

21/3/24 A2¹⁻³ Week-1

Import & Export a csv file with pandas in python:

Code:

```
import pandas as pd  
df = pd.read_csv("file-path.csv")  
df.head()
```

OUTPUT :

	Id	Sepal Length	Sepal Width	Petal Length	Petal Width
0	[0]	[0.50]	[0.4]	[1.4]	[0.2]
1	[1]	[0.79]	[0.37]	[1.4]	[0.2]
2	[2]	[1.3]	[0.2]	[4.3]	[1.3]
3	[3]	[1.4]	[0.2]	[3.9]	[1.3]
4	[4]	[1.3]	[0.2]	[4.7]	[1.8]

Export code:

```
df.to_csv("newPath\newName.csv")
```

Kaggle Completed : 3 chapters out of 6

1) Creating, Reading, Writing

→ pd.DataFrame({ 'Yes' : [0,1], 'No' : [2] })

\rightarrow pd.DataFrame ([], index=[])

\rightarrow pd.Series ([], index=[], name= ' ',)

\rightarrow pd.read_csv (" ")

\rightarrow df.new_file.reshape()

\rightarrow new_file.read()

\rightarrow new_file.to_csv(' ')

(*)

2) Indexing, Selecting, Assigning

\rightarrow df.iloc []

\rightarrow df.iloc [[row_index], [col_index]]

\rightarrow df.loc [[row_name/index], [col_name/index]]

\rightarrow df.set_index (" ")

\rightarrow df.specific_col_name = 'name'

\rightarrow df.loc [df.specific_col_name == 'name']

\rightarrow df[df.specific_col_name.isin (' ',)]

\rightarrow l \rightarrow & , & \rightarrow and

\rightarrow not null to make sure that a row is not null.

3) Summary Functions & Maps:

\rightarrow

df.specific_col_name.describe()

* gives count, unique, top & most frequent.

\rightarrow df.number_col.mean()

\rightarrow df.col.unique()

\rightarrow df.col.value_counts()

* To see unique sets & how often they appear.

28/03/24

Week-2 Project End-End ML Project :

3. Data ~~is discovered~~ & Visualize the data to gain insights.

Create the Data :

```

import os
import tarfile
import urllib

DOWNLOAD_ROOT = "https://raw.githubusercontent.com/lasagni/
handom-wiki/master/"

HOUSING_PATH = os.path.join("data", "ol")
HOUSING_URL = DOWNLOAD_ROOT + "datasets/housing/housing.
tgz"

```

```

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()

```

To download data, call the function

```
fetch_housing_data()
```

Latitude / longitude info can be represented using .plot()

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

```

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

```

```
plt.show()
```

```

housing.plot(kind='scatter', x='longitude',
              y='latitude', alpha=0.1)

```

```
plt.show()
```

```

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)

```

To check the relation between two or more data info, 'corr' method is used.

```

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

```

This generates a weak correlation b/w the specified parameters.

For a better correlation, we need to do,

```
housing['corr_matrix'] = housing.corr()
```

```
(corr_matrix['median_house_value'])
```

```
list_values(ascending=False)
```

To download data, call the function

4. Prepare Data for ML Algorithms.

from sklearn.impute import SimpleImputer

imputer = SimpleImputer(strategy='median')

This can handle only numbers

housing_num = housing.drop(['ocean-proximity', 'axis=1'])

Now we fit ~~fit~~ imputer ~~over~~ our data ~~of~~ over

fitting is basically training.

imputer.fit(housing_num)

Now we can use this to transform the

numbers by replacing their missing

values with corresponding medians.

X = imputer.transform(housing_num)

X is a numpy array with the transform -ed features. We add it back to the

data frame.

housing_rf = pd.DataFrame(data=X, index=housing_num.index,

columns=housing_num.columns)

linear regression has better performance but

lower accuracy when compared to Decision tree & Random Forest.

linear regression is represented as $y = mx + c$

5. Select & Train a Model

from sklearn.linear_model import LinearRegression

lin_reg = LinearRegression()

lin_reg.fit(X=housing_prep, y=housing_labels)

some_data = housing.iloc[:5] # lets take

some_labels = housing_labels.iloc[:5] # small amount

of the data

some_data_prep = full_pipeline.transform(some_data)

some_labels.tolist() is not accurate. So, we use another method:

from sklearn.metrics import mean_squared_error.

housing_pred = lin_reg.predict(housing_prep)

lin_mse = mean_squared_error(housing_labels, housing_predictions)

lin_mse = np.sqrt(lin_mse)

this is the root-mean squared error which is more accurate.

