

Name _____

Prajwal Bhagat

Sub.

Std.:

Div.

Roll No.

Telephone No.

E-mail ID.

Blood Group.

Birth Day.

Sr.No.	Title	Page No.	Sign./Remarks
1.	Import & exporting data		
2.	End to End ml project		
3.	Linear Regression		
	multiple Linear Regression		
4.	Decision tree		
5.	Logistic Regression KNN		
6.	K means SVM PCA		Not yet 30/5 To complete To know concepts
7.	Ensemble Learning Random forest Boosting		

21/3/24

Import and Export Pandas library further

import pandas as pd

col_names = ["sepal-length-in-cm", "sepal-width-in-cm", "petal-length-in-cm", "petal-width-in-cm", "Class"]

data = pd.read_csv("C:/Users/Admin/Desktop/Iris/Iris.csv");

data.columns = col_names.

Output :

sepal length in - Cm	sepal width	petal length	petal width
4.9	3.0	1.4	0.1

4.7	3.2	1.3	0.2
-----	-----	-----	-----

~~No value~~

Class

Iris-setosa

Iris-virginica

27/3/24

End to End ml Project.

1. Problem Formulation

formulation formulate the problem by considering what input has to be given what is expected output and features available

2. get the Data

- * import OS

- provide functions for interacting with OS

- * import URL lib

- provides functions for fetching URLs

- import Pandas as pd

- provides functions to work and manipulate data sets.

housing.head() \Rightarrow return top 5 rows of file

housing.info() \Rightarrow gives the count of non null values in each attribute.

* housing.describe() gives the statistical measures such as rows mean, std, min, max and other attribute features.

split train - test()

it divides the dataset into train set and test set the ratio can be specified.

train set. shape

gives the shape of train set.

train['id'].value_count()

gives the count of values.

3. Discover and visualize the data to gain insights

- * train - test - set. shape()

- * we plot scatterplot by using housing.plot (kind= 'scatter', n-long y-bt)

plt. show()

- * using scattermatrix plot correlation.

scatter_matrix (from = housing [attribute]. f9
size = (12, 7))

4. prepare data for ml algorithms.

→ Data cleaning

~~Not~~ get rid of whole attributes
housing. drop (total_bedrooms, axis=0)

- * set missing values to some value
zero, mean, median

→ Handling test and categorical attribute

→ Custom Transformer

→ we use fit, transform

5. feature Mating

min_max scales & standardization

5. Select a train model

we use linear regression model

Decision tree regression - it is a machine learning algo used for regression task where the goal is to predict a continuous target variable

6. Tune your model

→ Grid search CV is a model provided by Scikit learn library in python for hyper parameter tuning of machine learning models.

7. Launch, monitor & maintain your system
we can automate the process by

→ collecting fresh data regularly & labeling it

→ writing script to train models & find best hyper parameters

→ writing script to evaluate modl.

Implementation of Linear Regression

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

```
def estimate_coef(x, y):
    n = np.size(x)
    mx = np.mean(x)
    my = np.mean(y)
    ss_xy = np.sum(y*x) - n*mx*my
    b_1 = ss_xy / ss_xx
    b_0 = my - b_1 * mx
    return (b_0, b_1)
```

```
def plot_regression_line(m, y, b):
    plt.. scatter(m, y, color = "m", marker = "o", s = 10)
    y_pred = b[0] + b[1]*x
    plt.plot(x, y_pred, color = "g")
    plt.xlabel("x")
    plt.ylabel("y")
```

N
9/5/2021

```
def main():
    m = np.array([0, 1, 2, ..., 9])
    y = np.array([1, 3, 4, ..., 12])
    b = estimate_coef(x, y)
    print(b)
    plot_regression_line(m, y, b)
```

Output:

$(b_0, b_1) : (1.2363, \dots, 1.69691\dots)$

~~multiple~~ multiple Linear Regression

```
from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt
import numpy as np
from sklearn import datasets, linear_model, metrics
```

```
data_url = "https://lib.stat.cmu.edu/datasets/boston"
raw_df = pd.read_csv(data_url, skiprows=22,
                     header=None)
```

```
X = np.stack([raw_df.values[:, 2], raw_df.values[:, 1]], axis=1)
y = raw_df.values[:, 5].values
```

```
reg = linear_model.LinearRegression()
reg.fit(X_train, y_train)
```

```
print("Coefficients:", reg.coef_)
```

```
print("variance score = {} ".format(reg.score(X-test, y-test)))
```

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

```
print("Co")
```

```
plt.scatter(reg.predict(X-train), reg.predict(X-train),
            y-train, color="green", s=10,
            label="Train-data")
```

```
plt.scatter(reg.predict(X-test), reg.predict(X-test),
            y-test, color="red",
            label="Test-data")
```

plt. union (y=0, xmin=0, xmax= 50, linewidth=2)

plt. legend (loc = "upper right")

plt. title ("Residual error")

plt. show()

Implementation of ID3.

```

import numpy as np
import pandas as pd
eps = np.info('float').eps
from numpy import log2 as log
from google.colab import drive
drive.mount('/content/drive')
path = 'drive/myDrive/mldataset/playtennis.csv'
df = pd.read_csv(path)

def find_entropy(df):
    target = df.keys()[-1]
    entropy = 0
    values = df[target].unique()

    for value in values:
        fraction = len(df[df[target] == value]) / len(df)
        entropy += -fraction * np.log2(fraction)

    return entropy
  
```

if average-information(df, attribute)

target = df.keys()[-1]

target-variable = df[target].unique()

variable = df[attribute].unique()

entropy = 0

for variable in variables

entropy = 0

for target-variable != target-variable

num = len(df[attribute][df[attribute] ==

`(method) [df(target) -> target(renamable)]`

der = len([df['attribute']] [df['attribute'] == variable])

fraction = num / (den + eps)

entropy + - fraction * log(fraction + ps)

