

21/3/24

Week - 1

Pandas (Reading csv)

import pandas as pd

df2 = pd.read_csv("/content/sample_data/california_housing_test.csv")

df2.head()

Output:

longitude	latitude	housing_median_age	total_rooms	total_bedrooms
-122.05	37.37	27.0	3885.0	661.0

population	households	median_income	median_house_value
1537.0	606.0	6.6085	344700.0

df2.to_csv("hi.csv")

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

url = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/
iris.data"

df = pd.read_csv(url, names = col_names)

df.head()

Output:

sepal-length-in-cm	sepal-width-in-cm	petal-length-in-cm
5.1	3.5	1.4

petal-width-in-cm class

0.2 Iris-setosa

21/3/24

28/3/24

Week - 2

1. Performance Measure

2. Get the data

2.1. Download the data

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

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

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

```
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()
```

2.2 Take a quick look at the Data Structure

```
housing = load_housing_data()
```

```
housing.head()
```

```
housing.info()
```

~~```
housing.describe()
```~~

## 2.3 Plotting graphs

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

```
import seaborn as sns
```

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

```
plt.show()
```

## 2.4 Create 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)
```

```
test_indices = shuffled_indices[:test_set_size]
```

```
train_indices = shuffled_indices[test_set_size:]
```

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

```
train_set, test_set = split_train_test(data=housing)
```

```
len(train_set), len(test_set)
```

```
from zlib import crc32
```

```
def test_set_check(identifier, test_ratio=0.2):
```

```
total_size = 2**32
```

```
hex_repr = crc32(np.int64(identifier)) & 0xffffffff
```

```
in_test = hex_repr < (test_ratio * total_size)
```

```
return in_test
```

```
[test_set_check(i) for i in range(10)]
```

## 2.5 Co-ordinates to Index

```
def from_z_to_N(z):
```

```
if z >= 0:
```

$$n = 2 * z$$

```
else:
```

$$n = -2 * z - 1$$

```
return n
```

```
def cantor_pairing(n1, n2):
```

$$n = ((n1+n2)*(n1+n2+1))/2 + n2$$

```
return n
```

```
def lat_lon_to_index(lat, lon):
```

$$lat, lon = \text{int}(lat * 100), \text{int}(lon * 100)$$

$$lat, lon = \text{from\_z\_to\_N}(lat), \text{from\_z\_to\_N}(lon)$$

$$\text{index} = \text{cantor\_pairing}(lat, lon)$$

```
return np.int64(index)
```

```
housing['id'] = housing.apply(lambda row: lat_lon_to_index
```

```
(row['latitude'], row['longitude']), axis=1)
```

~~```
housing['id'].value_counts()
```~~~~```
df[28] / 3 / 24
```~~~~```
[(0) appear in i not (i) block - true - test]
```~~~~```
not in at zero block - 2
```~~

28/3/24

Week - 2

### 3. Discover & visualise the data to gain insights

#### 3.1

```
start_test_set.reset_index().to_json(frame='data/oj/start-test-set')
```

housing = start\_train\_set.copy

housing.shape

#### 3.2 Visualising Geographical Data

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

#### 3.3 Looking for correlations

corr\_matrix = housing.corr()

corr\_matrix['median\_house\_value'].sort\_values(ascending=False)

#### 3.4 Experimenting with attribute combinations

housing['rooms-per-household'] = housing['total\_rooms']/housing['households']

housing['bedrooms-per-room'] = housing['total\_bedrooms']/housing['total\_rooms']

### 4. Prepare the data for machine learning

#### 4.1 Data Cleaning

imputer = SimpleImputer(strategy='median')

housing\_num = housing.drop("ocean\_proximity", axis=1)

imputer.fill(housing\_num)

## 4.2 Handling Text and Categorical Attributes

housing-cat = housing[["ocean-proximity"]]

housing-cat.head(10)

housing-cat[["ocean-proximity"]].value\_counts()

## 4.3 Custom Transforms

class combinedAttributesAdder(BaseEstimator, TransformerMixin):

def \_\_init\_\_(self, add\_bedrooms\_per\_room=True):

self.add\_bedrooms\_per\_room = add\_bedrooms\_per\_room

def fit(self, X, y=None):

return self

def transform(self, X, y=None):

rooms\_per\_households = X[:, "rooms\_ix"] / X[:, "households\_ix"]

population\_per\_household = X[:, "population\_ix"] / X[:, "households\_ix"]

return np.c\_[X[:, "rooms\_per\_household"], population\_per\_household]

