

LAB-1

Q1 Write a python program to import and export data using pandas library functions.

→ import pandas as pd.

```
airbnb_data = pd.read_csv("data/listings_austin.csv")  
airbnb_data.head()
```

output1

id	name	host_id	host_name
2265	Ken-retreat	2466	Poddy
5245	Ecocolonial	2466	Poddy
3265	Neighbourhood	7892	Private room
1234	latitude	3234	Executive apt
5298	longitude	7892	Private room

Reading data from URL

```
url = "https://ardrive-ics-uci.edu/full/machine-learning-database/iris/iris-data"
```

```
url_names = ["sepal.length-in-cm", "sepal.width-in-cm",  
             "sepal.length-in-cm", "class"]
```

```
iris_data = url.read_csv(url, names = col_names)  
iris_data.read()
```

→ output

sepal.length-in-cm	sepal.width-in-cm	
0 5.1	3.5	
1 6.9	3	
sepal.length-in-cm	sepal.width-in-cm	class
1.4	0.2	iris-setosa
1.4	0.2	iris-setosa

→ End-to-end project

① Look at the big picture

Dataset consists of California census data to build a model for housing prices in the state

- features:
 - + population
 - + Median Income
 - + Median housing price for each block group in California

Block group is smallest geographical area.

The problem is a supervised learning problem, and because we're predicting median housing price for a district, then we can use regression task.

Performance measure: (Root mean square method).

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (h(x^{(i)}) - y^{(i)})^2}$$

② Get the data:

→ importing the data using pandas

→ `.head()` [loads the dataset into the file representing all the columns & rows]

→ Attribute

- Longitude
- Latitude
- housing-median-age
- total-rooms
- total-bedrooms
- population
- households
- median-income
- median-house-value
- ocean-proximity.

(10 columns,
20,640 rows)

→ describe() function shows all numerical values of dataset for each attribute
eg. count, mean, std, min, max, IQR.
(ignores null values)

→ plotting histograms for each attributes which shows the frequency of distribution over a continuous range of values.

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

→ dataset is divided to 80-20 for training & testing data: [testing data: 20%]

→ housing['id'].value_counts()

→ housing.groupby(by=['longitude', 'latitude']).count()

from sklearn.model_selection import train_test_split

train_set, test_set = train_test_split(housing,
test_size=0.2, random_state=42)

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')  
plt.show()
```

```
alpha = 0.4, s = housing['population'] / 100, (add='population',  
figsize=(10,7), c='median-house-value',  
cmap=plt.get_cmap(name='jet'), colorbar=True)  
plt.legend()
```


④ 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()
```

⑤ Data cleaning:-

```
housing.dropna(subset = ['total_bedrooms'])  
computer = SimpleImputer(strategy = 'median')  
housing_num = housing.drop("ocean_proximity")  
input_fit(housing_num)
```

⑥ Handling text & categorical Attributes.

```
housing_cat = housing[["ocean_proximity"]]  
housing_cat.head(10)
```

* Select & train model.

```
from sklearn.linear_model import LinearRegression
```

```
lin_reg = LinearRegression()
```

```
tree_reg = DecisionTreeRegressor
```

```
forest_rmse = RandomForestRegressor()
```

Find - Tune your model.

```
* forest_reg = RandomForestRegressor
```

```
grid_search = GridSearchCV(estimator = forest_reg,  
                             param_grid = param_grid, scoring = 'neg_mean_squared_error')
```

- Launch, monitor & maintain your system.

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.

WEEK-4 (LAB 4)

18/4/2021

Simple Linear Regression

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

```
import data
df_sal = pd.read_csv('content/Salary-Data.csv')
df_sal.head()
```

Analyze data

```
df_sal.describe() // mean, count, std, min, max,
                    IQR.
```

Distribution:

```
plt.title("Salary Distribution plot")
sns.distplot(df_sal['Salary'])
plt.show()
```

Scatter plot [Salary vs Experience]

```
plt.scatter(df_sal['Years Experience'], df_sal['Salary'],
            color='lightcoral')
```

```
plt.title('Salary vs Experience')
```

```
plt.xlabel('Years of Experience')
```

```
plt.ylabel('Salary')
```

```
plt.box(False)
```

```
plt.show()
```


Split Data :

