

21/10/24

classmate

Date _____

Page _____

Lab - 1

- a. Write a Python program to import and export csv files using Pandas.

A) import pandas as pd

```
df = pd.read_csv('Iris.csv')  
df.head()
```

OUT:

	Id	Sepal Length cm	Sepal Width cm	Petal Length cm
1		5.1	3.5	1.4
		Petal Width cm	Species	
	0.2		Iris-setosa	

```
url = 'https://...'  
colnames = ['Length in cm', 'Width in cm',  
            'Petal Length in cm', 'Petal Width in cm',  
            'Class']
```

```
data = pd.read_csv(url, names=colnames)  
data.head()
```

OUT:

	Sepal Length in cm	Sepal Width in cm	Petal Length in cm	Petal Width in cm	Class
1	5.1	3.5	1.4	0.2	Iris-setosa

data.to_csv('NewIris.csv') //export

Lab - 2

CLASSMATE

Date _____

Page _____

Housing Data

Main steps involved are :

- 1) Defining the problem
- 2) Import the data
- 3) Cleaning of data ; analysing of data
- 4) Selecting an ML model
- 5) Training the Model
- 6) Evaluating the Model's performance
- 7) Deploying the model.

The data has been imported via the url.

3) Cleaning and Description

To describe the dataset, we used the
• `describe` and `head` function
• `info` method can be used to get info on the
dataset

The dataset has 10 attributes. and 20,640 instances.

The null values in each column can be found out using `isnull()` method to get count of NAN values.

The data types per column is also visible.

Visualisation:

By making use of `matplotlib` library histograms on each column can be drawn to

display the continuous data.

For example

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

This makes use of hist function from seaborn library.

4) Creating Test Data

```
shuffled_indices = np.random.permutation(len(data))  
test_size = int(len(data) * test_ratio)  
test_indices = shuffled_indices[:test_size]  
train_indices = shuffled_indices[test_size:]
```

```
return data.iloc[train_indices], data.iloc[test_indices]
```

~~This divides the data set into test and train sets with 20% of data in test dataset. The indices are selected randomly.~~

Sp.1
20% 2M

Lab - 3

classmate

Date _____
Page _____

- * Describe and visualise the data to gain insights.

```
from zlib import crc32
```

```
def test_set_check(identifier, test_nation=.2):  
    total_size = 2**32  
    hex_nepr = crc32(mp.int64(identifier)) & 0xfffffff  
    im_test = hex_nepr < (test_nation * total_size)  
    return im_test
```

```
def split_train_test_by_id(data, test_nation, id):  
    ids = data[id]  
    im_test_set = id.apply(lambda id: test_set_check(id, test_nation))  
    return data.loc[~im_test_set], data.loc[im_test_set]
```

```
train_set, test_set = split_train_test_by_id(data=housing_with_id, test_nation=0.2, id="index")  
train_set.shape, test_set.shape  
(16512, 11), (4128, 11)
```

1) housing = stat_train.set.copy();
housing.shape

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

3) housing[['population', 'median_house_value']].corr()

population	population	median_house_value
population	1.000	-0.026882
median_house_value	-0.026882	1.00000

4) housing.describe()

5) housing.hist(bins=50, figsize=(20,15))

plt.show()

The descriptive statistics of dataset are presented including mean, standard deviation and percentiles of numerical attributes.

• Select & Train Model

```
from sklearn.linear_model import LinearRegression
lin_neg = LinearRegression()
lin_neg.fit(x=housing_prepared, y=housing_labels)
```

```
from sklearn.metrics import mean_squared_error
```

```
housing_predictions = lin_neg.predict(housing_prepared)
```

~~& Tune the model~~

```
from sklearn.model_selection import GridSearchCV
```

~~forest_neg = RandomForestRegressor()~~

```
grid_search.fit(x=housing, y=housing_labels)
```