## 4.4 Transformation Pipelines

num\_pipeline = Pipeline([("imputer", SimpleImputer(strategy="median"))])

("attribs-adder", combinedAttributesAdder()),

("std-scaler", StandardScaler()))

(imputer->imputer) when I print = return

(I=mean, "ocean\_proximity\_mean") with print = mean-future

(mean-future) // print return

## 5. Train Model

```
def display_scores(scores):
 print("Scores:", scores)
 print("Mean:", scores.mean())
 print("Std:", scores.std())
```

## 6. Fine Tune Model

```
param_grid = [
 {'n_estimators': [3, 10, 30], 'max_features': [2, 4, 6, 8]},
 {'bootstrap': [False], 'n_estimators': [3, 10],
 'max_features': [2, 3, 4]}]
```

## 7. Evaluate on test set

```
final_model = grid_search.best_estimator_
```

```
final_predictions = final_model.predict(x=x_test_prepared)
```

~~final\_mse = mean\_squared\_error(y\_true=y\_test, y\_pred=~~

~~(test\_x) library: final\_predictions~~

```
final_rmse = np.sqrt(final_mse)
```

```
final_rmse
```

~~(final\_rmse = rmse, test\_y, test\_x) library: rmse~~

~~(final\_rmse = rmse, test\_y, test\_x) library: rmse~~

~~{'final\_rmse': rmse} = dictify~~

~~{'final\_rmse': rmse} = final\_rmse~~

: tuple

~~[('final\_rmse', rmse)] = dictify~~

~~[('accuracy', accuracy)] = final\_rmse~~

accuracy

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from pandas.core.common import random_state
from sklearn.linear_model import LinearRegression

df_sal = pd.read_csv("SalaryData.csv")
df_sal.head()
x = df_sal.iloc[:, :-1]
y = df_sal.iloc[:, -1]

x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2,
 random_state=0)
model = LinearRegression()
model.fit(x_train, y_train)

y_pred = model.predict(x_test)
y_pred_train = model.predict(x_train)

plt.scatter(x_train, y_train, color='lightcoral')
plt.scatter(x_train, y_pred_train, color='firebrick')
plt.scatter(x_test, y_test, color='lightcoral')

print(f"Coefficient: {model.coef_}")
print(f"Intercept: {model.intercept_}")
```

Output:

Coefficient: [9312.57512673]

Intercept: [26780.09915063]

Ans 25/4/24

18/4/24

Week - 3

Decision Tree

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn import datasets
from sklearn.model_selection import train_test_split
iris = datasets.load_iris()
```

```
X_train, X_test, y_train, y_test = train_test_split(iris.data,
iris.target, test_size=0.33, random_state=42)
```

```
from sklearn.tree import DecisionTreeClassifier
```

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

```
from sklearn import tree
tree.plot_tree(model)
```

```
from sklearn.metrics import accuracy_score, classification_report
y_pred = model.predict(X_test)
print(classification_report(y_test, y_pred, target_names=['0', '1', '2']))
```

Output:

0.98

$x[2] \leq 2.45$   
 $\text{gini} = 0.666$   
 $\text{samples} = 100$   
 $\text{values} = [31, 35, 34]$

$\text{gini} = 0.0$   
 $\text{samples} = 31$   
 $\text{value} = [31, 0, 0]$

$x[3] \leq 1.75$   
 $\text{gini} = 0.5$   
 $\text{samples} = 69$   
 $\text{values} = [0, 35, 34]$

$\text{gini} = 0.188$   
 $\text{samples} = 38$   
 $\text{value} = [0, 34, 4]$

$\text{gini} = 0.062$   
 $\text{samples} = 31$   
 $\text{value} = [0, 1, 30]$

~~Day 25/12/24~~

25/4/24

Week - 4

## Logistic Regression

```
import pandas as pd
```

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

```
%matplotlib inline
```

```
from sklearn.model_selection import train_test_split
```

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

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

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

```
df.head()
```

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

```
X_train, X_test, y_train, y_test = train_test_split(df[['age']],
df.bought_insurance, train_size=0.7, random_state=42)
```

```
model = LogisticRegression()
```

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

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

```
model.predict_proba(X_test)
```

```
model.predict(X_test)
```

```
model.score(X_test, y_test)
```

