

VISVESVARAYA TECHNOLOGICAL UNIVERSITY

“JnanaSangama”, Belgaum -590014, Karnataka.



LAB REPORT

On

MACHINE LEARNING

Submitted by

Mohammad Adnan Khan (1BM21CS107)

**in partial fulfilment for the award of the degree of
BACHELOR OF ENGINEERING**

in

COMPUTER SCIENCE AND ENGINEERING



**B.M.S. COLLEGE OF ENGINEERING
(Autonomous Institution under VTU)**

BENGALURU-560019

March 2024 to June 2024

B. M. S. College of Engineering,

Bull Temple Road, Bangalore 560019

**(Affiliated To Visvesvaraya Technological University, Belgaum) Department
of Computer Science and Engineering**

CERTIFICATE



This is to certify that the Lab work entitled “**MACHINE LEARNING**” carried out by **Mohammad Adnan Khan (1BM21CS107)**, who is bonafide student of **B. M. S. College of Engineering**. It is in partial fulfilment for the award of **Bachelor of Engineering in Computer Science and Engineering** of the Visvesvaraya Technological University, Belgaum during the year 2023-24. The Lab report has been approved as it satisfies the academic requirements in respect of Machine Learning Lab - (**22CS6PCMAL**) work prescribed for the said degree.

Sonika S

Assistant Professor

Department of CSE
BMSCE, Bengaluru

Dr. Jyothi S Nayak

Professor and Head

Department of CSE
BMSCE, Bengaluru

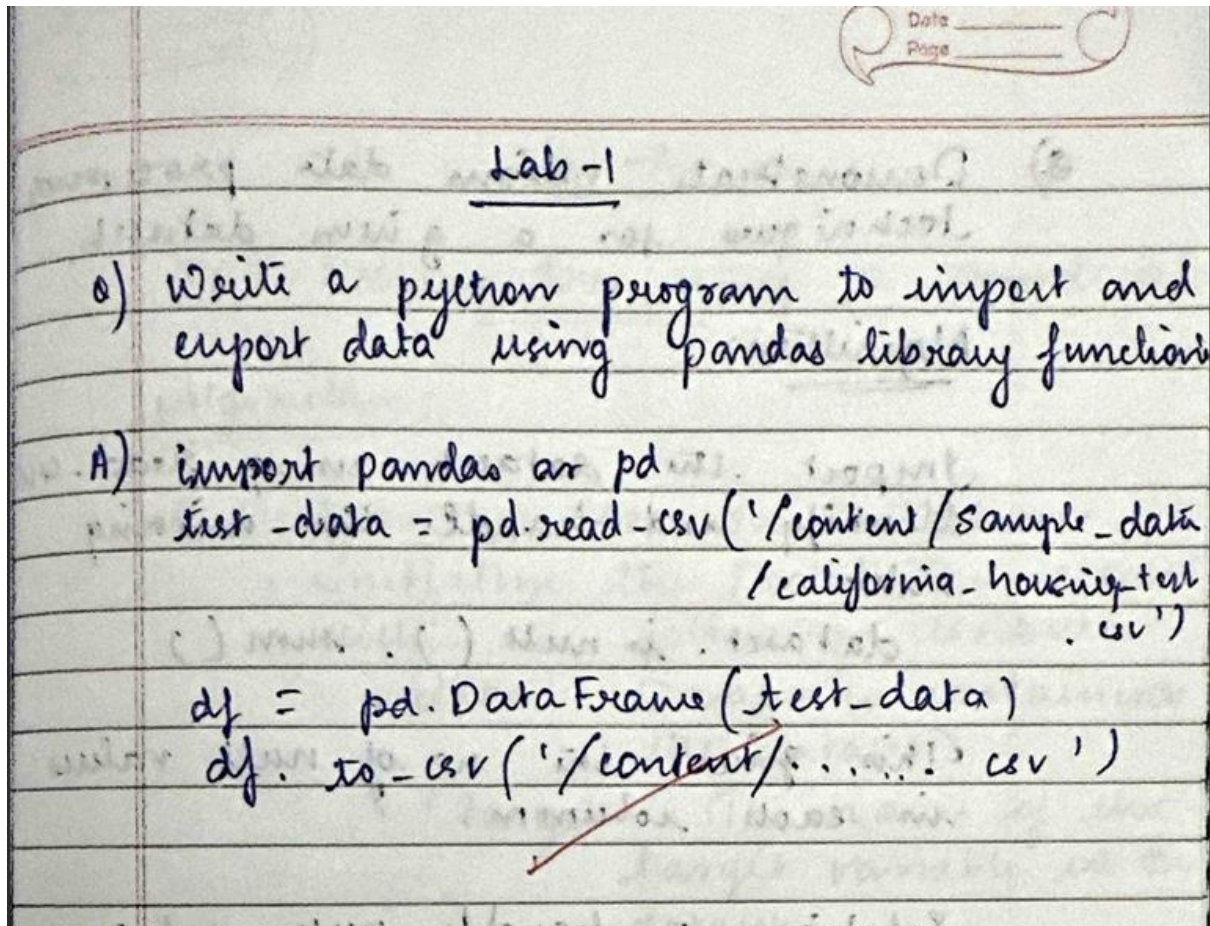
Index

Sl. No.	Experiment Title	Page No.
1	Write a python program to import and export data using Pandas library functions	1
2	Demonstrate various data pre-processing techniques for a given dataset	4
3	Use an appropriate data set for building the decision tree (ID3) and apply this knowledge to classify a new sample.	6
4	Build KNN Classification model for a given dataset.	13
5	Implement Linear and Multi-Linear Regression algorithm using appropriate dataset	20
6	Build Logistic Regression Model for a given dataset	27
7	Build Support vector machine model for a given dataset	37
8	Build k-Means algorithm to cluster a set of data stored in a .CSV file.	43
9	Implement Dimensionality reduction using Principal Component Analysis (PCA) method.	47
10	Build Artificial Neural Network model with back propagation on a given dataset	51
11	a) Implement Random Forest ensemble method on a given dataset. b) Implement Boosting ensemble method on a given dataset.	54

Course outcomes:

CO1	Apply machine learning techniques in computing systems
CO2	Evaluate the model using metrics
CO3	Design a model using machine learning to solve a problem
CO4	Conduct experiments to solve real-world problems using appropriate machine learning techniques

Write a python program to import and export data using Pandas library functions



```
import pandas as pd
df=pd.read_csv("austinHousingData.csv")
print(df.head())
df.to_csv("exported_dataset.csv")
```

Output:

```
   zipid  city      streetAddress  zipcode  ...  numOfBathrooms  numOfBedrooms  numOfStories  homeImage
0  111373431  pflugerville  14424 Lake Victor Dr  78660  ...      3.0          4          2  111373431_ffce26843283d3365c11d81b8ebdc6f-p_f...
1  120900430  pflugerville  1104 Strickling Dr  78660  ...      2.0          4          1  120900430_8255c127be8dcf0a1a18b7563d987088-p_f...
2  2084491383  pflugerville  1408 Fort Dessau Rd  78660  ...      2.0          3          1  2084491383_a2ad649e1a7a098111dcea084a11c855-p_...
3  120901374  pflugerville  1025 Strickling Dr  78660  ...      2.0          3          1  120901374_b469367a619da85b1f5ceb69b675d88e-p_f...
4   60134862  pflugerville  15005 Donna Jane Loop  78660  ...      3.0          3          2  60134862_b1a48a3df3f111e005bb913873e98ce2-p_f.jpg

[5 rows x 47 columns]
```

2. Demonstrate various data pre-processing techniques for a given dataset

Q) Demonstrate various data processing techniques for a given dataset.

Algorithm

Import the dataset using read_csv.
Identify and handle the missing values.

`dataset.isnull().sum()`

This gives the no of null values in each column.

Solution to handle null values:

- 1) Use Dropna to drop column having high no of null values
- 2) Use Fillna to replace a NULL value with a specified value.

Encoding categorized data using `pd.get_dummies()` which convert categorical data into dummy or indicator values.


```
D: %matplotlib inline
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import sklearn

(1) Python

dataset = pd.read_csv("/content/Data.csv")

(2) Python

df = pd.DataFrame(dataset)

(3) Python

df

(4) Python
Outputs are collapsed --

# = df.iloc[:, 1:1].values
y = df.iloc[:, -1].values
Add to CodeDPT Chat (Ctrl + Shift + E). Add to chat (Ctrl+L) | Edit highlighted code (Ctrl+I).

(5) Python

print(X)

(6) Python
Outputs are collapsed --

print(y)

(7) Python

... ['No' 'Yes' 'No' 'No' 'Yes' 'Yes' 'No' 'Yes' 'No' 'Yes']

(8) Python

df.isnull().sum()

(9) Python
Outputs are collapsed --
```

Dropna

```
(10) Python

df1 = df.copy()

(11) Python

# summarize the shape of the raw data
print('Before:',df1.shape)

# drop rows with missing values
df1.dropna(inplace=True)

# summarize the shape of the data with missing rows removed
print('After:',df1.shape)

(12) Python
Outputs are collapsed --

# summarize the shape of the raw data
print('Before:',df1.shape)

# drop rows with missing values
df1.dropna(inplace=True)

# summarize the shape of the data with missing rows removed
print('After:',df1.shape)

(13) Python
Outputs are collapsed --
```

Scikit-Learn

```
(14) Python

X

(15) Python
Outputs are collapsed --

from sklearn.impute import SimpleImputer
imputer = SimpleImputer(missing_values=np.nan, strategy='mean')
imputer.fit(X[:, 1:3])
X[:, 1:3] = imputer.transform(X[:, 1:3])

(16) Python

print(X)

(17) Python
Outputs are collapsed --
```

Using Dummies

```
(18) Python

df2

(19) Python
Outputs are collapsed --

pd.get_dummies(df2)
Add to CodeDPT Chat (Ctrl + Shift + E). Add to chat (Ctrl+L) | Edit highlighted code (Ctrl+I).

(20) Python
Outputs are collapsed --
```

Output:

```
df
[1] Python
Country Age Salary Purchased
0 France 44.0 72000.0 No
1 Spain 27.0 46000.0 Yes
2 Germany 30.0 54000.0 No
3 Spain 36.0 61000.0 No
4 Germany 40.0 NaN Yes
5 France 35.0 58000.0 Yes
6 Spain NaN 52000.0 No
7 France 48.0 79000.0 Yes
8 Germany 50.0 63000.0 No
9 France 37.0 67000.0 Yes

# = df.iloc[:, 1:3].values
y = df.iloc[:, 3].values
[1] Python
Add to CodeDPT Chat (Ctrl + Shift + E), Add to chat (Ctrl+4) | Edit highlighted code (Ctrl+1).

print(X)
[1] Python
[(['France', 44.0, 72000.0],
 ['Spain', 27.0, 46000.0],
 ['Germany', 30.0, 54000.0],
 ['Spain', 36.0, 61000.0],
 ['Germany', 40.0, nan],
 ['France', 35.0, 58000.0],
 ['Spain', nan, 52000.0],
 ['France', 48.0, 79000.0],
 ['Germany', 50.0, 63000.0],
 ['France', 37.0, 67000.0])]
```

```
print(y)
[1] Python
['No' 'Yes' 'No' 'No' 'Yes' 'Yes' 'No' 'Yes' 'No' 'Yes']

df.isnull().sum()
[1] Python
Country      0
Age           1
Salary        1
Purchased     0
dtype: int64
```

```
# summarize the shape of the raw data
print('Before:', df1.shape)

# drop rows with missing values
df1.dropna(inplace=True)

# summarize the shape of the data with missing rows removed
print('After:', df1.shape)
[1] Python
Before: (10, 4)
After: (8, 4)
```

```
Scikit-Learn

X
[1] Python
array([[ 'France', 44.0, 72000.0],
 [ 'Spain', 27.0, 46000.0],
 [ 'Germany', 30.0, 54000.0],
 [ 'Spain', 36.0, 61000.0],
 [ 'Germany', 40.0, nan],
 [ 'France', 35.0, 58000.0],
 [ 'Spain', nan, 52000.0],
 [ 'France', 48.0, 79000.0],
 [ 'Germany', 50.0, 63000.0],
 [ 'France', 37.0, 67000.0]], dtype=object)

from sklearn.impute import SimpleImputer
imputer = SimpleImputer(missing_values=np.nan, strategy='mean')
imputer.fit(X[:, 1:3])
X[:, 1:3] = imputer.transform(X[:, 1:3])
[1] Python

print(X)
[1] Python
[(['France', 44.0, 72000.0],
 ['Spain', 27.0, 46000.0],
 ['Germany', 30.0, 54000.0],
 ['Spain', 36.0, 61000.0],
 ['Germany', 40.0, 63777.77777777778],
 ['France', 35.0, 58000.0],
 ['Spain', 36.77777777777778, 52000.0],
 ['France', 48.0, 79000.0],
 ['Germany', 50.0, 63000.0],
 ['France', 37.0, 67000.0])]
```


Using Dummies

9F2

(M)

Python

	Country	Age	Salary	Purchased
0	France	44.0	72000.0	No
1	Spain	27.0	48000.0	Yes
2	Germany	30.0	54000.0	No
3	Spain	36.0	61000.0	No
4	Germany	40.0	NaN	Yes
5	France	35.0	58000.0	Yes
6	Spain	NaN	52000.0	No
7	France	48.0	79000.0	Yes
8	Germany	50.0	69000.0	No
9	France	37.0	67000.0	Yes

D:

(M)

pd.get_dummies(9F2)

Add to CodeSPT Chat (Ctrl + Shift + E).