X = df.sal.iloc[:, 1] # independent
y = df.sal.iloc[:, 2] # dependent

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

regressor = LinearRegression()
regressor.fit(X_train, y_train)

y_pred_test = regressor.predict(X_test)

y_pred_train = regressor.predict(X_train)

prediction on training test(set)

plt.scatter(X_train, y_train)

plt.plot(X_train, y_pred_train)

plt.title('Salary vs Experience')

plt.xlabel('Years of Experience')

plt.ylabel('Salary')

plt.show()

prediction on test set.

plt.scatter(X_test, y_test)

plt.plot(X_train, y_pred_train)

plt.title('Salary vs Experience')

plt.xlabel('Years of Experience')

plt.ylabel('Salary')

plt.show()

Lab 4

Simple Linear Regression

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

def estimate-coeff(x, y):
    n = np.size(x)
    # mean of x, y vector
    # calculating cross-deviation & deviation about x
    # calculating regression coefficient

def plot-regression-line(x, y, b):
    plt.scatter(x, y, color='m', marker='o', s=20)
    # plotting the reg line
    plt.show()

def main():
    # observation / data
    # estimating coefficient
    plot-reg-line

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

Multiple Linear Regression

```
import sklearn
import sklearn.model_selection
import sklearn.metrics
```

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

```
from sklearn import datasets, linear_model
```

- * Import data from url
- * `reg = linear_model.LinearRegression()`
- * calculate coefficient
- * calculate variance
- * plot the ~~Residual~~ error

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

$$I = \sum \frac{p_i + n_i}{p+n} \text{Entropy}(\text{attribute})$$

Information Gain

$$\text{Gain} = \text{Entropy}(S) - I(A+).$$

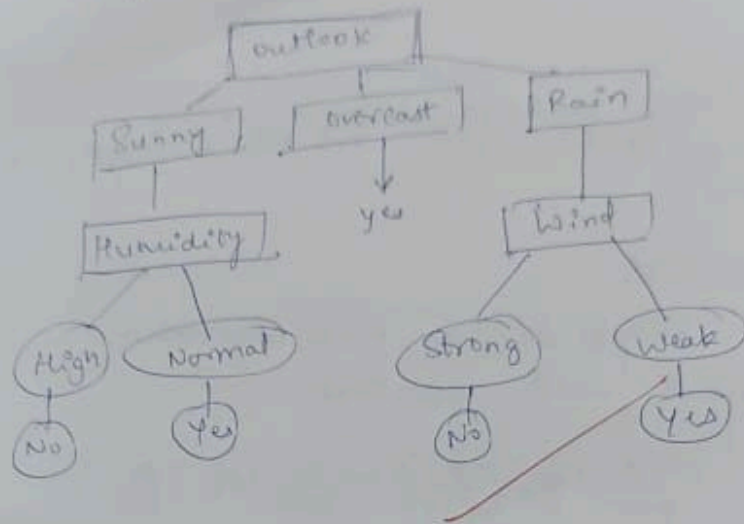
Build Decision tree

Build Tree function construct a decision tree recursively.

Select the attribute with high information gain

process continue till subsets are pure.

split the tree



LAB-6 (KNN)

```
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import train_test_split
from sklearn.datasets import load_iris
// loading datasets
```

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

```
irisData = load_iris()
```

```
# create feature and target arrays.
```

```
x = irisData.data
```

```
y = irisData.target
```

```
# split into training and test set
```

```
x_train, x_test, y_train, y_test = train_test_split(
    x, y, test_size=0.2, random_state=42)
```

```
neighbors = np.arange(1, 9)
```

```
test_acc = np.empty(len(neighbors))
```

```
test_acc = np.empty(len(neighbors))
```

```
# loop over k values
```

```
for i, k in enumerate(neighbors):
```

```
    knn = KNeighborsClassifier(n_neighbors=k)
```

```
    knn.fit(x_train, y_train)
```

```
# compute training and test data accuracy
```

```
train_accuracy[i] = knn.score(x_train, y_train)
```

```
test_accuracy[i] = knn.score(x_test, y_test)
```

```
# generate plot
plt.plot(neighbors, test_acc)
plt.plot(neighbors, train_acc)
```

```
plt.legend()
plt.xlabel("n-neighbors")
plt.ylabel("Accuracy")
plt.show()
```

<u>o/p</u>	<u>Sepal-length</u>	<u>Sepal width</u>	<u>Petal length</u>	<u>width</u>	<u>Species</u>
0	5.1	3.5	1.4	0.2	Iris-Setosa
1	4.9	3.0	1.4	0.2	"
2	4.7	3.2	1.3	0.2	"
3	4.6	2.1	1.5	0.2	"
4	5.0	3.6	1.4	0.2	"

