~~len(train set), len(test set)~~

Get  
20/21

## lab - 3

\* Discover & Visualize the Data to gain Insights

## Visualizing the geographical Data

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

```
housing.plot(kind='scatter', x='longitude', y='latitude',
alpha=0.7, s=housing['population']/100, label='population',
figsize=(10, 7), c='median_house_value',
cmap=plt.get_cmap('jet'), colorbar=True)
plt.legend()
```

## Looking for correlations

```
corr_matrix = housing.corr()
```

## Prepare the Data for Machine learning Algorithm

```
housing = strat_train_set.drop("median_house_value")
housing_labels = strat_train_set["median_house_value"].copy()
```

## \* Data Cleaning

~~housing.dropna(subset=['total\_bedrooms'])~~
~~imputer = SimpleImputer(strategy='median')~~
~~housing\_num = housing.drop(['ocean\_proximity'])~~
~~imputer.fit(housing\_num)~~

# Handling test of Categorical Attribute

housing.cat: housing[[' Ocean\_proximity']]  
housing.cat.head(10)

## \* Select & Train a Model

from sklearn.linear\_model import LinearRegres

lin\_Reg = linearRegression()

tree\_Reg = DecisionTreeRegressor()

forest\_Reg = RandomForestRegressor()

## \* Fine-Tune your Model

forest\_Reg = RandomForestRegressor()

grid\_search = GridSearchCV(estimator=forest\_Reg,  
param\_grid=param\_grid, scoring='neg\_mean\_squared\_error', cv=5, return\_train\_score=True, n\_jobs=-1)

## \* Train Launch, Monitor, & Maintain your sys

We can hot-load the model with in web application or alternatively, wrap the model around its own API end-point and design web component separately

See  
Other

lab - 4

## Simple Linear Regression

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

```
def estimate_slope(x, y)
```

```
n = np.size(x)
```

```
# mean of x, y vector
```

```
# calculating cross-deviation &  
# deviations about x
```

```
# calculating regression coefficient
```

```
def plot_regression_line(x, y, b):
```

```
plt.scatter(x, y, color='m', marker='o', s=30)  
# plotting the reg line  
plt.show()
```

```
def main();
```

```
# observations / data
```

```
# estimating coefficient  
plot Reg line
```

```
if __name__ == "__main__":  
    main()
```

# Multiple Linear Regression

```
Import sklearn.model_selection import  
train_test_split  
import matplotlib.pyplot as plt  
import numpy as np  
import pandas as pd  
from sklearn import dataset, linear_model  
metrics
```

\* Import data from `iris`

- \* `reg = linear_model.LinearRegression()`
- \* calculate coefficient
- \* calculate variance

\* plot the Residual errors

Q &  
Ans.

## Lab 5

## Decision tree ID3      Algorithm Implementation

## Steps

- \* Reading the data sets
- \* Importing packages numpy, pandas
- \* Reading the data set and set the indexes