Add to chat (Ctrl+L) | Edit highlighted code (Ctrl+I).

Python

(M)

(M)

	Age	Salary	Country_France	Country_Germany	Country_Spain	Purchased_No	Purchased_Yes
0	44.0	72000.0	True	False	False	True	False
1	27.0	48000.0	False	False	True	False	True
2	30.0	54000.0	False	True	False	True	False
3	36.0	61000.0	False	False	True	True	False
4	40.0	NaN	False	True	False	False	True
5	35.0	58000.0	True	False	False	False	True
6	NaN	52000.0	False	False	True	True	False
7	48.0	79000.0	True	False	False	False	True
8	50.0	69000.0	False	True	False	True	False
9	37.0	67000.0	True	False	False	False	True

12/04/2024

Decision Tree Algorithm:

Lab - 2ID3 Decision tree using a sample set.Algorithm:1. Decision Tree Class Initialization:

- Initialize the Decision Tree class with the following attributes

- 'data': Dataframe containing the dataset

- 'target': The name of the target variable in the dataset

- 'positive': Value representing the positive class in the target variable

- 'parent_val': The value of the parent node in the decision tree

- 'parent': The name of the parent node

- 'childs': List to store child nodes

- 'decision': Empty string to store decision made at node

2. -get-entropy Method :

- Calculate the entropy of a given dataset based on the target variable
- Formula - $p, \text{atio} \times \log_2(p, \text{atio}) + n, \text{atio} \times \log_2(n, \text{atio})$
 where p, atio is the ratio of the positive instances and n, atio is the ratio of negative instances

3. -get-gain Method :

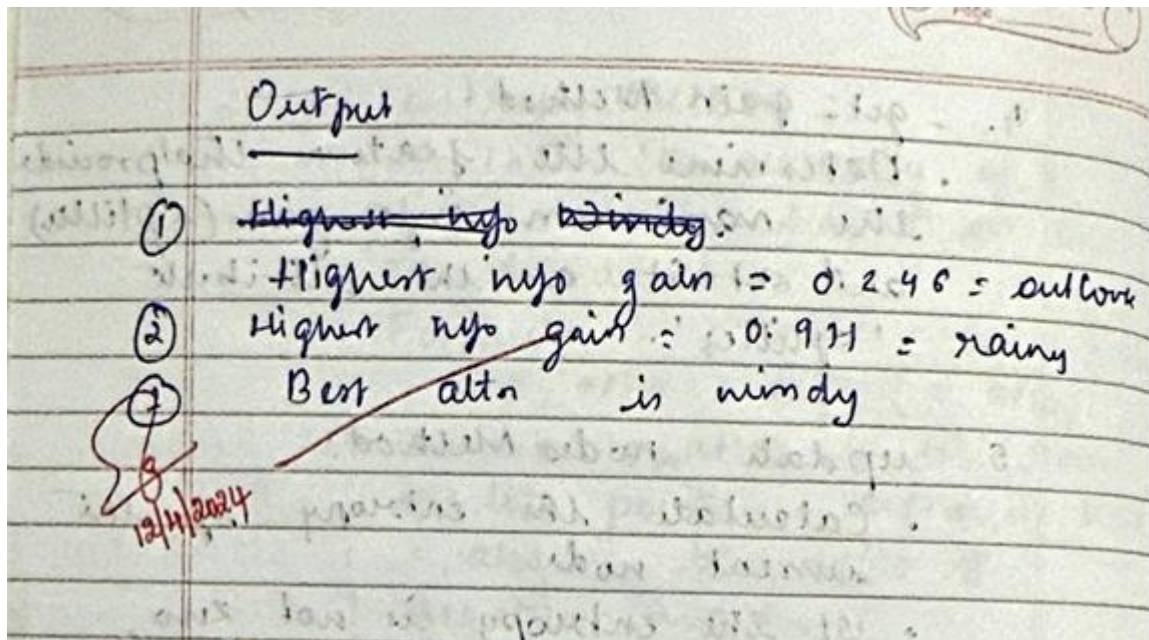
- Calculate the information gain for a given feature
- Iterate through unique values of the feature and calculate the weighted average of entropy for each value
- Formula: $\text{Entropy}(D) - \sum_{v \in \text{Values}(f)} \frac{|D_v|}{|D|} \times \text{Entropy}(D_v)$

4. - get - gain Method :

- Determine the feature that provides the maximum info gain (splitter) and set it as the attribute 'splitter'.

5. update nodes Method:

- Calculate the entropy of the current node
- If the entropy is not zero, calculate the info gain for each feature.
- Determine the feature that provides the maximum info gain (splitter)
- Iterate through unique values of the splitter feature and child nodes
- Recursively call the 'update nodes' method for each child node.



```
# Importing the required libraries
import pandas as pd
import numpy as np
import math

# Reading the dataset (Tennis-dataset)
data = pd.read_csv("../content/PlayTennis.csv")

from google.colab import drive
drive.mount("../content/drive")

def highlight(cell_value):
    """
    Highlight yes / no values in the dataframe
    """
    color_1 = 'background-color: pink;'
    color_2 = 'background-color: lightgreen;'

    if cell_value == 'no':
        return color_1
    elif cell_value == 'yes':
        return color_2

data.style.applymap(highlight)
.set_properties(subset=data.columns, **{'width': '100px'})
.set_table_styles([{'selector': 'th', 'props': [('background-color', 'lightgray'), ('border', '1px solid gray'), ('font-weight', 'bold')]},
                   {'selector': 'tr:hover', 'props': [('background-color', 'white'), ('border', '1.5px solid black')]}])
```

```

def find_entropy(data):
    """
    Returns the entropy of the class or Features
    Formula: - Σ P(X)logP(X)
    """
    entropy = 0
    for i in range(data.nunique()):
        x = data.value_counts()[i]/data.shape[0]
        entropy += (- x * math.log(x,2))
    return round(entropy,3)

def information_gain(data, data_):
    """
    Returns the information gain of the Features
    """
    info = 0
    for i in range(data_.nunique()):
        df = data[data_ == data_.unique()[i]]
        w_avg = df.shape[0]/data.shape[0]
        entropy = find_entropy(df.play)
        x = w_avg * entropy
        info += x
    ig = find_entropy(data.play) - info
    return round(ig, 3)

def entropy_and_infogain(datax, feature):
    """
    Grouping features with the same class and computing their
    entropy and information gain for splitting
    """
    for i in range(data[feature].nunique()):
        df = datax[datax[feature]==data[feature].unique()[i]]
        if df.shape[0] < 1:
            continue

        display(df[(feature, 'play')].style.applymap(highlight)\
                .set_properties(subset=[feature, 'play'], **{'width': '80px'})\
                .set_table_styles([(('selector': 'th', 'props': [('background-color', 'lightgray'),\
                ('border', '1px solid gray'),\
                ('font-weight', 'bold')]),\
                ('selector': 'td', 'props': [('border', '1px solid gray')]),\
                ('selector': 'tr:th', 'props': [('background-color', 'white'),\
                ('border', '1.5px solid black')]))]))

        print(f'Entropy of {feature} - {data[feature].unique()[i]} = {find_entropy(df.play)}')
        print(f'Information Gain for {feature} = {information_gain(datax, datax[feature])}')

```

```

print(f'Entropy of the entire dataset: {find_entropy(data.play)}')

```

Entropy of the entire dataset: 0.94

```

entropy_and_infogain(data, 'temp')

```

Outputs are collapsed --

```

entropy_and_infogain(data, 'humidity')

```

Outputs are collapsed --

```

entropy_and_infogain(data, 'windy')

```

Outputs are collapsed --

#Rainy-outlook

```

rainy = data[data['outlook'] == 'rainy']
rainy.style.applymap(highlight)\
    .set_properties(subset=data.columns, **{'width': '100px'})\
    .set_table_styles([(('selector': 'th', 'props': [('background-color', 'lightgray'), ('border', '1px solid gray'),\
    ('font-weight', 'bold')]),\
    ('selector': 'tr:th', 'props': [('background-color', 'white'), ('border', '1.5px solid black')]))]))

```

Outputs are collapsed --

```

print(f'Entropy of the Rainy dataset: {find_entropy(rainy.play)}')

```

Entropy of the Rainy dataset: 0.571

```

entropy_and_infogain(rainy, 'temp')

```

Outputs are collapsed --

```

entropy_and_infogain(rainy, 'humidity')

```

Outputs are collapsed --

```

entropy_and_infogain(rainy, 'temp')

```

Outputs are collapsed --

```

entropy_and_infogain(rainy, 'humidity')

```

Outputs are collapsed --

```

entropy_and_infogain(rainy, 'windy')

```

Outputs are collapsed --

wind has highest information gain

Output:

```
> def highlight(cell_value):
    """
    Highlight yes / no values in the dataframe
    """
    color_1 = 'background-color: pink;'
    color_2 = 'background-color: lightgreen;'

    if cell_value == 'no':
        return color_1
    elif cell_value == 'yes':
        return color_2

data.style.applymap(highlight)
.set_properties(subset=data.columns, **{'width': '100px'})
.set_table_styles([{'selector': 'th', 'props': [('background-color', 'lightgray'), ('border', '1px solid gray'), ('font-weight', 'bold')]},
                  {'selector': 'tr:hover', 'props': [('background-color', 'white'), ('border', '1.5px solid black')]])
```

[1]

sunny	hot	high	False	no
sunny	hot	high	True	yes
overcast	hot	high	False	no
rainy	mild	high	False	yes
rainy	cool	normal	False	yes
rainy	cool	normal	True	yes
overcast	cool	normal	True	yes
sunny	mild	high	False	no
sunny	cool	normal	False	yes
rainy	mild	normal	False	yes
sunny	mild	normal	True	yes
overcast	mild	high	True	yes
overcast	hot	normal	False	yes
rainy	mild	high	True	yes

```
print(f'Entropy of the entire dataset: {find_entropy(data.play)}')
```

[1] Entropy of the entire dataset: 0.94

```
entropy_and_infoGain(data, 'temp')
```

[4]

hot	no
hot	yes
hot	yes
hot	no

Entropy of temp - hot = 1.0

mild	yes
mild	no
mild	yes
mild	yes
mild	yes
mild	yes
mild	yes

Entropy of temp - mild = 0.918

cool	yes
cool	no
cool	yes
cool	yes

Entropy of temp - cool = 0.811
Information Gain for temp = 0.029

```
entropy_and_infoGain(data, 'humidity')
```