"Decision Tree" is a combination of multiple Decision Trees. So, it is better at capturing non-linear relationships.

Week-9

Linear Regression

Without

```
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import sklearn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from random import seed
from sklearn.linear_model import LinearRegression
df_sal = pd.read_csv("content/Sal-Data.csv")
df_sal.describe()
plt.title("Salary Distribution Plot")
sns.displot(df_sal['Salary'])
plt.show()
plt.scatter(df_sal['YearsExp'], df_sal['Sal'],
            color='lightblue')
plt.title('Sal vs Exp')
plt.xlabel('Exp')
plt.ylabel('Salary')
plt.box(False)
plt.show()

Split Data:
X = df_sal.iloc[:, :-1].values
y = df_sal.iloc[:, -1].values
train_size = int(0.8 * len(X))
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)
```

Split into Train test sets:

```
= the train-test-split(x, y, test_size = 0.2,
```

```
random_state = 0)
```

Train Model:

```
regressor = linearRegression()
```

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

Predict Results:

```
y_pred_test = regressor.predict(x_test)
```

```
y_pred_train = regressor.predict(x_train)
```

Visualizing Predictions:

```
plt.scatter(x_train, y_train, 'lightcoral')
```

```
plt.plot(x_train, y_pred_train, color='firebrick')
```

```
plt.title('Salary vs Exp (Training Set)')
```

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

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

```
plt.legend([x_train['pred(y-test)',
```

```
'x_train / y_train'], title='Sal(Exp',
```

```
loc='best', facecolor='white')
```

```
plt.show(Fake)
```

```
plt.show()
```

plt.show()

Coefficient and Intercept

```
print({'Coeff': regressor.coef_3})
```

```
print({'Intercept': regressor.intercept_3})
```

Output:

```
Coeff: [26312.57]
```

```
Intercept: [26780.099]
```

Multiple Linear Regression:

```

import pandas as pd
import numpy as np

import matplotlib.pyplot as plt
import relearn as sns
from sklearn.model_selection import train_test_split as tts.

from sklearn.linear_model import LinearRegression

df_start = pd.read_csv ('laptop_startup.csv')
df_start.head()
df_start.describe()

Distribution:
plt.title ('Project Distribution Plot')
sns.distplot (df_start ['Profit'])
plt.show()

Relation b/w profit & R&D spend
plt.scatter (df_start ['R&D spend'], df_start ['Profit'],
             color='light coral')

plt.title ('Profit vs R&D spend')
plt.xlabel ('R&D spend')
plt.ylabel ('Profit')
plt.box (False)
plt.show()

```

Split into Independent | Dependent variables :

x = $\overline{df_start.iloc[:, :-1].values}$
y = $\overline{df_start.iloc[:, -1].values}$

One-hot encoding :
 $ct = \text{ColumnTransformer} (\text{transformers} = [(\text{encoder},$
 $\text{ct} = \text{ColumnTransformer} (\text{transformers} = [(\text{encoder},$
 $\text{y} = \text{df_start.iloc[:, :-1].values}$
 $\text{OneHotEncoder} \text{ OneHotEncoder} (\text{passthrough} = \text{True}))], \text{remainder} = \text{Id}$

```

x = np.array (ct .fit_transform (x))

```

Split into train | test Data:

```

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

```

Distribution:

```

plt.title ('Project Distribution Plot')
sns.distplot (df_start ['Profit'])
plt.show()

```

Train Model :
 $regressor = \text{LinearRegression} ()$

regressor.fit (X_train, Y_train)

Predict Results :

$\hat{Y}_{pred} = regressor.predict (X_{test})$

Compare Predictions:

```

mp. rint - print options (precision=2)
result = mp. concatenate ([y_pred.reshape (len (y_pred), 1),
                           y_true.reshape (len (y_true), 1)], 1)

```

Week-4

Decision Tree

Implementation

Plot the Tree graph:

`tree.plot_tree(model)`

`zeros = np.sum(y-train == 0)`

`ones = np.count_nonzero(y-train)`

$$\text{val} = 1 - ((\text{zeros}/90)^{**2} + (\text{ones}/90)^{**2})$$

from sklearn.datasets import make_classification, load_iris
 from sklearn import tree
 from sklearn.model_selection import train_test_split
 from numpy as np

load the dataset

`iris = load_iris()`

`X = iris.data[:, 2:]`

`y = iris.target`

`X_train, y_train` \rightarrow `X-test, y-test =`

~~from~~ `train(X, y, test_size=0.3, shuffle=True,`
`random_state=1)`

DecisionTreeClassifier

`model = tree.DecisionTreeClassifier()`

`model = model.fit(X_train, y_train)`

Predicting:

`predicted_value = model.predict(X-test)`

`print(f"Pred Val : {predicted_value}")`

`Pred Val : [0 1 0 2 1 0 0 2 1 0 2 1 1 0 1`
`1 0 0 1 2 0 2 1 0 0 1 2 1 2 1 2`
~~2 0 1 2 0 1 2 2 0 1 2 1 7~~

25/4/24

Week-5

Create & fit the Model :

```
model = LR()
model.fit(X_train, y_train)
```

```
X-test # display
```

Imports

```
# import pandas as pd
from matplotlib import pyplot as plt
from sklearn.model_selection import train_test_split as tts.
```

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

```
import math
```

```
!matplotlib inline
```

```
→ 0.9090
```

Read the Dataset (csv file)

```
df = pd.read_csv("insurance.csv")
```

```
df.head()
```

Sigmoid functions:

```
model.coef_
→ array([[-0.12335856]])
```

Data Exploration & Visualization:

```
# α
plt.scatter(df['age'], df['bought_insurance']
            marker = '+', color = 'red')
```

```
model.intercept_
→ array([-4.89081951])
```

Analyze Sigmoid:

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

X-train, X-test, y-train, y-test = tts(df[['age']], df['bought_insurance'], train_size=0.6, random_state=42)
```

Perform Split:

Predict function:

```
def predict_function(age):  
    g = 0.042 * age - 1.53  
    y = sigmoid(g)
```

return y

Check value with probability:

age = 35

prediction-function (age)

→ 0.485

age = 88

prediction-function (age)

→ 0.89715

Show the Plot:

```
plt.show()  
plt.figure(figsize=(10, 6))
```

sns.scatterplot(data=df, x=df[["Petal Length"]]

y = df[["Sepal Width"]]

legend = False, hue = df[["Species"]]

Add labels & title to Plot:

plt.xlabel("Petal Length")

plt.ylabel("Sepal Width")

plt.title("Petal length vs Sepal Width")

plt.show()

KNN:

```
from sklearn.model_selection import train_test_split as tts  
from sklearn.neighbors import KNeighborsClassifier
```

X_train, X_test, Y_train, Y_test = tts(X, y, 0.6, 0)
neigh = KNeighborsClassifier(n_neighbors = 2)

~~iris~~ Iris/5 KNN & SVM

import pandas as pd
from matplotlib import pyplot as plt
import seaborn as sns
%matplotlib inline

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

acc = accuracy_score(y-pred, y-test)

Print (acc)

$\Rightarrow 1.0$

SVM :

from sklearn.svm import SVC

model = SVC()

model.fit(x, y)

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

acc = accuracy_score(y-pred, y-test)

Print (acc)

$\Rightarrow 1.0$

23/5/24

23-5-M Week-7
ANN :

import numpy as np

X = np.array([[2, 9], [1, 5], [3, 6]], dtype = float)

y = np.array([1, 0, 1], dtype = float)

X = X/np.amax(X, axis=0)

y = y/100

Initialize All Variables:

epoch = 100

lr = 0.1

inputlayer_neurons = 2

hiddenlayer_neurons = 3

output_neurons = 1

Initializing Weights and Bias:

wh = np.random.uniform(size=(inp-new, hid-new))

bh = np.random.uniform(size=(1, hid-new))

wout = np.random.uniform(size=(hid-new, out-new))

bout = np.random.uniform(size=(1, out-new))

Sigmoid Function:

def sigmoid(x):

return 1/(1 + (np.exp(-x)))

* derivative

def derivatives_sigmoid(x):

return x * (1 - x)

Driver Code:

```

for i in range(epoch):
    hinp = np.dot(X, wh)
    hinp = hinp + bh
    hlayer_act = sigmoid(hinp)
    out_p1 = np.dot(hlayer_act, wout)
    out_p = out_p1 out_p1 + bout
    output = sigmoid(out_p)