~~grid_search_cv~~

```
estimator = RandomForestRegressor()
```

Linear Regression

```
import numpy as np
```

```
import matplotlib.pyplot as plt
```

```
def estimate_coeff(x, y):
```

```
    m = np.mean(x)
```

```
    m_x = np.mean(y)
```

```
    ss_xy = np.sum(y * x) - m * m_y + m * m_x,
```

```
    ss_xx = np.sum(x * x) - m * m_x + m * m_x
```

```
    b_1 = ss_xy / ss_xx
```

```
    b_0 = m_y - b_1 * m_x
```

```
    return(b_0, b_1)
```

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

```
    plt.scatter(x, y, color='m', marker='o', 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, 3, 4, 5, 6, 7, 8, 9])
```

```
y = np.array([1, 3, 2, 5, 7, 8, 9, 10, 12])
```

```
b = estimate_coeff(x, y)
```

```
print("Estimate coeff : In b_0 = {} & b_1 = {}")
```

```
format(b))
```

Output:

(b_0, b_1) = (1.2, 36.3, ... 1.16969...)

with imbuil

2 - 2020

```
from sklearn.model_selection import train_test_split  
from sklearn.compose import ColumnTransformer  
from sklearn.preprocessing import OneHotEncoder  
from sklearn.linear_model import LinearRegression
```

The above can work for simple linear as well as multiple linear regression as shown

In multiple linear regression we predict dependent variable based on multiple independent instead of single one.

20/01/2024

classmate

Date _____
Page _____

Week - 5

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))

df = pd.read_csv('path')

df = df.drop('day', axis=1)

df.head()

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

print(df.columns)

df.info()

df.describe()

Handling missing data

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

- 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 :

$$I(\text{Attribute}) = \frac{\sum p_i m_i}{n+m} \cdot \text{Entropy}(\text{Attribute})$$

Information gain

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

Building decision tree

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

Building Decision Tree

```
def buildTree(df, tree=None):
```

```
    target = df.keys()[-1]
```

```
    if tree is None:
```

```
        tree = {}
```

```
        tree[Node] = {}
```

for value in attribute >= value:

```
    subtable = get_subtable(df, mode_value)
```

```
    if len(counts) == 1:
```

```
        tree[mode][value] = class_value[0]
```

```
    else
```

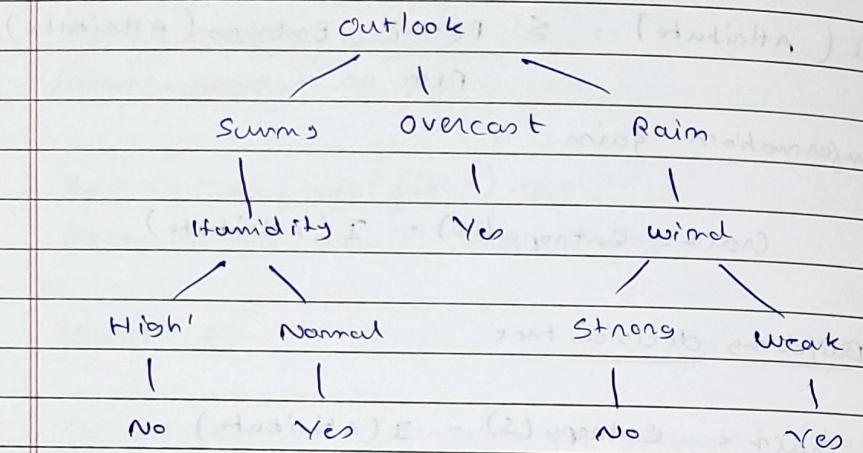
```
        tree[mode][value] = buildTree(subtable)
```

```
    return tree
```

tree = tree = buildTree(df)

print(tree)

02/05/24

OutputCode Output

```

{'outlook': 'Sunny', 'outlook': 'Overcast', 'outlook': 'Rain', 'wind': 'Strong', 'wind': 'Weak', 'humidity': 'Yes', 'humidity': 'No', 'sunrise': 'High', 'sunrise': 'Normal'}
  
```

f58

Lab GLogistic Regression and KNN

Dataset used: Social_Network_Ads

from sklearn.linear_model import LogisticRegression

df = LabelEncoder()