[1]

high	no
high	no
high	yes
high	yes
high	yes
high	yes
high	yes

Entropy of humidity - high = 0.985

normal	yes
normal	no
normal	yes
normal	yes
normal	yes
normal	yes
normal	yes

Entropy of humidity - normal = 0.592
Information Gain for humidity = 0.151

```
> entropy_and_infogain(data, 'windy')
[10] Python
...


|       | False | True |
|-------|-------|------|
| False | 100   | 0    |
| False | 100   | 0    |
| False | 100   | 0    |
| False | 100   | 0    |
| False | 100   | 0    |
| False | 100   | 0    |
| False | 100   | 0    |
| False | 100   | 0    |
| False | 100   | 0    |


...
Entropy of windy - False = 0.811



|      | False | True |
|------|-------|------|
| True | 0     | 100  |
| True | 0     | 100  |
| True | 0     | 100  |
| True | 0     | 100  |
| True | 0     | 100  |
| True | 0     | 100  |
| True | 0     | 100  |
| True | 0     | 100  |
| True | 0     | 100  |


...
Entropy of windy - True = 1.0
Information Gain For windy = 0.048

#Rainy -outlook
```

```
rainy = data[data['outlook'] == 'rainy']
rainy.style.applymap(hidehighlight)
rainy.set_properties(subset=data.columns, **{'width': '100px'})
rainy.set_table_styles([{'selector': 'th', 'props': [('background-color', 'lightgray'), ('border', '1px solid gray'), ('font-weight', 'bold')]}, {'selector': 'tr:hover', 'props': [('background-color', 'white'), ('border', '1.5px solid black')]}])
[1] Python
...


| rainy | mild | high   | False | True |
|-------|------|--------|-------|------|
| rainy | cool | normal | False | 100  |
| rainy | cool | normal | True  | 100  |
| rainy | mild | normal | False | 100  |
| rainy | mild | high   | True  | 100  |


...
print(f'Entropy of the Rainy dataset: {find_entropy(rainy.play)}')
[10] Python
...
Entropy of the Rainy dataset: 0.971
```

```
entropy_and_infogain(rainy, 'temp')
[11] Python
...


|      | mild | high |
|------|------|------|
| mild | 100  | 0    |
| mild | 100  | 0    |
| mild | 100  | 0    |


...
Entropy of temp - mild = 0.918



|      | cool | normal |
|------|------|--------|
| cool | 100  | 0      |
| cool | 100  | 0      |


...
Entropy of temp - cool = 1.0
Information Gain For temp = 0.02
```

```
entropy_and_infogain(rainy, 'humidity')
[12] Python
...


|      | high | normal |
|------|------|--------|
| high | 100  | 0      |
| high | 100  | 0      |


...
Entropy of humidity - high = 1.0



|        | normal | high |
|--------|--------|------|
| normal | 100    | 0    |
| normal | 100    | 0    |
| normal | 100    | 0    |


...
Entropy of humidity - normal = 0.918
Information Gain for humidity = 0.02
```

```
> entropy_and_infogain(rainy, 'windy')
[13] Python
...


|       | False | True |
|-------|-------|------|
| False | 100   | 0    |
| False | 100   | 0    |
| False | 100   | 0    |


...
Entropy of windy - False = 0.0



|      | False | True |
|------|-------|------|
| True | 0     | 100  |
| True | 0     | 100  |


...
Entropy of windy - True = 0.0
Information Gain for windy = 0.971
```

wind has highest information gain

03/05/2024

Build KNN Classification model for a given dataset.:

Implement KNN

1) Import libraries necessary for KNN

```
import pandas as pd
```

```
import numpy as np
```

```
from sklearn.neighbors import
```

```
KNeighborsClassifier
```

2) Read dataset, file path

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

```
text = df[0:10, :]
```

```
train = df[1:10, :]
```

3) Create a model and predict

```
clf = KNeighborsClassifier
```

```
(n_neighbors=5, metric='uniform')
```

```
clf.fit(X=train, y=train)
```

```
clf.predict(X=test)
```

880
8/5/24

```
import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
import matplotlib.pyplot as plt # for data visualization purposes
import seaborn as sns # for data visualization
%matplotlib inline
```

```
data = '../content/cancer_detector.txt'
df = pd.read_csv(data, header=None)
```

```
df.shape
```

Outputs are collapsed --

```
col_names = ['Id', 'Clump_thickness', 'Uniformity_Cell_Size', 'Uniformity_Cell_Shape', 'Marginal_Adhesion',
             'Single_Epithelial_Cell_Size', 'Bare_Nuclei', 'Bland_Chromatin', 'Normal_Nucleoli', 'Mitoses', 'Class']

df.columns = col_names
df.columns
```

Outputs are collapsed --

```
df.head()
```

Outputs are collapsed --

```
df.drop('Id', axis=1, inplace=True)
```

```
# view summary of dataset
df.info()
```

Outputs are collapsed --

```
for var in df.columns:
    print(df[var].value_counts())
```

Outputs are collapsed --

```
df['Bare_Nuclei'] = pd.to_numeric(df['Bare_Nuclei'], errors='coerce')
```

```
df.dtypes
```

Outputs are collapsed --

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

Outputs are collapsed --

```
# check 'na' values in the dataframe
df.isna().sum()
```

Outputs are collapsed --

```
# check frequency distribution of 'Bare_Nuclei' column
df['Bare_Nuclei'].value_counts()
```

Outputs are collapsed --

```
# check unique values in 'Bare_Nuclei' column
df['Bare_Nuclei'].unique()
```

```
array([ 1., 10.,  2.,  4.,  3.,  9.,  7., nan,  5.,  6.])
```

```
# check for nan values in 'Bare_Nuclei' column --
```

```
...
```

```
16
```

```
# view frequency distribution of values in 'Class' variable
df['Class'].value_counts()

Outputs are collapsed --

import numpy as np

# view summary statistics in numerical variables
print(round(df.describe(),2))

Outputs are collapsed --

X = df.drop(['Class'], axis=1)
y = df['Class']

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 = 0)

X_train.shape, X_test.shape

Outputs are collapsed --

for col in X_train.columns:
    if X_train[col].isnull().mean()>0:
        print(col, round(X_train[col].isnull().mean(),4))

... Bare_Nuclei 0.0233

...

cols = X_train.columns

from sklearn.preprocessing import StandardScaler--

X_train = pd.DataFrame(X_train, columns=cols)

X_test = pd.DataFrame(X_test, columns=cols)

X_train.head() --
...


|   | Clump_thickness | Uniformity_Cell_Size | Uniformity_Cell_Shape | Marginal_Adhesion | Single_Epithelial_Cell_Size | Bare_Nuclei | Bland_Chromatin | Normal_Nucleoli | Mitoses   |
|---|-----------------|----------------------|-----------------------|-------------------|-----------------------------|-------------|-----------------|-----------------|-----------|
| 0 | 2.028983        | 0.29506              | 0.28573               | 1.119077          | -0.546543                   | 1.858357    | -0.577774       | 0.041241        | -0.324258 |
| 1 | 1.669451        | 2.257680             | 2.304569              | -0.622471         | 3.108879                    | 1.297389    | -0.159953       | 0.041241        | -0.324258 |
| 2 | -1.202005       | -0.679581            | -0.717925             | 0.074148          | -1.003220                   | -0.104329   | -0.995595       | -0.608165       | -0.324258 |
| 3 | -0.125209       | -0.026896            | -0.046260             | -0.622471         | -0.546543                   | -0.665096   | -0.159953       | 0.041241        | -0.324258 |
| 4 | 0.233723        | -0.353219            | -0.382092             | -0.274161         | -0.546543                   | -0.665096   | -0.577774       | -0.283462       | -0.324258 |



# import KNeighbors Classifier from sklearn
from sklearn.neighbors import KNeighborsClassifier

# instantiate the model
knn = KNeighborsClassifier(n_neighbors=3)

# fit the model to the training set
knn.fit(X_train, y_train)

Outputs are collapsed --

y_pred = knn.predict(X_test)
y_pred

Outputs are collapsed --

knn.predict_proba(X_test)[:,:0]

Outputs are collapsed --

from sklearn.metrics import accuracy_score
print('Model accuracy score: (0:0.4f)'.format(accuracy_score(y_test, y_pred)))

Outputs are collapsed --

y_pred_train = knn.predict(X_train) --

print('Training-set accuracy score: (0:0.4f)'.format(accuracy_score(y_train, y_pred_train)))

Outputs are collapsed --
```

Output:

```
D> df.shape

[0]
... (699, 11)

col_names = ['Id', 'Clump_thickness', 'Uniformity_Cell_Size', 'Uniformity_Cell_Shape', 'Marginal_Adhesion',
             'Single_Epithelial_Cell_Size', 'Bare_Nuclei', 'Bland_Chromatin', 'Normal_Nucleoli', 'Mitoses', 'Class']

df.columns = col_names

df.columns

Index(['Id', 'Clump_thickness', 'Uniformity_Cell_Size',
      'Uniformity_Cell_Shape', 'Marginal_Adhesion',
      'Single_Epithelial_Cell_Size', 'Bare_Nuclei', 'Bland_Chromatin',
      'Normal_Nucleoli', 'Mitoses', 'Class'],
      dtype=object')

df.head()

   Id  Clump_thickness  Uniformity_Cell_Size  Uniformity_Cell_Shape  Marginal_Adhesion  Single_Epithelial_Cell_Size  Bare_Nuclei  Bland_Chromatin  Normal_Nucleoli  Mitoses  Class
0  1000025           5                1                1                1                2                1                3                1                1                2
1  1000945           5                4                4                5                7                10                3                2                1                2
2  1015425           3                1                1                1                2                2                3                1                1                2
3  1016277           6                8                8                1                3                4                3                7                1                2
4  1017023           4                1                1                3                2                1                3                1                1                2

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 699 entries, 0 to 698
Data columns (total 10 columns):
#   Column                Non-Null Count  Dtype
---  --
0   Clump_thickness        699 non-null    int64
1   Uniformity_Cell_Size   699 non-null    int64
2   Uniformity_Cell_Shape  699 non-null    int64
3   Marginal_Adhesion      699 non-null    int64
4   Single_Epithelial_Cell_Size  699 non-null    int64
5   Bare_Nuclei            699 non-null    object
6   Bland_Chromatin        699 non-null    int64
7   Normal_Nucleoli        699 non-null    int64
8   Mitoses                699 non-null    int64
9   Class                  699 non-null    int64
dtypes: int64(9), object(1)
memory usage: 54.7+ KB

for var in df.columns:
    print(df[var].value_counts())

Clump_thickness
1    145
5    130
3    108
4     80
10    69
2     50
8     46
6     34
7     23
9     14
Name: count, dtype: int64
Uniformity_Cell_Size
1    384
10   67
3    52
2    45
4    40
5    30
8    29
6    27
7    19
9     6
Name: count, dtype: int64
Uniformity_Cell_Shape
...
Class
2    458
4    241
Name: count, dtype: int64
```

```
[8] df['Bare_Nuclei'] = pd.to_numeric(df['Bare_Nuclei'], errors='coerce') Python

[9] df.dtypes Python

Clump_thickness      int64
Uniformity_Cell_Size  int64
Uniformity_Cell_Shape int64
Marginal_Adhesion     int64
Single_Epithelial_Cell_Size int64
Bare_Nuclei           float64
Bland_Chromatin       int64
Normal_Nucleoli       int64
Mitoses              int64
Class                 int64
dtype: object

[10] df.isnull().sum() Python

Clump_thickness      0
Uniformity_Cell_Size 0
Uniformity_Cell_Shape 0
Marginal_Adhesion    0
Single_Epithelial_Cell_Size 0
Bare_Nuclei          16
Bland_Chromatin      0
Normal_Nucleoli      0
Mitoses              0
Class                0
dtype: int64

[11] # check 'na' values in the dataframe
df.isna().sum() Python

Clump_thickness      0
Uniformity_Cell_Size 0
Uniformity_Cell_Shape 0
Marginal_Adhesion    0
Single_Epithelial_Cell_Size 0
Bare_Nuclei          16
Bland_Chromatin      0
Normal_Nucleoli      0
Mitoses              0
Class                0
dtype: int64

[12] # check frequency distribution of 'Bare_Nuclei' column
df['Bare_Nuclei'].value_counts() Python

Bare_Nuclei
1.0    402
10.0   132
2.0     30
5.0     30
3.0     28
8.0     21
4.0     19
9.0      9
7.0      8
6.0      4
Name: count, dtype: int64

[13] # check unique values in 'Bare_Nuclei' column
df['Bare_Nuclei'].unique() Python

array([ 1., 10.,  2.,  4.,  3.,  9.,  7., nan,  5.,  6.,  6.])

[14] # check for nan values in 'Bare_Nuclei' column--
16

[15] # view frequency distribution of values in 'Class' variable
df['Class'].value_counts() Python

Class
2     458
4     241
Name: count, dtype: int64

[16] import numpy as np Python

[17] # view summary statistics in numerical variables
print(round(df.describe(),2)) Python

Clump_thickness  Uniformity_Cell_Size  Uniformity_Cell_Shape \
count      699.00      699.00      699.00
mean         4.42         3.13         3.21
std          2.82         3.05         2.97
min           1.00         1.00         1.00
25%           2.00         1.00         1.00
50%           4.00         1.00         1.00
75%           6.00         5.00         5.00
max          10.00        10.00        10.00

Marginal_Adhesion  Single_Epithelial_Cell_Size  Bare_Nuclei \
count      699.00      699.00      699.00
mean         2.81         3.22         3.54
std          2.86         2.21         3.64
min           1.00         1.00         1.00
25%           1.00         2.00         1.00
50%           1.00         2.00         1.00
75%           4.00         4.00         6.00
max          10.00        10.00        10.00

Bland_Chromatin  Normal_Nucleoli  Mitoses  Class
count      699.00      699.00      699.00      699.00
mean         3.44         2.87         1.59         2.69
std          2.44         3.05         1.72         0.95
min           1.00         1.00         1.00         2.00
25%           2.00         1.00         1.00         2.00
50%           3.00         1.00         1.00         2.00
75%           5.00         4.00         1.00         4.00
max          10.00        10.00        10.00         4.00
```


Python

Python

Python

Python

Python

Python

Python

Python

Python

Python

Python

Python

Python

Abstract

Abstract

Abstract

Python

Python

03/05/2024

Linear Regression

Lab - 3

Implement linear regression algo using appropriate dataset.