~~```
model = LinearRegression()
```~~~~```
model.fit(X_train, y_train)
```~~~~```
model.score(X_test, y_test)
```~~

```
import math
```

1 - 1/391

```
def sigmoid(z):
```

```
    return 1/(1+math.exp(-z))
```

```
def predi(age):
```

$z = 0.042 * \text{age} - 1.53$

$y = \text{sigmoid}(z)$

```
return y
```

```
print(predi(35))
```

```
print(predi(43))
```

Output:

Prediction: array([1, 0, 1, 0, 0, 0, 1, 0])

Score: 0.8888

Linear Reg Score: 0.584321

~~Predictions: 0.485~~

0.5685

0725412u

(1 - 1/391) = 1/391
(test_B, test_X) \rightarrow 1/391
(test_B, test_X) \rightarrow 1/391

9/15/24

Week - 6

```
import pandas as pd  
import numpy as np  
import matplotlib.pyplot as plt  
import seaborn as sns  
from sklearn.svm import SVC  
from sklearn.model_selection import train_test_split  
from sklearn.metrics import confusion_matrix  
from sklearn.metrics import accuracy_score, classification_report  
from sklearn.neighbors import KNeighborsClassifier
```

cd_names = ["sepal-length-in-cm", "Sepal-width-in-cm", "petal-length-in-cm",
"petal-width-in-cm", "class"]

dataset = pd.read_csv("https://archive.ics.uci.edu/ml/datasets/Iris", header =
None, names = colnames)

dataset = dataset.replace({"class": {"Iris-setosa": 1, "Iris-versicolor": 2,
"Iris-virginica": 3}})

dataset.head()

sns.heatmap(dataset.corr())
plt.title('Correlation On Iris Classes')

X = dataset.iloc[:, :-1]

y = dataset.iloc[:, -1].values

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25,
random_state = 0)

```
classifier = SVC(kernel='linear', random_state=0)
```

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

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

```
cm = confusion_matrix(y-test, y-pred)
```

```
print(cm)
```

```
# print(y-pred)
```

```
print(classification_report(y-test, y-pred, target_names=[1, 2, 3]))
```

```
classifier2 = KNeighborsClassifier(n_neighbors=2)
```

```
classifier2.fit(X-train, y-train)
```

```
y-pred = classifier2.predict(X-test)
```

```
print(y-pred)
```

```
cm = confusion_matrix(y-test, y-pred)
```

```
print(cm)
```

```
print(classification_report(y-test, y-pred, target_names=[1, 2, 3]))
```

Accuracy:

SVM : 0.98

KNN : 0.98

23/5/24 Week - 7

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

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

```
df.isnull().sum()
```

```
df['Prediction'].fillna(df.Prediction.mode()[0], inplace=True)
```

```
df.isnull().sum()
```

```
df.fillna(value='None', inplace=True)
```

```
df = df.drop('Allergen's', axis=1)
```

```
y = df['Prediction']
```

```
x = df.drop('Prediction', axis=1)
```

```
# x = x.drop('Food Product', axis=1)
```

```
x = x.apply(lambda y: "-".join(map(str, y)), axis=1)
```

```
from sklearn.model_selection import train_test_split
```

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

```
from sklearn.feature_extraction.text import TfidfVectorizer
```

```
from sklearn.pipeline import Pipeline
```

```
from sklearn.ensemble import AdaBoostClassifier
```

```
from sklearn.metrics import accuracy_score
```

```
clf = Pipeline([
```

```
    ('vectorizer_tfidf', TfidfVectorizer()),
```

```
    ('AdaBoost', AdaBoostClassifier(n_estimators=150, learning_rate=0.8))
```

```
])
```

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

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

```
print("Accuracy is ", accuracy_score(y-test, y-pred))
```

```
from sklearn.ensemble import RandomForestClassifier
```

```
clf2 = Pipeline([
```

```
    ('vectorizer_tfidf', TfidfVectorizer()),
```

```
    ('RF', RandomForestClassifier(n_estimators=100, criterion='entropy'))
```