Fraction 2: $\text{dexpf}[\ln(\text{df})]$

entropy & + = fraction & entropy

Defin abs(entropy)

tree = buildTree(df)

import pprint

$\text{Pf point} = \text{Ppp point (tree)}$

11984

Implementation of Logistic Regression

```
import pandas as pd
```

```
import matplotlib.lib as plt
```

```
%matplotlib inline
```

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

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

```
df.head()
```

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

```
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=0.8)
```

```
print(X_test)
```

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

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

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

```
print(y_predict)
```

```
print(X_test)
```

from sklearn.linear_model import LinearRegression

model = LinearRegression()

model.fit(x_train, y_train)

print("coefficients (m):", model.coef_)

print("intercept (b):", model.intercept_)

import math

def sigmoid(x):

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

def prediction_func(x):

$$m = 0.042$$

$$b = -1.53$$

$$Z = m \cdot \text{age} + b$$

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

Return y.

N
9/5/24

Implementation of KNN

```
import numpy as np
```

```
import pandas as pd
```

```
dataset = pd.read_csv('C:/content/notebooks/ml-Lab 1/iris.csv')
```

```
dataset.shape
```

```
dataset.head()
```

```
dataset.describe()
```

```
dataset.groupby('species').size()
```

```
feature_columns = ['Sepal length cm', 'Sepal width cm',  
'Petal length cm', 'Petal width cm']
```

```
X = dataset[feature_columns].values
```

```
y = dataset['species'].values
```

```
from sklearn.preprocessing import LabelEncoder
```

```
le = LabelEncoder()
```

```
y = le.fit_transform(y)
```

```
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, random_state=20)
```

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

```
from sklearn.neighbors import KNeighborsClassifier
```

```
from sklearn.metrics import confusion_matrix,  
accuracy_score
```

```
from sklearn.model_selection import cross_val_score
```

classifier = KNeighborsClassifier(n_neighbors=3)

classifier.fit(x_train, y_train)

y_pred = classifier.predict(x_test)

accuracy = accuracy_score(y_test, y_pred)*100

print(accuracy)

N
also

SVM

```
from sklearn.datasets import load_breast_cancer
import matplotlib.pyplot as plt
from sklearn import svm
from sklearn import SVC
```

Display

```
from sklearn.svm import SVC
```

Cancer = load_breast_cancer()

X = cancer.data[:, :-1]

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.get_cmap

alpha = 0.8,

X_label = cancer.feature_names[0],

Y_label = cancer.feature_names[1],

plt.scatter(X[:, 0], X[:, 1], C=y, s=20, edgecolor='color')

plt.show()

PCA

import pandas as pd
 import numpy as np
 from sklearn.decomposition import PCA.
 from sklearn.preprocessing import StandardScaler

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

scaling = StandardScaler()
 scaling.fit(data)

Scaled_data = scaling.transform(data)

principal = PCA(n_components=3)

principal.fit(scaled_data)

x2 = principal.transform(scaled_data)

plt.figure(figsize=(10, 10))

plt.scatter(x[:, 0], x[:, 1], c=data['target'],
 p1 = x2[:, 0])

plt.ylabel('PC1')

K-means

import pandas as pd

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

data.head()

import numpy as np

from sklearn.datasets import load_iris
 from sklearn.cluster import KMeans

Y, y = load_iris(return_X_y=True)

Kmeans = KMeans(n_clusters=3, random_state=2)

Kmeans.fit(X)

pred = Kmeans.predict(X)

pred.

Ensemble Boosting

from sklearn.datasets import load_iris
 from sklearn.model_selection import train_test_split

from sklearn.ensemble import AdaBoostClassifier

from sklearn.metrics import accuracy_score

iris = load_iris()

X = iris.data

y = iris.target

X-train, X-test, y-train, y-test = train_test_split(X, y, test_size=0.4, random_state=42)

AdaBoostClassifier(n_estimators=10, learning_rate=1.0, random_state=42)

adaBoost_clf.fit(X-train, y-train)

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

accuracy = accuracy_score(y-test, y-pred)
 print("Accuracy", accuracy)

Ensemble Random Forest

import numpy as np

import pandas as pd

- import matplotlib.pyplot as plt

- import seaborn as sns

import sklearn

from sklearn import datasets

- from sklearn.model_selection import train_test_split

from sklearn.ensemble import RandomForestClassifier

from sklearn.metrics import accuracy_score,
confusion_matrix, classification_report

iris = datasets.load_iris()

type(iris)

iris. target_names.

iris. feature_names.

iris. target.

iris. data = pd. DataFrame(2).

'sepal length' : iris. data[:, 0],

'sepal width' : iris. data[:, 1],

'petal length' : iris. data[:, 2],

'petal width' : iris. data[:, 3],

'species' : iris. target

3)

iris_data.head()

X = iris_data.iloc[:, :-1].values

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

X-train, X-test, Y-train, Y-test =

train_test_split(X, y, test_size=0.33, random_state=42)

model = RandomClassifier()

model.fit(X-train, y-train)

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

accuracy_score(y-test, y-pred)

model1 = RandomForestClassifier(n_estimators=10)

model1.fit(X-train, y-train)

feature_inp = pd.Series(data=model1.feature_importances_, index=iris['feature_names'])

feature_inp = feature_inp.sort_values(ascending=False)

feature_inp.

Sns. barplot(x=feature_inp, y=feature_inp)

plt.xlabel('feature importance score')

plt.ylabel('feature')

plt.title('Visualizing important features')

plt.show()

Not done
as