Algorithm

1. Import necessary libraries
2. Import dataset
3. Visualization of dataset using different plots like heatmap, distribution plot, scatterplot etc.
4. Preprocess the data, convert or encode categorical data
5. Split the dataset into training set and test set from sklearn model - selection import train-test splitter

$X_{train}, X_{test}, y_{train}, y_{test}$
 $= \text{train_test_split}(X, y,$
 $\text{test_size=0.2},$
 $\text{random_state=42})$

6. Build Model

from sklearn.linear_model import

LinearRegression

lin_reg = LinearRegression()

7. Fit the dataset to the model and train it

lin. reg - fit (x_{train} , y_{train})

8. Calculate the accuracy using mean square error.

Output

~~meansquarederr = 0.9121132468~~

Implement multilinear regression

① Import necessary libraries like linear regression

② Import dataset

③ Visualise the dataset using matplotlib.pyplot

④ Encode categorical data
ct = ColumnTransformer (transformers=[('encode', OneHotEncoder, [3])], remainder='passthrough')

⑤ split dataset to training and test dataset

⑥ We can fit multiple independent variables

⑦ create regressor model
regressor = LinearRegression()

⑧ Fit the train set

⑨ Test the model using test set

⑩ Compare the actual value and predicted value

Code and output:

Linear Regression

[1] import pandas as pd

Python

[2] salary = pd.read_csv('https://github.com/ybifoundation/Dataset/raw/main/SalaryX200ata.csv')

Python

[3] y = salary['Salary']

Python

[4] X = salary[['Experience Years']]

Python

[5] from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.7, random_state=2529)

Python

[6] X_train.shape, X_test.shape, y_train.shape, y_test.shape

Python

... ((28, 1), (12, 1), (28,), (12,))

[7] from sklearn.linear_model import LinearRegression
model = LinearRegression()

Python

[8] model.fit(X_train, y_train)

Python

... * LinearRegression
LinearRegression()

[9] model.coef_

Python

... array([9405.62])

[10] y_pred = model.predict(X_test)

Python

[11] y_pred

Python

... array([90555.15, 58916.62, 106544.7 , 64219.43, 68922.24, 123474.81,
84911.78, 63278.87, 65159.99, 61397.74, 37883.7 , 50111.])

[12] from sklearn.metrics import mean_absolute_error, mean_absolute_percentage_error, mean_squared_error

Python

[13] import numpy as np

Python

[14] mean_absolute_error(y_test, y_pred)

Python

... 4025.9263101681768

Multiple Regression

[15] house = pd.read_csv('https://github.com/ybifoundation/Dataset/raw/main/Boston.csv')

Python

[16] house.head()

Python

...

	CRIM	ZN	INDUS	CHAS	NIX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	MEDV
0	0.00632	18.0	2.31	0	0.538	6.575	65.2	4.0900	1	296.0	15.3	396.90	4.98	24.0
1	0.02731	0.0	7.07	0	0.469	6.421	78.9	4.9671	2	242.0	17.8	396.90	9.14	21.6
2	0.02729	0.0	7.07	0	0.469	7.185	61.1	4.9671	2	242.0	17.8	392.83	4.03	34.7
3	0.02127	0.0	2.18	0	0.458	6.998	45.8	6.0622	3	222.0	18.7	394.63	2.94	33.4
4	0.06905	0.0	2.18	0	0.458	7.147	54.2	6.0622	3	222.0	18.7	396.90	5.33	36.2

```
D> house.describe()

[1] Python
...
count    CRIM      ZN      INDUS      CHAS      NX      RM      AGE      DIS      RAD      TAX      PTRATIO      B      LSTAT      MEDV
mean    3.61324    11.36436    11.13679    0.069170    0.554695    6.184634    68.574901    3.795043    9.549407    408.237154    18.455534    356.674012    12.653063    22.532806
std     8.601545    23.322453    6.860353    0.253994    0.115878    0.702617    28.148861    2.105710    8.707259    168.537116    2.164346    91.294864    7.141062    9.197104
min     0.006320    0.000000    0.460000    0.000000    0.395000    3.561000    2.900000    1.129600    1.000000    187.000000    12.600000    0.310000    1.710000    5.000000
25%     0.009045    0.000000    5.190000    0.000000    0.449000    5.895500    45.025000    2.100175    4.000000    279.000000    17.400000    375.375500    6.910000    17.025800
50%     0.256910    0.000000    9.690000    0.000000    0.538000    6.109500    77.500000    3.207450    5.000000    330.000000    19.050000    391.440000    11.360000    21.200000
75%     3.677000    12.500000    16.100000    0.000000    0.624000    6.623500    94.075000    5.188425    24.000000    666.000000    20.200000    386.225000    16.955000    25.000000
max    88.970200    100.000000    27.740000    1.000000    0.871000    8.780000    100.000000    12.128500    24.000000    711.000000    22.000000    396.900000    37.970000    50.000000

house.columns

[2] Python
...
Index(['CRIM', 'ZN', 'INDUS', 'CHAS', 'NX', 'RM', 'AGE', 'DIS', 'RAD', 'TAX',
      'PTRATIO', 'B', 'LSTAT', 'MEDV'],
      dtype='object')

y = house['MEDV']

[3] Python

X = house.drop(['MEDV'],axis=1)

[4] Python

from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X,y, train_size=0.7, random_state=2529)

[5] Python

X_train.shape, X_test.shape, y_train.shape, y_test.shape

[6] Python
((354, 13), (152, 13), (354,), (152,))

from sklearn.linear_model import LinearRegression
model = LinearRegression()

[7] Python

# Step 4 : train or fit model
model.fit(X_train,y_train)

[8] Python

y = LinearRegression()
LinearRegression()

model.intercept_

[9] Python
34.21916368862993

model.coef_

[10] Python
array([-1.29e+01,  3.65e+02,  1.54e+02,  2.35e+00, -2.04e+01,  4.41e+00,
        4.61e-03, -1.59e+00,  2.51e-01, -9.60e-03, -9.64e-01,  1.01e-02,
        -5.43e-01])

# Step 5 : predict model
y_pred = model.predict(X_test)

[11] Python

y_pred

[12] Python
array([31.72, 22.02, 21.17, 39.78, 20.1 , 22.86, 18.36, 14.79, 22.56,
       21.35, 18.38, 27.97, 29.86,  6.45, 10.68, 26.25, 21.89, 25.23,
       3.62, 36.22, 24.08, 22.94, 14.27, 20.79, 24.23, 16.74, 18.75,
       20.97, 28.51, 20.86,  9.23, 17.07, 22.07, 22.23, 39.26, 26.17,
       42.5 , 19.35, 34.52, 14.07, 13.81, 23.28, 11.79,  9.01, 21.05,
       25.55, 18.17, 16.82, 14.66, 14.86, 33.79, 33.27, 15.49, 24.08,
       27.64, 19.58, 45.02, 20.97, 20.07, 27.67, 34.59, 12.71, 23.66,
       31.65, 28.97, 32.46, 13.93, 35.49, 19.36, 19.6 ,  1.44, 24.1 ,
       33.02, 20.62, 26.89, 21.29, 31.95, 29.74, 13.93, 13.82, 19.76,
       21.54, 20.87, 23.63, 28.8 , 23.64,  6.95, 22.2 , -6.82, 16.97,
       16.77, 25.44, 14.95,  3.72, 15.03, 16.91, 21.46, 31.66, 30.72,
       23.79, 22.19, 13.76, 18.47, 18.15, 30.6 , 27.49, 11. , 17.26,
       22.49, 16.53, 29.49, 22.86, 28.68, 20.38, 19.69, 22.55, 27.23,
       24.86, 20.2 , 29.14,  7.43,  5.85, 25.35, 38.73, 23.94, 25.28,
       20.11, 19.75, 25.07, 35.16, 27.32, 27.26, 31.4 , 16.55, 14.3 ,
       23.77,  7.65, 23.35, 21.37, 26.12, 25.32, 13.12, 17.67, 36.2 ,
       20.5 , 27.95, 22.46, 10.15, 31.24, 20.85, 27.36, 30.53])

# Step 6 : model accuracy
from sklearn.metrics import mean_absolute_error, mean_absolute_percentage_error, mean_squared_error

[13] Python

mean_absolute_error(y_test,y_pred)

[14] Python
3.155018927602485

[15] Python
```


03/05/2024

Logistic Regression

Lab - 4

Implement Logistic Regression

- 1) Import necessary libraries.
`import pandas as pd`
`import numpy as np`
`from sklearn, linear model`
`import Logistic Regression`
- 2) Load the dataset
`data = load - iris`
`X = iris . data`
`y = iris . target`
- 3) Create model and train
`model = LogisticRegression()`
`model . fit (df . data , df . target)`
- 4) Predict Value
`model . predict (df . data)`

```
D import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)

import matplotlib.pyplot as plt

# Input data files are available in the "../input/" directory.
# For example, running this (by clicking run or pressing Shift+Enter) will list the files in the input directory

import os


data = pd.read_csv("../content/data.csv")


data.drop(['Unnamed: 32', "id"], axis=1, inplace=True)
data.diagnosis = [1 if each == "M" else 0 for each in data.diagnosis]
y = data.diagnosis.values
x_data = data.drop(['diagnosis'], axis=1)


# Assuming x_data is a numpy array or pandas Dataframe
X = (x_data - np.min(x_data)) / (np.max(x_data) - np.min(x_data))


from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.15, random_state=42)

X_train = X_train.T
X_test = X_test.T
y_train = y_train.T
y_test = y_test.T

print("X train: ", X_train.shape)
print("X test: ", X_test.shape)
print("y train: ", y_train.shape)
print("y test: ", y_test.shape)

... X train: (30, 483)
X test: (30, 86)
y train: (483,)
y test: (86,)
```

```
D def initialize_weights_and_bias(dimension):
    w = np.full((dimension,1),0.01)
    b = 0.0
    return w, b


def sigmoid(z):
    y_head = 1/(1+np.exp(-z))
    return y_head


def forward_backward_propagation(w,b,X_train,y_train):
    # Forward propagation
    z = np.dot(w.T,X_train) + b
    y_head = sigmoid(z)
    loss = -y_train*np.log(y_head)-(1-y_train)*np.log(1-y_head)
    cost = (np.sum(loss))/X_train.shape[1] # X_train.shape[1] is for scaling
    # backward propagation
    derivative_weight = (np.dot(X_train,((y_head-y_train).T)))/X_train.shape[1] # X_train.shape[1] is for scaling
    derivative_bias = np.sum(y_head-y_train)/X_train.shape[1] # X_train.shape[1] is for scaling
    gradients = {'derivative_weight': derivative_weight, 'derivative_bias': derivative_bias}
    return cost,gradients


def update(w, b, X_train, y_train, learning_rate,number_of_iteraton):
    cost_list = []
    cost_list2 = []
    index = []
    # updating (learning) parameters is number_of_iteraton times
    for i in range(number_of_iteraton):
        # make forward and backward propagation and find cost and gradients
        cost,gradients = forward_backward_propagation(w,b,X_train,y_train)
        cost_list.append(cost)
        w = w - learning_rate * gradients['derivative_weight']
        b = b - learning_rate * gradients['derivative_bias']
        if i % 10 == 0:
            cost_list2.append(cost)
            index.append(i)
            print ("Cost after iteration %i: %f" % (i, cost))
    # we update (learn) parameters weights and bias
    parameters = {"weight": w,"bias": b}
    plt.plot(index,cost_list2)
    plt.xticks(index,rotations='vertical')
    plt.xlabel("Number of Iteration")
    plt.ylabel("Cost")
    plt.show()
    return parameters, gradients, cost_list
```

```
D def predict(w,b,X_test):
    # X_test is a input for forward propagation
    z = sigmoid(np.dot(w.T,X_test)+b)
    Y_prediction = np.zeros((1,X_test.shape[1]))
    # if z is bigger than 0.5, our prediction is sign one (y_head=1),
    # if z is smaller than 0.5, our prediction is sign zero (y_head=0),
    for i in range(z.shape[1]):
        if z[0,i]<= 0.5:
            Y_prediction[0,i] = 0
        else:
            Y_prediction[0,i] = 1

    return Y_prediction


(pip install mymodule

Collecting mymodule
  Downloading myModule-1.0.0.tar.gz (787 bytes)
  Preparing metadata (setup.py) ... done
Building wheels for collected packages: mymodule
  Building wheel for mymodule (setup.py) ... done
  Created wheel for mymodule: filename=mymodule-1.0.0-py3-none-any.whl size=1413 sha256=2cdcfae4d015f4dc631a533b7e726b48c29061c95d1a6a58b070071cb04054f3f
  Stored in directory: /root/.cache/pip/wheels/F4/64/ed/4ac2e7514f3e85b8aea309dff4c236817b10c4d0809e551108
Successfully built mymodule
Installing collected packages: mymodule
Successfully installed mymodule-1.0.0
```

```

def sigmoid(z):
    return 1 / (1 + np.exp(-z))

def initialize_weights_and_bias(dim):
    w = np.zeros(dim, 1)
    b = 0
    return w, b

def compute_cost(w, b, x, y):
    m = x.shape[1]
    a = sigmoid(np.dot(w.T, x) + b)
    cost = -1 / m * np.sum(y * np.log(a) + (1 - y) * np.log(1 - a))
    return cost

def propagate(w, b, x, y):
    m = x.shape[1]
    a = sigmoid(np.dot(w.T, x) + b)
    dw = 2 / m * np.dot((a - y).T)
    db = 1 / m * np.sum(a - y)
    return dw, db

def logistic_regression(x_train, y_train, x_test, y_test, learning_rate, num_iterations):
    # Initialization
    dimension = x_train.shape[0] # Number of features
    w, b = initialize_weights_and_bias(dimension)
    costs = []

    # Gradient Descent
    for i in range(num_iterations):
        # Forward and Backward Propagation
        dw, db = propagate(w, b, x_train, y_train)

        # Update parameters
        w += learning_rate * dw
        b += learning_rate * db

        # Record the costs
        if i % 200 == 0:
            cost = compute_cost(w, b, x_train, y_train)
            costs.append(cost)
            print("Cost after iteration (%i): (%f)" % (i, cost))

    # Evaluate model
    y_prediction_train = predict(w, b, x_train)
    y_prediction_test = predict(w, b, x_test)

    train_accuracy = 100 - np.mean(np.abs(y_prediction_train - y_train)) * 100
    test_accuracy = 100 - np.mean(np.abs(y_prediction_test - y_test)) * 100

    print("Train accuracy: (%f) %%" % format(train_accuracy))
    print("Test accuracy: (%f) %%" % format(test_accuracy))

    return w, b

# Assuming you have defined the predict function
# def predict(w, b, x):
# ...

# Assuming you have defined x_train, y_train, x_test, y_test, learning_rate, and num_iterations
logistic_regression(x_train, y_train, x_test, y_test, learning_rate=0.01, num_iterations=100)