Lab-7

→ SVM (Support vector Machine Model)

(Linear)

→ Importing necessary libraries.

iris = load_iris()

x, y = iris.data, iris.target

x_train, x_test, y_train, y_test = train_test_split(x, y,
test_size=0.2, random_state=42)

Scaler = StandardScaler()

x_train_scaled = Scaler.fit_transform(x_train)

x_test_scaled = Scaler.transform(x_test)

Svm_model = SVC(kernel='rbf', C=1.0, gamma='scale',
random_state=42)

Svm_model.fit(x_train_scaled, y_train)

⇒ SVC

SVC(random_state=42)

y_pred = Svm_model.predict(x_test_scaled)

accuracy = accuracy_score(y_test, y_pred)

~~print~~ print("Accuracy", accuracy)

Accuracy: 1.0.

```
print("\n Classification Report:")
print(classification_report(y_test, y_pred, target_names=
                             iris.target_names))
```

	Precision	recall	f1-score	support
Setosa	1.00	1.00	1.00	10
Versicolour	1.00	1.00	1.00	9
Virginica	1.00	1.00	1.00	11

	accuracy	macro avg	weighted avg
	1.00	1.00	1.00
	30	30	30

implement dimensionality reduction using
principal component analysis (pca) method.

→ Import the library.

→ Load the iris dataset

→ Split the dataset into training & testing sets

→ preprocess the data.

Apply PCA for dimensionality Reduction.

pca = PCA(n_components=2) # Reduce to 2 principal components

X_train_pca = pca.fit_transform(X_train_scaled)

X_test_pca = pca.fit_transform(X_test_scaled)

svm_model = SVC(kernel='rbf', c=1.0, gamma='scale',
random_state=42)

y_pred = svm_model.predict(X_test_pca)

accuracy = accuracy_score(y_test, y_pred)

print("Accuracy with pca:", accuracy)

print(classification_report(y_test, y_pred, target_names=
=iris.target_names))

Accuracy with PCA = 0.9

Classification Report with PCA:

	Precision	recall	f1-score	Support
all avg			0.90	30
avg	0.90	0.90	0.90	30
weighted avg	0.90	0.90	0.90	30

Q Build k-means algorithm to cluster a set of data stored in a csv file

⇒ import numpy as np

class KMeans:

def __init__(self, n_clusters, max_iter=300):

self.n_clusters = n_clusters

self.max_iter = max_iter

def fit(self, x):

centroids_indices = np.random.choice(len(x),
size = self.n_clusters, replace=False)

self.centroids = x[centroids_indices]

for _ in range(self.max_iter):

labels = self._assign_clusters(x)

new_centroids = np.array([x[label == k],
mean(x, axis=0) for k in range(self.n_clusters)])

if np.all(self.centroids == new_centroids):

break

self.centroids = new_centroids

np.random.seed(42)

x = np.random.rand(100, 2)

kmeans = KMeans(n_clusters=3)

~~kmeans.fit(x)~~

label = kmeans._assign_clusters(x)

print("Cluster labels:", label)

LAB-08

0) Random Forest Classifier

Importing necessary libraries

Load the Iris dataset

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_model = RandomForestClassifier(n_estimators = 100,
random_state = 42)

rf_model.fit(X_train, y_train)

y_pred = rf_model.predict(X_test)

evaluate the model

accuracy = accuracy_score(y_test, y_pred)

print("Accuracy of Random Forest model:",
accuracy)

2. Build Logistic Regression Model

```
from sklearn.datasets import load_diabetes.  
from sklearn.linear_model import LogisticRegression.  
from sklearn.model_selection import train_test_split
```

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

```
X_train, X_test, y_train, y_test =  
train_test_split(X, y, test_size=0.2,  
random_state=23)
```

```
clf = LogisticRegression(random_state=0)
```

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

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

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

```
print("Logistic Regression model accuracy %  
acc * 100)
```

Output

~~Accuracy = 95.61~~

2/3/24

— X —

g) Boosting Implementation

import necessary libraries

loading the iris dataset

iris = load_iris()

X = iris.data

Y = iris.target

split the dataset into training & testing sets.

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=1, random_state = 42)

ada_model = AdaBoostClassifier(base_estimator = base_estimator,

n_estimators = 50, random_state = 42)

train the model

ada_model.fit(X_train, y_train)

y_pred = ada_model.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)

print("Accuracy of AdaBoost model:", accuracy)

30/5/2024