```
])
```

```
clf2.fit(x-train, y-train)
```

```
y-pred = clf2.predict(x-test)
```

```
print("Accuracy is ", accuracy_score(y-test, y-pred))
```

Output :

AdaBoost : 0.95

Random Forest : 0.9375

23/5/24

ANN

import numpy as np

$x = \text{np.array}([2, 9], [1, 5], [3, 6]), \text{dtype} = \text{float})$

$y = \text{np.array}([82], [86], [89]), \text{dtype} = \text{float})$

$x = X / \text{np.amax}(X, \text{axis} = 0)$

$y = y / 100$

epoch = 500

lr = 0.1

input_layer_neurons = 2

hidden_layer_neurons = 3

output_neurons = 1

$wh = \text{np.random.uniform}(\text{size} = (\text{input_layer_neurons}, \text{hidden_layer_neurons}))$

$bh = \text{np.random.uniform}(\text{size} = (1, \text{hidden_layer_neurons}))$

$wout = \text{np.random.uniform}(\text{size} = (\text{hidden_layer_neurons}, \text{output_neurons}))$

$bout = \text{np.random.uniform}(\text{size} = (1, \text{output_neurons}))$

def sigmoid(x):

return $1 / (1 + \text{np.exp}(-x))$

def derivatives_sigmoid(x):

return $x * (1 - x)$

for i in range(epoch):

~~hinp~~ = np.dot(X, wh)

hinp = hinp + bh

hlayer_act = sigmoid(hinp)

out~~inp~~ = np.dot(hlayer_act, wout)

out~~inp~~ = out~~inp~~ + bout

output = sigmoid(out~~inp~~)

$EO = y - \text{output}$

DATA

$\text{outgrad} = \text{derivatives_sigmoid}(\text{output})$

$d\text{-output} = EO * \text{outgrad}$

$EH = d\text{-output} \cdot \text{dot}(w_{\text{out}}, T)$

$\text{hiddengrad} = \text{derivatives_sigmoid}(h_{\text{layer_act}})$

$d\text{-hidden_layer} = EH * \text{hiddengrad}$

$w_{\text{out}}^T = h_{\text{layer_act}} \cdot T \cdot \text{dot}(d\text{-output}) * b_r$

$w_h^T = X \cdot T \cdot \text{dot}(d\text{-hidden_layer}) * b_r$

`print("Input: " + str(X))`

`print("Actual Output: " + str(y))`

`print("Predicted Output: " + str(output))`

Input:

$\begin{bmatrix} 0.6667 & 1. \end{bmatrix}$

$\begin{bmatrix} 0.3333 & 0.5556 \end{bmatrix}$

$\begin{bmatrix} 1. & 0.6667 \end{bmatrix}$

Actual Output:

$\begin{bmatrix} 0.92 \end{bmatrix}$

$\begin{bmatrix} 0.96 \end{bmatrix}$

$\begin{bmatrix} 0.89 \end{bmatrix}$

Predicted Output:

$\begin{bmatrix} 0.85099637 \end{bmatrix}$

$\begin{bmatrix} 0.83531048 \end{bmatrix}$

$\begin{bmatrix} 0.85823955 \end{bmatrix}$

$\text{array} [0.572]$

30/5/24

K-Means

```
import matplotlib.pyplot as plt  
from sklearn import datasets  
from sklearn.cluster import KMeans  
import pandas as pd  
import numpy as np
```

```
iris = datasets.load_iris()
```

```
X = pd.DataFrame(iris.data)
```

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

```
y = pd.DataFrame(iris.target)
```

```
y.columns = ['Targets']
```

```
model = KMeans(n_clusters = 3)
```

```
model.fit(X)
```

```
plt.figure(figsize=(14,14))
```

```
colorMap = np.array(['red', 'lime', 'black'])
```

```
plt.subplot(2,2,1)
```

```
plt.scatter(X.Petal-Length, X.Petal-Width, c = colorMap[y.Targets], s=40)
```

```
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[Model.labels_], s=40)
```

```
plt.title('K-Means Clustering')
```

```
plt.xlabel('Petal Length')
```

```
plt.ylabel('Petal Width')
```

30/5/24

PCA

import (

matplotlib.pyplot as plt,
pandas as pd,
numpy as np,
seaborn as sns)

%matplotlib inline

from sklearn.datasets import load_breast_cancer

Cancer = load_breast_cancer()

Cancer.keys()

print(Cancer['DESCR'])

from sklearn.preprocessing import StandardScaler,

scaler = StandardScaler()

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

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

plt.scatter(x_pca[:, 0], x_pca[:, 1], c=cancer['target'], cmap='plasma')

plt.xlabel('First Principle Component')

plt.ylabel('Second Principle Component')

Arg 30/5/24