```

```

... Cost after iteration 0: 0.6782740168052536
Train accuracy: 80.74534161490683 %
Test accuracy: 81.3953488372093 %

```

```

... (array([[ 1.77806654e-02],
 [ 1.10160388e-02],
 [ 1.27806976e-01],
 [ 1.95749649e+00],
 [ 1.85931875e-05],
 [ 2.68863405e-04],
 [ 4.89020438e-04],
 [ 2.63106883e-04],
 [ 3.49357933e-05],
 [-2.02145931e-05],
 [ 1.25690784e-03],
 [-3.98285024e-04],
 [ 8.96937014e-03],
 [ 2.02426962e-01],
 [-3.60718647e-06],
 [ 4.19190446e-05],
 [ 6.03411729e-05],
 [ 2.00740406e-05],
 [-6.24803672e-06],
 [ 6.24944700e-07],
 [ 2.79506973e-02],
 [ 1.99326360e-02],
 [ 1.98774929e-01],
 [ 3.39189908e+00],
 [ 5.79135019e-05],

...
 [ 1.25862280e-03],
 [ 4.60695564e-04],
 [ 1.89671301e-04],
 [ 3.52490835e-05]]),
-1.5161875221606185)

```

24/05/2024

Build Support vector machine model for a given dataset

Algorithm:

classmate
Date _____
Page _____

Lab 5

Q) Build an SVM model for a given dataset

Steps

- (i) Input the dataset to be trained (or dataset)
- (ii) Convert the dataset into dataframe (note). Add the target column to the dataframe
- (iii) Print the first few rows of dataframe to input the data. Create a scatter plot to visualize the relationship between target line and model
- (iv) Split the dataset into training and testing sets. Use 70% for training and 30% for testing
- (v) Create and train the SVM classifier with a linear kernel. Train the dataset using training data


```
from google.colab import drive
drive.mount('/content/drive')

#routed at /content/drive

import seaborn as sns
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import plotly.express as px

df = pd.read_csv('/content/drive/MyDrive/breast-cancer.csv')
df.head()
```

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	...	radius_worst	texture_worst	perimeter_worst	area_worst	smoothness_worst	compactness_worst	concavity_worst	concave points_worst	symmetry_worst	fractal_dimension_worst
0	842302	M	17.99	10.39	122.80	1001.0	0.11840	0.27760	0.3001	0.1410	...	25.38	17.33	184.80	2019.0	0.1622	0.6555	0.7119	0.2554	0.4801	0.11890
1	842517	M	20.57	17.77	132.90	1216.0	0.08474	0.07964	0.0069	0.07017	...	34.90	23.41	198.80	1956.0	0.1238	0.1896	0.2416	0.1890	0.2750	0.09002
2	8430003	M	16.69	21.25	130.00	1003.0	0.10960	0.1076	0.1599	0.12790	...	23.57	25.53	152.50	1739.0	0.1444	0.4245	0.4504	0.2430	0.3513	0.06758
3	8434301	M	11.42	20.39	77.58	361.1	0.14250	0.26390	0.2814	0.10520	...	14.91	26.50	98.87	597.7	0.2068	0.8563	0.6889	0.2575	0.6638	0.17300
4	8435802	M	20.29	14.34	135.10	1257.0	0.10030	0.13280	0.1280	0.10430	...	22.54	16.67	152.20	1575.0	0.1374	0.2090	0.4000	0.1625	0.2364	0.07678

5 rows x 22 columns

```
df.drop('id', axis=1, inplace=True) #drop redundant columns

df.drop('id', axis=1, inplace=True) #drop redundant columns

df.describe().T
```

	count	mean	std	min	25%	50%	75%	max
radius_mean	560	14.127292	3.540469	6.951000	11.700000	13.200000	15.700000	26.11000
texture_mean	560	19.28949	4.201026	9.710000	14.100000	16.840000	21.800000	39.30000
perimeter_mean	560	91.699033	24.59991	43.700000	75.100000	86.240000	104.100000	198.80000
area_mean	560	654.89104	351914.129	143.500000	420.300000	551.100000	762.700000	2501.00000
smoothness_mean	560	0.096300	0.014064	0.052500	0.086370	0.095970	0.105300	0.16340
compactness_mean	560	0.104341	0.052013	0.019300	0.066200	0.092630	0.130400	0.34540
concavity_mean	560	0.008799	0.079720	0.000000	0.029590	0.061540	0.130700	0.42480
concave points_mean	560	0.046919	0.038003	0.000000	0.020710	0.033200	0.074000	0.26180
symmetry_mean	560	0.181162	0.027414	0.100000	0.161500	0.175000	0.195700	0.34600
fractal_dimension_mean	560	0.062798	0.007090	0.048960	0.057700	0.061540	0.066100	0.09744
radius_worst	560	0.405172	0.277313	0.111900	0.238400	0.324300	0.478900	2.87300
texture_worst	560	1.216853	0.551648	0.340200	0.833900	1.100000	1.474000	4.88900
perimeter_worst	560	236.6059	202.1855	0.757000	1.800000	2.887000	3.357000	21.90000
area_worst	560	48.179779	45.491096	6.850000	19.890000	24.130000	48.190000	542.25000
smoothness_worst	560	0.007041	0.003000	0.001713	0.005169	0.005380	0.006146	0.03113
compactness_worst	560	0.025478	0.017908	0.002252	0.010300	0.020480	0.032480	0.13540
concavity_worst	560	0.031894	0.030186	0.000000	0.010500	0.025890	0.042050	0.36800
concave points_worst	560	0.011796	0.006170	0.000000	0.007638	0.010630	0.014710	0.05279
symmetry_worst	560	0.020542	0.008206	0.007882	0.015180	0.018730	0.023480	0.07995
fractal_dimension_worst	560	0.003795	0.003246	0.000000	0.002288	0.003187	0.004558	0.02684
radius_worst	560	15.181900	4.032842	7.220000	10.010000	14.970000	19.790000	26.04000
perimeter_worst	560	107.261213	33.802542	50.410000	84.110000	97.660000	125.400000	251.20000
area_worst	560	880.881128	569.356993	185.200000	515.300000	685.500000	1084.000000	4254.00000
smoothness_worst	560	0.132399	0.022832	0.071170	0.116600	0.131300	0.146000	0.22280
compactness_worst	560	0.254265	0.187336	0.027580	0.147200	0.211900	0.339100	1.05800
concavity_worst	560	0.071488	0.248664	0.000000	0.114500	0.229700	0.352900	1.25500
concave points_worst	560	0.134806	0.065732	0.000000	0.064030	0.092630	0.161300	0.26100
symmetry_worst	560	0.200076	0.061897	0.159500	0.203600	0.382300	0.317900	0.66380
fractal_dimension_worst	560	0.083045	0.018081	0.055040	0.071480	0.080040	0.092080	0.20750

```
df['diagnosis'] = (df['diagnosis'] == 'M').astype(int) #encode the label into 0/1

corr = df.corr()
```

```

b
# Get the absolute value of the correlation
cor_target = abs(corf['diagnosis'])

# Select highly correlated features (threshold = 0.2)
relevant_features = cor_target[cor_target>0.2]

# Collect the names of the features
names = [index for index, value in relevant_features.items()]

# Drop the target variable from the results
names.remove('diagnosis')

# Display the results
print(names)

Python

[ 'radius_mean', 'texture_mean', 'perimeter_mean', 'area_mean', 'smoothness_mean', 'compactness_mean', 'concavity_mean', 'concave points_mean', 'symmetry_mean', 'radius_se', 'perimeter_se', 'area_se', 'compactness_se', 'concavity_se', 'concave points_se', 'radius_worst', 'texture_worst', 'perimeter_worst', 'area_worst', 'smoothness_worst', '
'

X = df[names].values
y = df['diagnosis']

Python

def scale(X):
    """
    Standardizes the data in the array X.

    Parameters:
    X (numpy.ndarray): Features array of shape (n_samples, n_features).

    Returns:
    numpy.ndarray: The standardized features array.
    """
    # Calculate the mean and standard deviation of each feature
    mean = np.mean(X, axis=0)
    std = np.std(X, axis=0)

    # Standardize the data
    X = (X - mean) / std
    return X

Python

X = scale(X)

Python

def train_test_split(X, y, random_state=1, test_size=0.2):
    """
    Splits the data into training and testing sets.

    Parameters:
    X (numpy.ndarray): Features array of shape (n_samples, n_features).
    y (numpy.ndarray): Target array of shape (n_samples,).
    random_state (int): Seed for the random number generator. Default is 42.
    test_size (float): Proportion of samples to include in the test set. Default is 0.2.

    Returns:
    Tuple(numpy.ndarray): A tuple containing X_train, X_test, y_train, y_test.
    """
    # Get number of samples
    n_samples = X.shape[0]

    # Set the seed for the random number generator
    np.random.seed(random_state)

    # Shuffle the indices
    shuffled_indices = np.random.permutation(np.arange(n_samples))

    # Determine the size of the test set
    test_size = int(n_samples * test_size)

    # Split the indices into test and train
    test_indices = shuffled_indices[test_size:]
    train_indices = shuffled_indices[:test_size]

    # Split the features and target arrays into test and train
    X_train, X_test = X[train_indices], X[test_indices]
    y_train, y_test = y[train_indices], y[test_indices]

    return X_train, X_test, y_train, y_test

Python

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state=42) # Split the data into training and validation

Python

class SVM:
    """
    A Support Vector Machine (SVM) implementation using gradient descent.

    Parameters:
    iterations : int, default=1000
        The number of iterations for gradient descent.
    lr : float, default=0.01
        The learning rate for gradient descent.
    lambda_ : float, default=0.01
        The regularization parameter.

    Attributes:
    lambda_ : float
        The regularization parameter.
    iterations : int
        The number of iterations for gradient descent.
    lr : float
        The learning rate for gradient descent.
    w : numpy array
        The weights.
    b : float
        The bias.

    Methods:
    """
    def __init__(self, iterations=1000, lr=0.01, lambda_=0.01):
        """
        Initializes the SVM model.

        Parameters:
        iterations : int, default=1000
            The number of iterations for gradient descent.
        lr : float, default=0.01
            The learning rate for gradient descent.
        lambda_ : float, default=0.01
            The regularization parameter.
        """
        self.lambda_ = lambda_
        self.iterations = iterations
        self.lr = lr
        self.w = None
        self.b = None

    def initialize_parameters(self, X):
        """
        Initializes the weights and bias.

        Parameters:
        X : numpy array
            The input data.
        """