## Decision tree

ID3 is simple decision tree algorithm

Implement used functions

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

## Average Information

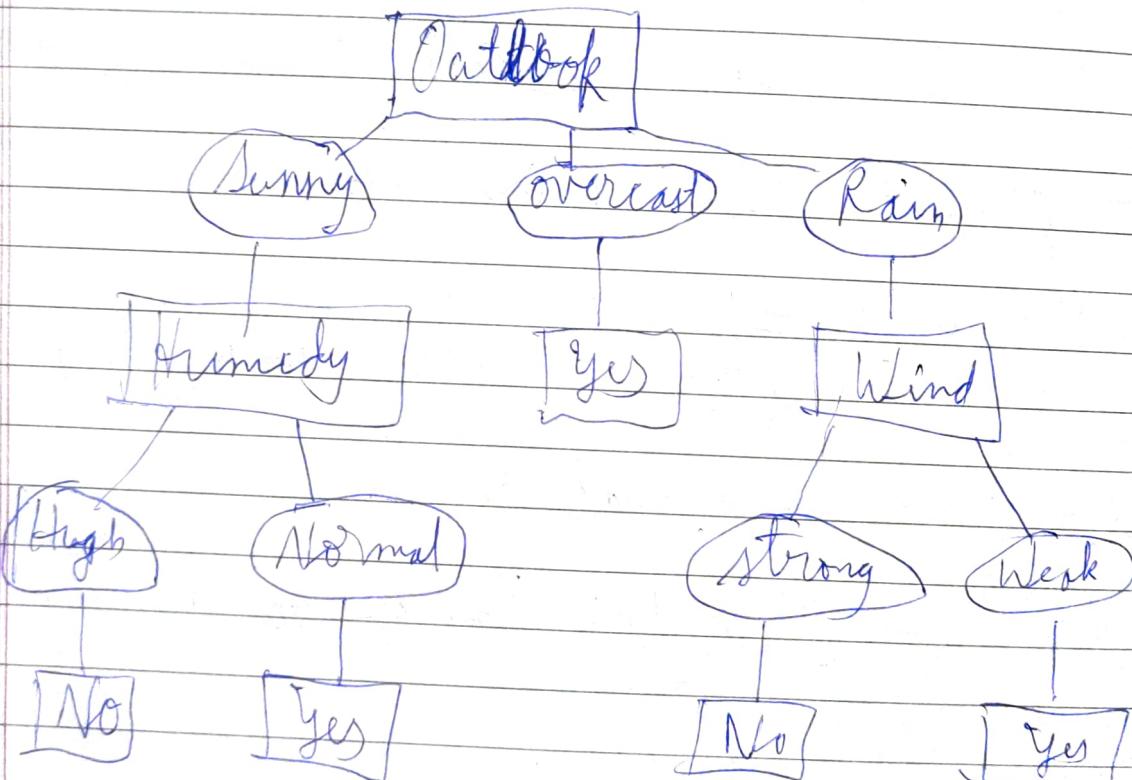
$$I = \sum_{i=1}^{p+n} p_i \text{Entropy}(\text{attribute})$$

## Information Gain

$$\text{Gain} = \text{Entropy}(S) - I(\text{attr})$$

## Build Decision Tree

- \* build tree function construct a decision tree recursively
- \* select the attribute with high information gain
- \* process continue till subsets are pure
- \* print the tree



Sept  
25/9/24

Lab 6

# Logistic Regression

## Key concept

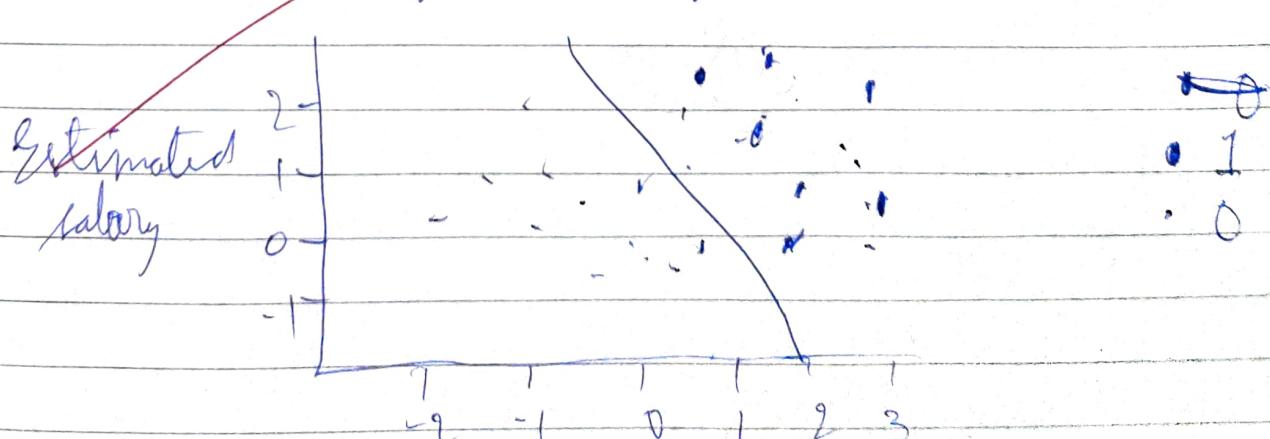
Sigmoid function :- It is a main function that ensures output are between 0 & 1 by converting a linear combination of input data into probabilities.

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

where  $z$  is a linear equation.