    EO = y - output
    outgrad = derivatives_sigmoid(output)
    d_output = EO * outgrad
    EH = d_output.dot(wout.T)

    hiddengrad = derivatives_sigmoid(hlayer_act)
    d_hidlay = EH * hiddengrad
    //end for.

    wout += hlayer_act.T.dot(d_output) * lr
    wh += X.T.dot(d_hidlay) * lr
    print("Input: {} \n Actual Output: {} \n Predicted Output: {} \n".format(X, y, output))

```

Output:

Input:	Actual Output:
$\begin{bmatrix} [0.66] \\ 1. \end{bmatrix}$	$\begin{bmatrix} 0.33 \\ 0.55 \end{bmatrix}$
1.0	0.66

Predicted Output: $\begin{bmatrix} 0.92 \\ 0.86 \end{bmatrix}$

Actual Output: $\begin{bmatrix} 0.89 \\ 0.8460 \end{bmatrix}$

Random Forest :

$y_pred = \text{model}.\text{predict}(X_test)$

from sklearn.ensemble import RandomForestClassifier as
rfc

import pandas as pd

from sklearn.metrics import accuracy_score as acc

LOAD DATASET:

df = pd.read_csv("content/drive/MyDrive/
adult-data.csv")

data = df.dropna(axis=0)

y = data.Price

m_features = [

'Rooms', 'Bathroom', 'Landsize',
'BuildingArea', 'Year Built', 'Latitude',
'Longitude'

X = data[m_features]

Driver Code:

from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = tts(X, y, train_size=0.8,
as tts

model = rfc()

model.fit(X, y)

$a = \text{acc}(y_pred, y_test)$

print(a)

$\rightarrow 0.9935$

AdaBoost

Week -8

A 30-5-24

Load Data same as Random Forest
 $X_{\text{train}}, X_{\text{test}}, Y_{\text{train}}, Y_{\text{test}} = \text{tts}(X, y, \text{train_size} = 0.9, \text{random} = 0)$
 from sklearn.ensemble import AdaBoostClassifier as ada

ada

model = ada(n_estimators=50)

model.fit(X, y)

y-pred = model.predict(X-test)

acc = acc(y-pred, y-test)

print(acc)

$\rightarrow 1.0\%$

iris = datasets.load_iris()

X = pd.DataFrame(iris.data)

Y = X.columns = ['Sepal Length', 'Petal Length', 'Sepal Width', 'Petal Width']

y = pd.DataFrame(iris.target)

y.columns = ['Targets']

Fit the model:

model = KMeans(n_clusters = 3)

model.fit(X)

Plot the figures:

plt.figure(figsize=(14,14))
 colormap = np.array(['red', 'lime', 'black'])

K-means

AdaBoost

plt.subplot(2, 2, 1)

plt.scatter(x.Petal_Length, x.Petal_Width, c=colormap[
y.Target])

plt.title('Real Clusters')

plt.xlabel('Petal Length')

plt.ylabel('Petal Width')

plt.subplot(2, 2, 2)

plt.scatter(x.Petal_Length, x.Petal_Width, c=colormap[
y.Target])

plt.title('K-Means Clustering')

plt.xlabel('Petal Length')

plt.ylabel('Petal Width')

Imports :

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

```
import pandas as pd
```

```
import numpy as np
```

```
import seaborn as sns
```

```
from matplotlib import style
```

```
style.use('fivethirtyeight')
```

```
from sklearn.datasets import load_breast_cancer
```

```
cancer = np.loadtxt('breast_cancer.csv')
```

```
df = pd.DataFrame(cancer[[data]], columns=  
cancer[[feature_name]])
```

PCA Visualisation:

```
from sklearn.preprocessing import StandardScaler  
as ss
```

```
scaler = ss()
```

```
scaler.fit(df)
```

```
scaled_data = scaler.transform(df)
```

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

```
pca = PCA(n_components=2)
```

```
pca.fit(scaled_data)
```

X-Pca = Pca.transform(scaled-data)

scaled-data.shape

X-Pca.shape

Plot the Graph:

plt.figure(figsize=(8, 6))

plt.scatter(x-pca[:, 0], x-pca[:, 1],

c=cancer['target'], cmap='plasma')

plt.xlabel('First Principal Component')

plt.ylabel('Second Principal Component')