```



```

y, b = y.shape
self.w = np.zeros(n)
self.b = 0

def gradient_descent(self, X, y):
    """
    Updates the weights and bias using gradient descent.

    Parameters:
    -----
    X : numpy array
        The input data.
    y : numpy array
        The target values.

    """
    X = np.where(X == 0, -1, 1)
    for i, x in enumerate(X):
        if x[i] * (np.dot(x, self.w) - self.b) > 1:
            dw = 2 * self.lmbdas * self.w
            db = 0
        else:
            dw = 2 * self.lmbdas * self.w - np.dot(x, x[i])
            db = x[i]
        self.update_parameters(dw, db)

def update_parameters(self, dw, db):
    """
    Updates the weights and bias.

    Parameters:
    -----
    dw : numpy array
        The change in weights.
    db : float
        The change in bias.

    """
    self.w = self.w + self.lr * dw
    self.b = self.b - self.lr * db
def fit(self, X, y):
    """
    Fits the SVM to the data.

    Parameters:
    -----
    X : numpy array
        The input data.
    y : numpy array
        The target values.

    """
    self.initialize_parameters(X)
    for i in range(self.iterations):
        self.gradient_descent(X, y)

def predict(self, X):
    """
    Predicts the class labels for the test data.

    Parameters:
    -----
    X : array-like, shape (n_samples, n_features)
        The input data.

    Returns:
    -----
    y_pred : array-like, shape (n_samples,)
        The predicted class labels.

    """
    # get the output
    output = np.dot(X, self.w) - self.b
    # get the sign of the labels depending on if it's greater/less than zero
    label_signs = np.sign(output)
    # set predictions to 0 if they are less than or equal to -1 also set them to 1
    predictions = np.where(label_signs < -1, 0, 1)
    return predictions

```

```

def accuracy(y_true, y_pred):
    """
    Computes the accuracy of a classification model.

    Parameters:
    -----
    y_true (numpy array): A numpy array of true labels for each data point.
    y_pred (numpy array): A numpy array of predicted labels for each data point.

    Returns:
    -----
    float: The accuracy of the model.

    """
    total_samples = len(y_true)
    correct_predictions = np.sum(y_true == y_pred)
    return correct_predictions / total_samples

```

```

model = SVM()
model.fit(X_train, y_train)
predictions = model.predict(X_test)

accuracy(y_test, predictions)

```

0.821366834957522

```

model.predict([-0.47090437, -0.10947104, -0.44328055, -0.49289378, 0.13411435,
-0.42769851, -0.38046141, 0.27623152, 0.41384307, -0.42724206,
-0.10269565, 0.12885552, 0.75531157, 0.16012869, -0.21193322,
0.10863051, 0.16046739, -0.13333337, 0.35626257, 0.44859425,
-0.18478885, -0.40342212, -0.10595318, 0.18738443, 0.15965764])

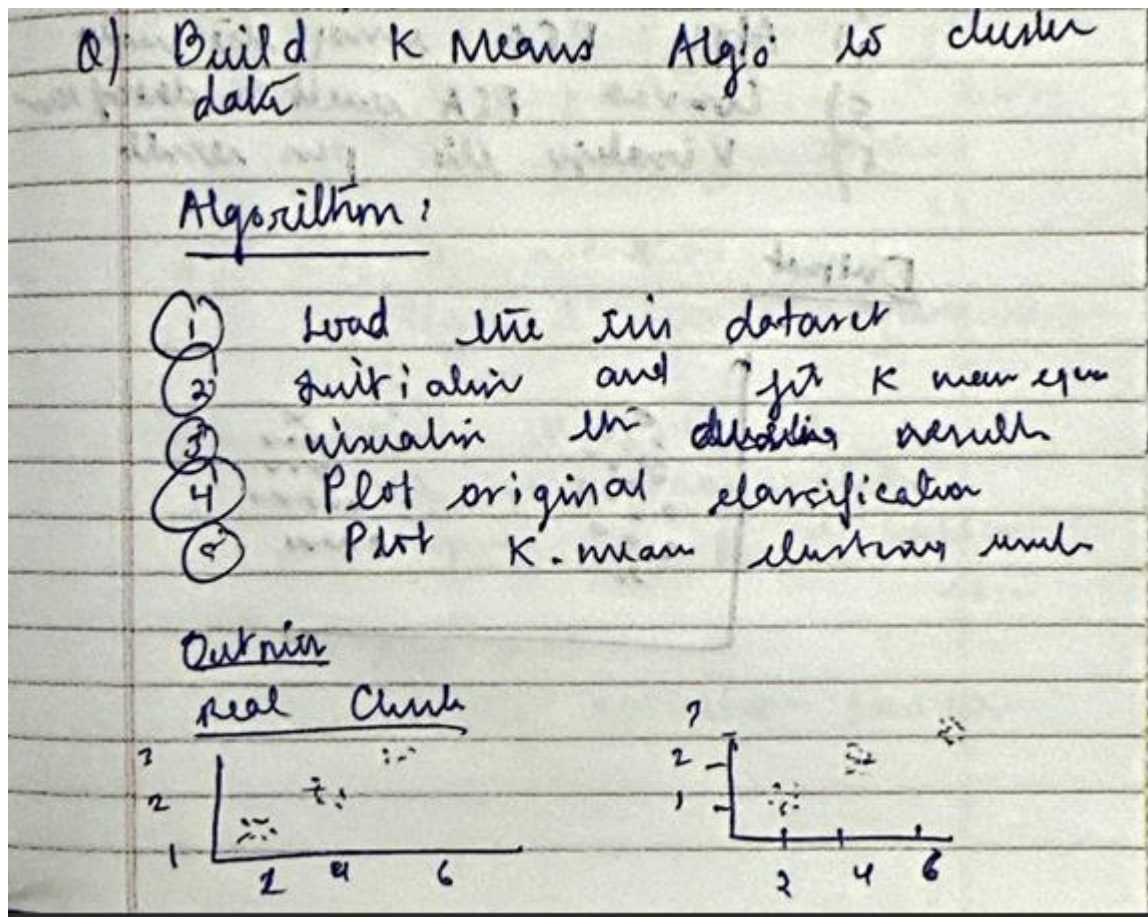
```

array(0)

24/05/2024

Build k-Means algorithm to cluster a set of data stored in a .CSV file.

Algorithm:



```
from glob import glob
import drive
drive.mount('/content/drive')
```

```
print os.listdir('/content/drive')
```

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import plotly.express as px
import sklearn as sns
import plotly.graph_objects as go
```

```
iris = pd.read_csv('/content/drive/MyDrive/iris.csv') #load data
iris.drop('id', inplace=True, axis=1) #drop id column
```

```
X = iris.iloc[:,1:5] #set our training data
```

```
y = iris.iloc[:,0] #we'll use this just for visualization as clustering doesn't require labels
```

```
class KMeans:
    """
    K-Means clustering algorithm implementation.

    Parameters:
    k (int): Number of clusters

    Attributes:
    k (int): Number of clusters
    centroids (numpy.ndarray): Array containing the centroids of each cluster

    Methods:
    __init__(self, k): Initialize the KMeans instance with the specified number of clusters.
    initialize_centroids(self, X): Initialize the centroids for each cluster by selecting k random points from the dataset.
    assign_points_centroids(self, X): Assigns each point in the dataset to the nearest centroid.
    compute_mean(self, X, points): Computes the mean of the points assigned to each centroid.
    fit(self, X, iterations=10): Clusters the dataset using the K-Means algorithm.

    def __init__(self, k):
        assert k > 0, "k should be a positive integer."
        self.k = k

    def initialize_centroids(self, X):
        assert X.shape[0] >= self.k, "Number of data points should be greater than or equal to k."

        randomized_X = np.random.permutation(X.shape[0])
        centroid_ids = randomized_X[self.k:] # get the indices for the centroids
        self.centroids = X[centroid_ids] # assign the centroids to the selected points

    def assign_points_centroids(self, X):
        """
        Assign each point in the dataset to the nearest centroid.

        Parameters:
        X (numpy.ndarray): dataset to cluster

        Returns:
        numpy.ndarray: array containing the index of the centroid for each point
        """
        X = np.expand_dims(X, axis=1) # expand dimensions to match shape of centroids
        distance = np.linalg.norm(X - self.centroids, axis=1) # calculate Euclidean distance between each point and each centroid
        points = np.argmin(distance, axis=1) # assign each point to the closest centroid
        assert len(points) == X.shape[0], "Number of assigned points should equal the number of data points."
        return points

    def compute_mean(self, X, points):
        """
        Compute the mean of the points assigned to each centroid.

        Parameters:
        X (numpy.ndarray): dataset to cluster
        points (numpy.ndarray): array containing the index of the centroid for each point

        Returns:
        numpy.ndarray: array containing the new centroids for each cluster
        """
        centroids = np.zeros((self.k, X.shape[1])) # initialize array to store centroids
        for i in range(self.k):
            centroid_mean = X[points == i].mean(axis=0) # calculate mean of the points assigned to the current centroid
            centroids[i] = centroid_mean # assign the new centroid to the mean of its points
        return centroids

    def fit(self, X, iterations=10):
        """
        Cluster the dataset using the K-Means algorithm.

        Parameters:
        X (numpy.ndarray): dataset to cluster
        iterations (int): number of iterations to perform (default=10)

        Returns:
        numpy.ndarray: array containing the final centroids for each cluster
        numpy.ndarray: array containing the index of the centroid for each point
        """
        self.initialize_centroids(X) # initialize the centroids
        for i in range(iterations):
            points = self.assign_points_centroids(X) # assign each point to the nearest centroid
            self.centroids = self.compute_mean(X, points) # compute the new centroids based on the mean of their points

            # assertions for debugging and validation
            assert len(self.centroids) == self.k, "Number of centroids should equal k."
            assert X.shape[1] == self.centroids.shape[1], "Dimensionality of centroids should match input data."
            assert min(points) >= 0, "Cluster index should be less than k."
            assert min(points) <= 0, "Cluster index should be non-negative."

        return self.centroids, points
```

```
X = X.values
```

```
kmeans = KMeans(3)
```

```
centroids, points = kmeans.fit(X, 1000)
```

```
fig = go.Figure()
fig.add_trace(go.Scatter(
    x0[points == 0, 0], y0[points == 0, 1],
    mode='markers', marker_color='indigo', name='iris-setosa'
)))

fig.add_trace(go.Scatter(
    x0[points == 1, 0], y0[points == 1, 1],
    mode='markers', marker_color='coral', name='iris-versicolour'
)))

fig.add_trace(go.Scatter(
    x0[points == 2, 0], y0[points == 2, 1],
    mode='markers', marker_color='forestgreen', name='iris-virginica'
)))

fig.add_trace(go.Scatter(
    x0[centroids[0, 0], 0], y0[centroids[0, 1], 1],
    mode='markers', marker_color='black', marker_size=14, name='Centroids'
)))

fig.update_layout(template='plotly_dark', width=1000, height=500)
```

24/05/2024

Implement Dimensionality reduction using Principle Component Analysis (PCA)

Algorithm:

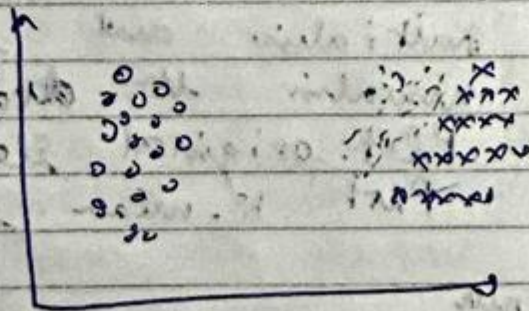
classmate
Date _____
Page _____

Q). Implement dimensionality reduction using principle component analysis method.

Algorithm:

- 1) Import necessary libraries for data handling
- 2) Load the iris dataset
- 3) Standardize the data
- 4) Apply PCA using the model
- 5) Convert PCA result to dataframe
- 6) Visualize the PCA result

Output