## steps

- \* Import the dataset from csv file
  - \* splitting the dataset into training and test set
  - \* splitting feature scaling
  - \* fitting the logistic regression into the training set
- ~~logistic reg~~

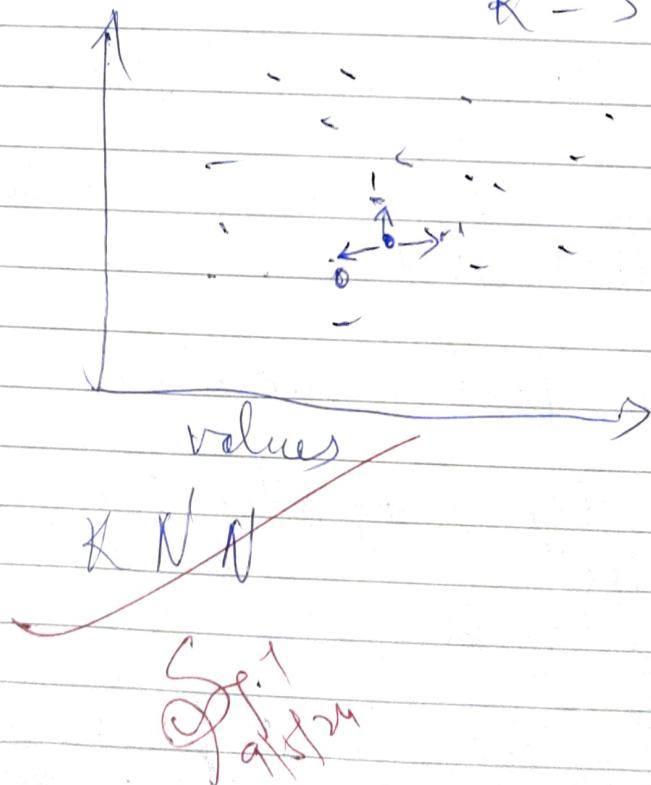


## KNN Model

### Key concept

- \* select the right K value
- \* train the model by feeding the training dataset
- \* select the test set value and find the Manhattan distance between the points
- \* select the min value of k data
- \* majority value containing will be the answer

$$K = 3$$



## Week - 7

## Support Vector Machine

```
import pandas as pd  
from sklearn import dataset  
from sklearn import *
```

```
Iris = dataset.load_iris()  
X = iris.data  
y = iris.target
```

```
X_train, X_test, y_train, y_test = train_test_split  
(X, y, test_size=0.2, random_state=42)
```

```
SVM_model = SVC(kernel='linear')  
SVM_model.fit(X_train, y_train)  
y_pred = SVM_model.predict(X_test)
```

```
accuracy = accuracy_score(y_test, y_pred)  
print("Accuracy of SVM model", accuracy)
```



## K-means Algorithm

- 1) import matplotlib.pyplot as plt  
from sklearn import \*
- 2) data = pd.read\_csv("Iota.csv")  
x = data.values
- 3) kmeans = KMeans(n\_clusters=3, random\_state=0)  
kmeans.fit(x)
- 4) data[["cluster"]] = label  
data.to\_csv('cluster Iota.csv', index=False)
- 5) plt.scatter(x[:, 0], x[:, 1], c=labels, cmap='viridis')  
plt.show()

## PCA

1) import pandas as pd  
import sklearn as sk

- 2) iris = dataset.load\_iris()  
x = iris.data  
y = iris.target

- 3) pca = PCA(n\_components=2)  
x\_pca = pca.fit\_transform(x)

4) `pea_df = pd.DataFrame[ data = X_pca, columns = [ 'principal Component 1', 'principal Component 2' ] ]`  
`pea_df[ 'Target' ] = y`

5) `plt.figure( figsize = ( 8, 6 ) )`  
`plt.scatter( pea_df[ 'PC1' ], pea_df[ 'PC2' ], c = pea_df[ 'Target' ], cmap = 'viridis' )`  
~~`plt.colorbar()`~~  
`plt.show()`

Sp.1  
23/8/24

lab 8

## Random Forest Boosting Implementation

```
import pandas as pd  
from sklearn import *
```

```
iris = load_iris()
```

```
X = iris.data
```

```
Y = iris.target
```

```
X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.2, random_state=42)
```

```
base_estimator = DecisionTreeClassifier(max_depth=3, random_state=42)
```

```
ada_model = AdaBoostClassifier(base_estimator=base_estimator, n_estimators=50, random_state=42)
```

```
ada_model.fit(X_train, y_train)
```

~~```
y_pred = ada_model.predict(X_train)
```~~~~```
acc = accuracy_score(y_test, y_pred)
```~~

# Random Forest Implementation

```
import pandas as pd  
from sklearn import *
```

```
iris = load_iris()  
X = iris.data()  
y = iris.target()
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y,  
test_size=0.2, random_state=42)
```

```
rf = RandomForestClassifier(n_estimators=100,  
random_state=42)
```

```
rf.fit(X_train, y_train)
```

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

```
ac = accuracy_score(y_test, y_pred)
```

Sp1  
20/3/20