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

plt.scatter(df['Age'], df['EstimatedSalary'])

x = df['Age'].values

y = df['EstimatedSalary'].values

xtrain, xtest, ytrain, ytest = train_test_split

(x, y, size=0.25, random_state)

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)

=> 0.83

print(f"\nF1 score: {f1_score(y_test, y_pred)}")

=> F1 score: 0.78

cf_matrix = confusion_matrix(y_test, y_pred)

sns.heatmap(cf_matrix)

For KNN fit feature scaling procedure is same

from sklearn.neighbors import KNeighborsClassifier as KNN

classifier = KNN(n_neighbors=5, metric='minkowski', p=2)

classifier.fit(x_train, y_train)

y_pred = classifier.predict(x_test)

accuracy = score(y_test, y_pred)

$\Rightarrow 0.88$

print(f" F1 score: {f1_score(y_test, y_pred)}")

$\Rightarrow F1 \text{ score} = 0.86$

cf_matrix = confusion_matrix(y_test, y_pred)
sns.heatmap(cf_matrix)

For sample data

Logistic Regression, age = 30, salary = 87000

Predicted result = 1 (purchased ads)

KNN, Age = 30, Salary = 87000

Predicted result = 0 (not purchased)

Correlation matrix

for LR:	52	6
y-pred	11	31
y-true		

for KNN:	50	8
y-pred	4	38
y-true		

F 105

0,0 → True Negative

0,1 → False Negative

1,0 → False Positive

1,1 → True Positive

Mv2

23/05/24

Lab 7

CLASSMATE
Date _____
Page _____

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)
```

Decision Boundary Display. from_estimator(

```
SVM, X, response_method="predict",
```

```
map = plt.cm.rainbow,
```

```
alpha = 0.8,
```

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

```
ylabel = cancer.feature_names[1];
```

```
)
```

```
plt.scatter(X[:, 0], X[:, 1], c=y, s=20, edgecolor="black")
plt.show()
```

NCA

import pandas as pd

import numpy as np

```
from sklearn.decomposition import NCA
```

```
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 = plt.cm.Paired)`

`plt.xlabel("PC1")`

`plt.ylabel("PC2")`

K-Means

`import pandas as pd`

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

`data.head()`

`import numpy as np`

`import seaborn as sns`

`from sklearn.datasets import load_iris`

`from sklearn.cluster import KMeans`

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

`kmeans = KMeans(n_clusters = 3, random_state = 2)`

`kmeans.fit(x)`

~~`pred = kmeans.fit_predict(x)`~~

30/07/24

Week 8

classmate

Date _____
Page _____

Ensemble Random Forest

```
from sklearn.datasets import make_moons  
from sklearn.metrics import accuracy_score
```

```
from sklearn.model_selection import train_test_split  
from sklearn.tree import DecisionTreeClassifier  
from sklearn.ensemble import RF, ABC, NBC
```

```
x, y = make_moons(n_samples=10000, noise=.5,  
random_state=0)
```

```
X-train, Y-train, X-test, y-test = train_test_split(x, y,  
test_size=0.2, random_state=0)
```

```
c1f = DecisionTreeClassifier()
```

```
c1f.fit(X-train, Y-train)
```

```
y-pred = c1f.predict(X-test)
```

```
accuracy-score(y-test, y-pred)
```

=> 0.753

```
c1f = RandomForestClassifier(n_estimators=100, max_features
```

```
"auto", random_state=0)
```

```
c1f.fit(X-train, Y-train)
```

```
y-pred = c1f.predict(X-test)
```

```
accuracy-score(y-test, y-pred)
```

=> 0.765

S. P. 30/07/24
~~c1f = AdaBoostClassifier(n_estimators=100)~~

~~```
c1f.fit(X-train, Y-train)
```~~~~```
y-pred = c1f.predict(X-test)
```~~~~```
accuracy-score(y-test, y-pred)
```~~~~=> 0.833~~