```
from google.colab import drive
drive.mount('/content/drive')

Mounted at /content/drive

import seaborn as sns
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import matplotlib.pyplot as plt
import plotly.express as px
import plotly.graph_objects as go
from plotly.subplots import make_subplots

df = pd.read_csv('/content/drive/mydrive/breast-cancer.csv')
df.head()


```

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	radius_worst	texture_worst	perimeter_worst	area_worst	smoothness_worst	compactness_worst	concavity_worst	concave points_worst	symmetry_worst	fractal_dimension_worst	
0	842302	M	17.99	10.38	122.80	1001.0	0.11840	0.27780	0.3001	0.14710	...	25.38	17.33	184.00	2019.0	0.1622	0.6555	0.7119	0.2654	0.4601	0.11890
1	842517	M	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	...	24.90	23.41	198.00	1956.0	0.1128	0.1866	0.2415	0.1880	0.2790	0.03902
2	8430063	M	16.69	21.25	130.00	1203.0	0.10090	0.15090	0.1074	0.10290	...	23.57	25.53	152.50	1709.0	0.1444	0.4245	0.4504	0.2430	0.3613	0.06758
3	8434001	M	11.42	20.38	77.58	386.1	0.14250	0.28300	0.2414	0.10520	...	14.91	26.50	98.87	587.7	0.2009	0.8663	0.8699	0.2575	0.6638	0.17300
4	8435802	M	20.50	14.34	135.10	1297.0	0.10030	0.13280	0.1060	0.10430	...	22.54	16.67	152.20	1575.0	0.1374	0.2090	0.4000	0.1625	0.2364	0.07476

5 rows x 20 columns

```
df.drop('id', axis=1, inplace=True) #drop redundant columns

df['diagnosis'] = (df['diagnosis'] == 'M').astype(int) #encode the label into 1/0

corr = df.corr()
```

```
# Get the absolute value of the correlation
cor_target = abs(corr['diagnosis'])

# select highly correlated features (threshold = 0.2)
relevant_features = cor_target[cor_target>0.2]

# Collect the names of the features
names = [index for index, value in relevant_features.items()]

# drop the target variable from the results
names.remove('diagnosis')

# display the results
print(names)

['radius_mean', 'texture_mean', 'perimeter_mean', 'area_mean', 'smoothness_mean', 'compactness_mean', 'concavity_mean', 'concave points_mean', 'symmetry_mean', 'radius_se', 'perimeter_se', 'area_se', 'compactness_se', 'concavity_se', 'concave points_se', 'radius_worst', 'texture_worst', 'perimeter_worst', 'area_worst', 'smoothness_worst', 'compactness_worst', 'concavity_worst', 'concave points_worst', 'symmetry_worst', 'fractal_dimension_worst']

X = df[names].values

class PCA:
    """
    Principal Component Analysis (PCA) class for dimensionality reduction.
    """
    def __init__(self, n_components):
        """
        Constructor method that initializes the PCA object with the number of components to retain.
        """
        self.n_components = n_components

    def fit(self, X):
        """
        Fits the PCA model to the input data and computes the principal components.
        """
        # Compute the mean of the input data along each feature dimension.
        mean = np.mean(X, axis=0)

        # Subtract the mean from the input data to center it around zero.
        X = X - mean

    def transform(self, X):
        """
        Transforms the input data by projecting it onto the principal components.
        """
        # Compute the explained variance ratio for each principal component.
        # Compute the total variance of the input data
        total_variance = np.sum(np.var(X, axis=0))

        # Compute the variance explained by each principal component
        self.explained_variances = eigenvalues[self.n_components:]

        # Compute the explained variance ratio for each principal component
        self.explained_variance_ratio_ = self.explained_variances / total_variance

    def fit_transform(self, X):
        """
        Fits the PCA model to the input data and computes the principal components then
        transforms the input data by projecting it onto the principal components.
        """
        self.fit(X)
        return self.transform(X)
```

```
# Compute the covariance matrix of the centered input data.
cov = np.cov(X.T)

# Compute the eigenvectors and eigenvalues of the covariance matrix.
eigenvalues, eigenvectors = np.linalg.eigh(cov)

# Reverse the order of the eigenvalues and eigenvectors.
eigenvalues = eigenvalues[::-1]
eigenvectors = eigenvectors[:,::-1]

# Keep only the first n_components eigenvectors as the principal components.
self.components = eigenvectors[:,0:self.n_components]

# Compute the explained variance ratio for each principal component.
# Compute the total variance of the input data
total_variance = np.sum(np.var(X, axis=0))

# Compute the variance explained by each principal component
self.explained_variances = eigenvalues[self.n_components:]

# Compute the explained variance ratio for each principal component
self.explained_variance_ratio_ = self.explained_variances / total_variance

def transform(self, X):
    """
    Transforms the input data by projecting it onto the principal components.
    """
    # Compute the explained variance ratio for each principal component.
    # Compute the total variance of the input data
    total_variance = np.sum(np.var(X, axis=0))

    # Compute the variance explained by each principal component
    self.explained_variances = eigenvalues[self.n_components:]

    # Compute the explained variance ratio for each principal component
    self.explained_variance_ratio_ = self.explained_variances / total_variance

    # Transform the input data by projecting it onto the principal components.
    transformed_data = np.dot(X, self.components)

    return transformed_data

def fit_transform(self, X):
    """
    Fits the PCA model to the input data and computes the principal components then
    transforms the input data by projecting it onto the principal components.
    """
    self.fit(X)
    return self.transform(X)
```

```
pc = PCA(2)

pc.fit(X)

pc.explained_variance_ratio_

array([0.8377429, 0.0363806])
```

```
X_transformed = pca.transform(X)
```

Python

```
X_transformed[0:1].shape
```

Python

... (step)

```
fig = plt.scatter(x=X_transformed[:,0], y=X_transformed[:,1])
fig.update_layout(
    title="PCA transformed data for breast cancer dataset",
    xaxis_title="PC1",
    yaxis_title="PC2"
)
fig.show()
```

Python

...

31.05.2024

Build Artificial Neural Network model with back propagation on a given dataset

Algorithm:

Lab - 6

- a) Build an ~~off~~ ANN model with back propagation.

Algorithm :-

- 1) Initialize parameters
 - Normalise input feature matrix 'x'
 - Normalise the output 'y'
 - set hyperparameters ' ' no of epochs, no of neurons
- 2) Define activation function
 - sigmoid function adjustment
- 3) Training the network
 - Forward propagation
 - Compute i/p to hidden layer
 - Add bias
 - apply activation function

4) Back Propagation

- compute error
- compute gradients
- compute delta

5) Update weights and bias

Output :- I/p $\begin{bmatrix} 0.667, 1 \end{bmatrix}$

$\begin{bmatrix} 0.333, 0.556 \end{bmatrix}$

$\begin{bmatrix} 0.1, 0.667 \end{bmatrix}$

Actual Output :- $\begin{bmatrix} 0.92 \end{bmatrix}$ $\begin{bmatrix} 0.86 \end{bmatrix}$

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

Predicted Output : $\begin{bmatrix} 0.809877 \end{bmatrix}$

$\begin{bmatrix} 0.792919 \end{bmatrix}$

$\begin{bmatrix} 0.701134 \end{bmatrix}$


```

import numpy as np
from sklearn.model_selection import train_test_split

db = np.loadtxt('content/duke-breast-cancer.txt')
print("database raw shape (%d,%d)" % np.shape(db))

# Database raw shape (86,7136)

np.random.shuffle(db)
x = db[:, 0]
y = db[:, 1]
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.1)
print(np.shape(x_train), np.shape(x_test))

# (77, 7125) (9, 7125)

hidden_layer = np.zeros(72)
weights = np.random.random((len(x[0]), 72))
output_layer = np.zeros(2)
hidden_weights = np.random.random((72, 2))

def sum_function(weights, index_locked_col, x):
    result = 0
    for i in range(0, len(x)):
        result += x[i] * weights[i][index_locked_col]
    return result

def activate_layer(layer, weights, x):
    for i in range(0, len(layer)):
        layer[i] = 1 / (1 + np.exp(-(2.0 * sum_function(weights, i, x) / 2.0)))

def softmax(layer):
    softmax_output_layer = np.zeros(len(layer))
    for i in range(0, len(layer)):
        denominator = 0
        for j in range(0, len(layer)):
            denominator += np.exp(layer[j] - np.max(layer))
        softmax_output_layer[i] = np.exp(layer[i] - np.max(layer)) / denominator
    return softmax_output_layer

def recalculate_weights(learning_rate, weights, gradient, activation):
    for i in range(0, len(weights)):
        for j in range(0, len(weights[i])):
            weights[i][j] = (learning_rate * gradient[i] * activation[i] + weights[i][j])

def back_propagation(hidden_layer, output_layer, one_hot_encoding, learning_rate, x):
    output_derivative = np.zeros(2)
    output_gradient = np.zeros(2)
    for i in range(0, len(output_layer)):
        output_derivative[i] = (1.0 - output_layer[i]) * output_layer[i]
    for i in range(0, len(output_layer)):
        output_gradient[i] = output_derivative[i] * (one_hot_encoding[i] - output_layer[i])
    hidden_derivative = np.zeros(72)
    hidden_gradient = np.zeros(72)
    for i in range(0, len(hidden_layer)):
        hidden_derivative[i] = (1.0 - hidden_layer[i]) * (1.0 + hidden_layer[i])
    for i in range(0, len(hidden_layer)):
        sum = 0
        for j in range(0, len(output_gradient)):
            sum += output_gradient[j] * hidden_weights[j][i]
        hidden_gradient[i] = sum
    recalculate_weights(learning_rate, hidden_weights, output_gradient, hidden_layer)
    recalculate_weights(learning_rate, weights, hidden_gradient, x)

one_hot_encoding = np.zeros((2,))
for i in range(0, len(one_hot_encoding)):
    one_hot_encoding[i] = 0
training_correct_answers = 0
for i in range(0, len(x_train)):
    activate_layer(hidden_layer, weights, x_train[i])
    activate_layer(output_layer, hidden_weights, hidden_layer)
    output_layer = softmax(output_layer)
    training_correct_answers += 1 if y_train[i] == np.argmax(output_layer) else 0
back_propagation(hidden_layer, output_layer, one_hot_encoding[int(x_train[i]), -1], x_train[i])
print("MLP Correct answers while learning: %d / %d (accuracy = %d)" % (training_correct_answers, len(x_train), "duke breast cancer"))

# MLP Correct answers while learning: 44 / 77 (accuracy = 0.5714285714285714) on duke breast cancer database.

testing_correct_answers = 0
for i in range(0, len(x_test)):
    activate_layer(hidden_layer, weights, x_test[i])
    activate_layer(output_layer, hidden_weights, hidden_layer)
    output_layer = softmax(output_layer)
    testing_correct_answers += 1 if y_test[i] == np.argmax(output_layer) else 0
print("MLP Correct answers while testing: %d / %d (accuracy = %d) on %d database" % (testing_correct_answers, len(x_test), "duke breast cancer"))

# MLP Correct answers while testing: 8 / 9 (accuracy = 0.8888888888888888) on duke breast cancer database

```

31.05.2024

Implement Random forest ensemble method on a given dataset.

Algorithm:

Q) Implement random forest ensemble method

Algorithm:

- 1) Import necessary libraries
- 2) Load and inspect data
- 3) Pre process the data
- 4) Split data to train and test
- 5) Initialize RF classifier and train using fit method
- 6) Make predictions on the sampling using method predict
- 7) Evaluate the model

Output : Accuracy : 0.98

Confusion Matrix $\begin{bmatrix} 23 & 0 & 0 \\ 0 & 19 & 0 \\ 0 & 1 & 12 \end{bmatrix}$

```

from google.colab import drive
drive.mount('/content/drive')

# Mounted at /content/drive

import math
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.express as px

iris = pd.read_csv('/content/drive/MyDrive/iris.csv') #load data
iris.drop('id', inplace=True, axis=1) #drop id column

iris.head().style.background_gradient(cmap=sns.light_palette('seagreen', as_cmap=True))

SepalLengthCm  SepalWidthCm  PetalLengthCm  PetalWidthCm  Species
0      5.100000      3.500000      1.400000      0.200000  Iris-setosa
1      4.900000      3.000000      1.400000      0.200000  Iris-setosa
2      4.700000      3.200000      1.300000      0.200000  Iris-setosa
3      4.600000      3.300000      1.500000      0.200000  Iris-setosa
4      5.000000      3.600000      1.400000      0.200000  Iris-setosa

X_train = iris.iloc[:, 1:-1] #set our training dataframe
y_train = iris.iloc[:, 1] # set our training labels dataframe

fig = px.scatter(iris, 'Species', color_discrete_sequence=['#38a83d', '#808080', '#d3d3d3'], title='Iris data distribution', template='plotly')
fig.show()

```

```

iris['Species'] = iris['Species'].astype('category')
codes = iris['Species'].cat.codes

def train_test_split(X, y, random_state=1, test_size=0.2):
    """
    Splits the data into training and testing sets.

    Parameters:
    X (numpy.ndarray): Features array of shape (n_samples, n_features).
    y (numpy.ndarray): Target array of shape (n_samples,).
    random_state (int): Seed for the random number generator. Default is 42.
    test_size (float): Proportion of samples to include in the test set. Default is 0.2.

    Returns:
    Tuple(numpy.ndarray): A tuple containing X_train, X_test, y_train, y_test.
    """
    # Set number of samples
    n_samples = X.shape[0]

    # Set the seed for the random number generator
    np.random.seed(random_state)

    # Shuffle the indices
    shuffled_indices = np.random.permutation(np.arange(n_samples))

    # Determine the size of the test set
    test_size = int(n_samples * test_size)

    # Split the indices into test and train
    test_indices = shuffled_indices[test_size:]
    train_indices = shuffled_indices[:test_size]

    # Split the features and target arrays into test and train
    X_train, X_test = X[train_indices], X[test_indices]
    y_train, y_test = y[train_indices], y[test_indices]

    return X_train, X_test, y_train, y_test

X = iris.iloc[:, 1:-1].values
y = iris.iloc[:, 1].values.reshape(-1,)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=1)

from sklearn.tree import DecisionTreeClassifier
m = DecisionTreeClassifier()

```

```

class RandomForest:
    """
    A random forest classifier.

    Parameters:
    n_trees : int, default=10
        The number of trees in the random forest.
    max_depth : int, default=5
        The maximum depth of each decision tree in the random forest.
    min_samples : int, default=2
        The minimum number of samples required to split an internal node
        of each decision tree in the random forest.

    Attributes:
    n_trees : int
        The number of trees in the random forest.
    max_depth : int
        The maximum depth of each decision tree in the random forest.
    min_samples : int
        The minimum number of samples required to split an internal node
        of each decision tree in the random forest.
    trees : list of DecisionTreeClassifier
        The decision trees in the random forest.
    """

    def __init__(self, n_trees=10, max_depth=5, min_samples=2):
        """
        Initialize the random forest classifier.

        Parameters:
        n_trees : int, default=10
            The number of trees in the random forest.
        max_depth : int, default=5
            The maximum depth of each decision tree in the random forest.
        min_samples : int, default=2
            The minimum number of samples required to split an internal node
            of each decision tree in the random forest.
        """
        self.n_trees = n_trees
        self.max_depth = max_depth
        self.min_samples = min_samples
        self.trees = []

    def fit(self, X, y):
        """
        Build a random forest classifier from the training set (X, y).

        Parameters:
        X : array-like of shape (n_samples, n_features)
            The training input samples.
        y : array-like of shape (n_samples,)
            The target values.

        Returns:
        """

```

```

    Returns
    -----
    self : object
    ...
    Returns self.

    # Create an empty list to store the trees.
    self.trees = []
    # Concatenate x and y into a single dataset.
    dataset = np.concatenate((x, y.reshape(-1, 1)), axis=1)
    # Loop over the number of trees.
    for _ in range(self.n_trees):
        # Create a decision tree instance.
        tree = DecisionTreeClassifier(max_depth=self.max_depth, min_samples_split=self.min_samples)
        # Sample from the dataset with replacement (bootstrapping).
        dataset_sample = self.bootstrap_sample(dataset)
        # Get the x and y samples from the dataset sample.
        x_sample, y_sample = dataset_sample[:, :-1], dataset_sample[:, -1]
        # Fit the tree to the x and y samples.
        tree.fit(x_sample, y_sample)
        # Store the tree in the list of trees.
        self.trees.append(tree)

    return self

def bootstrap_sample(self, dataset):
    """
    Bootstrap the dataset by sampling from it with replacement.

    Parameters
    -----
    dataset : array-like of shape (n_samples, n_features + 1)
        The dataset to bootstrap.

    Returns
    -----
    dataset_sample : array-like of shape (n_samples, n_features + 1)
        The bootstrapped dataset sample.
    """
    # Get the number of samples in the dataset.
    n_sample = dataset.shape[0]
    # Generate random indices to index into the dataset with replacement.
    np.random.seed(2)
    indices = np.random.choice(n_sample, n_sample, replace=True)
    # Return the bootstrapped dataset sample using the generated indices.
    dataset_sample = dataset[indices]
    return dataset_sample

def most_common_label(self, y):
    """
    Return the most common label in an array of labels.

    Parameters
    -----
    y : array-like of shape (n_samples,)
        The array of labels.

    Returns
    -----
    most_occurring_value : str or float
        The most common label in the array.

    y = list(y)
    # get the highest present class in the array
    most_occurring_value = max(y, key=y.count)
    return most_occurring_value

def predict(self, X):
    """
    Predict class for X.

```

```

    Parameters
    -----
    X : array-like of shape (n_samples, n_features)
        The input samples.

    Returns
    -----
    majority_predictions : array-like of shape (n_samples,)
        The predicted classes.
    """
    # Get prediction from each tree in the tree list on the test data
    predictions = np.array([tree.predict(X) for tree in self.trees])
    # get prediction for the same sample from all trees for each sample in the test data
    preds = np.swapaxes(predictions, 0, 1)
    # Get the most voted value by the trees and store it in the final prediction array
    majority_predictions = np.array([self.most_common_label(pred) for pred in preds])
    return majority_predictions

```

```

def accuracy(y_true, y_pred):
    """
    Computes the accuracy of a classification model.

    Parameters
    -----
    y_true (numpy array): A numpy array of true labels for each data point.
    y_pred (numpy array): A numpy array of predicted labels for each data point.

    Returns
    -----
    float: The accuracy of the model, expressed as a percentage.
    """
    y_true = y_true.flatten()
    total_samples = len(y_true)
    correct_predictions = np.sum(y_true == y_pred)
    return (correct_predictions / total_samples)

```

```

from sklearn.preprocessing import LabelEncoder

label_encoder = LabelEncoder()
x_train_encoded = label_encoder.fit_transform(x_train)
y_test_encoded = label_encoder.transform(y_test)

```

```

/usr/local/lib/python3.10/dist-packages/sklearn/preprocessing/_label.py:116: DataConversionWarning:
A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel().

/usr/local/lib/python3.10/dist-packages/sklearn/preprocessing/_label.py:124: DataConversionWarning:
A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel().

```

```

from sklearn.preprocessing import LabelEncoder

label_encoder = LabelEncoder()
x_train_encoded = label_encoder.fit_transform(x_train.ravel())
y_test_encoded = label_encoder.transform(y_test.ravel())
model = RandomForestClassifier(n_estimators=10, max_depth=2)
model.fit(x_train, y_train_encoded)

predictions = model.predict(x_test)
accuracy(y_test_encoded, predictions)

```

```

... 0.9333333333333333

```

```

from sklearn.tree import DecisionTreeClassifier

# Create and train the decision tree model
dt = DecisionTreeClassifier()
dt.fit(x_train, y_train_encoded)

# New predictions on the test data
predictions = dt.predict(x_test)

# Calculate accuracy
accuracy(y_test_encoded, predictions)

```

```

100
... 0.9

```


31.05.2024

Implement Boosting ensemble method on a given dataset.

Algorithm:

Q) Implement ensemble method

Algo:

- Import libraries
- Load the dataset
- Data preprocessing under supervision at features and labels
- Split data into train, test
- Initialize the model as classifier with specified no of estimators and base optimizer
- Train model using training data
- Make prediction using trained model
- Evaluate the model

Result

Model's accuracy: 0.9833

```

# Compute error rate, along with:
def compute_error(y, y_pred, w_1):
    """
    Calculate the error rate of a weak classifier. Arguments:
    y: actual target values
    y_pred: predicted value by weak classifier
    w_1: individual weights for each observation

    Note that all arrays should be the same length
    """
    return (sum(w_1 * (np.not_equal(y, y_pred)).astype(int))) / sum(w_1)

def compute_alpha_error():
    """
    Calculate the weight of a weak classifier. This is called
    alpha in chapter 34 of The Elements of Statistical Learning. Arguments:
    error: error rate from weak classifier.
    """
    return np.log(1 - error) / error

def update_weights(w_1, alpha, y, y_pred):
    """
    Update individual weights w_1 after a boosting iteration. Arguments:
    w_1: individual weights for each observation
    y: actual target values
    y_pred: predicted value by weak classifier
    alpha: weight of weak classifier used to estimate y_pred
    """
    return w_1 * np.exp(alpha * (np.not_equal(y, y_pred)).astype(int))

```

Python

```

# Define AdaBoost class
class AdaBoost:
    def __init__(self):
        self.alphas = []
        self.n = 1
        self.nf = None
        self.training_errors = []
        self.predictions_errors = []

    def fit(self, X, y, n = 100):
        """
        Fit model. Arguments:
        X: independent variables - array-like matrix
        y: target variable - array-like vector
        n: number of boosting rounds. Default is 100 - integer
        """
        # Clear before calling
        self.alphas = []
        self.training_errors = []
        self.nf = n

        # Iterate over n weak classifiers
        for m in range(1, n):
            # Set weights for current boosting iteration
            if m == 0:
                w_1 = np.ones(len(y)) * 1 / len(y) # at m = 0, weights are all the same and equal to 1 / N
            else:
                # (b) Update w_1
                w_1 = update_weights(w_1, alpha_m, y, y_pred)

            # Fit weak classifier on current boosting iteration
            if m == 0:
                w_1 = np.ones(len(y)) * 1 / len(y) # at m = 0, weights are all the same and equal to 1 / N
            else:
                # (b) Update w_1
                w_1 = update_weights(w_1, alpha_m, y, y_pred)

            # (c) Fit weak classifier and predict labels
            Q, R = DecisionTreeClassifier(max_depth = 3)
            Q.fit(X, w_1 * y, sample_weight = w_1)
            y_pred = Q.predict(X)

            self.d_n.append(Q) # Save to list of weak classifiers

            # (d) Compute error
            error_m = compute_error(y, y_pred, w_1)
            self.training_errors.append(error_m)

            # (e) Compute alpha
            alpha_m = compute_alpha_error_m()
            self.alphas.append(alpha_m)

        assert len(self.d_n) == len(self.alphas)

    def predict(self, X):
        """
        Predict using fitted model. Arguments:
        X: independent variables - array-like
        """
        # Initialize dataframe with weak predictions for each observation
        weak_preds = pd.DataFrame(index = range(len(X)), columns = range(self.nf))

        # Predict class label for each weak classifier, weighted by alpha_m
        for m in range(self.nf):
            y_pred_m = self.d_n[m].predict(X) * self.alphas[m]
            weak_preds.loc[:, m] = y_pred_m

        # Calculate final predictions
        y_pred = (1 * np.sign(weak_preds.sum(axis = 1).astype(int)))

        return y_pred

```

Python

Python

```

import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier

# Load data
df = pd.read_csv('content/spambase.data', header = None)

# Column names
names = pd.read_csv('content/spambase.names', sep = '|', skiprows=(9, 11), header = None)
col_names = list(names[0])
col_names.append('spam')

# Review of columns
df.columns = col_names

```

```

# Column names
names = pd.read_csv('content/spambase.names', sep = '|', skiprows=(9, 11), header = None)
col_names = list(names[0])
col_names.append('spam')

# Review of columns
df.columns = col_names

# Convert classes in target variable to {0, 1}
df['spam'] = df['spam'] * 2 - 1

# Train - test split
X_train, X_test, y_train, y_test = train_test_split(df.drop(columns = 'spam'), subset,
                                                df['spam'].values,
                                                train_size = 90%,
                                                random_state = 0)

```

Python

```

# Fit model
ab = AdaBoost()
ab.fit(X_train, y_train, n = 100)

# Predict on test set
y_pred = ab.predict(X_test)

# Review of metrics
from sklearn.metrics import accuracy_score

# Calculate accuracy
accuracy = accuracy_score(y_test, y_pred)
print('Accuracy:', accuracy)

```

Python

```

--- Accuracy: 0.9448341255555555

```