

# **VISVESVARAYA TECHNOLOGICAL UNIVERSITY**

**“JnanaSangama”, Belgaum -590014, Karnataka.**



## **LAB REPORT**

**On**

### **MACHINE LEARNING**

**Submitted by**

**Nishant S(1BM21CS118)**

**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 **Nishant S (1BM21CS118)**, who is bona fide 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

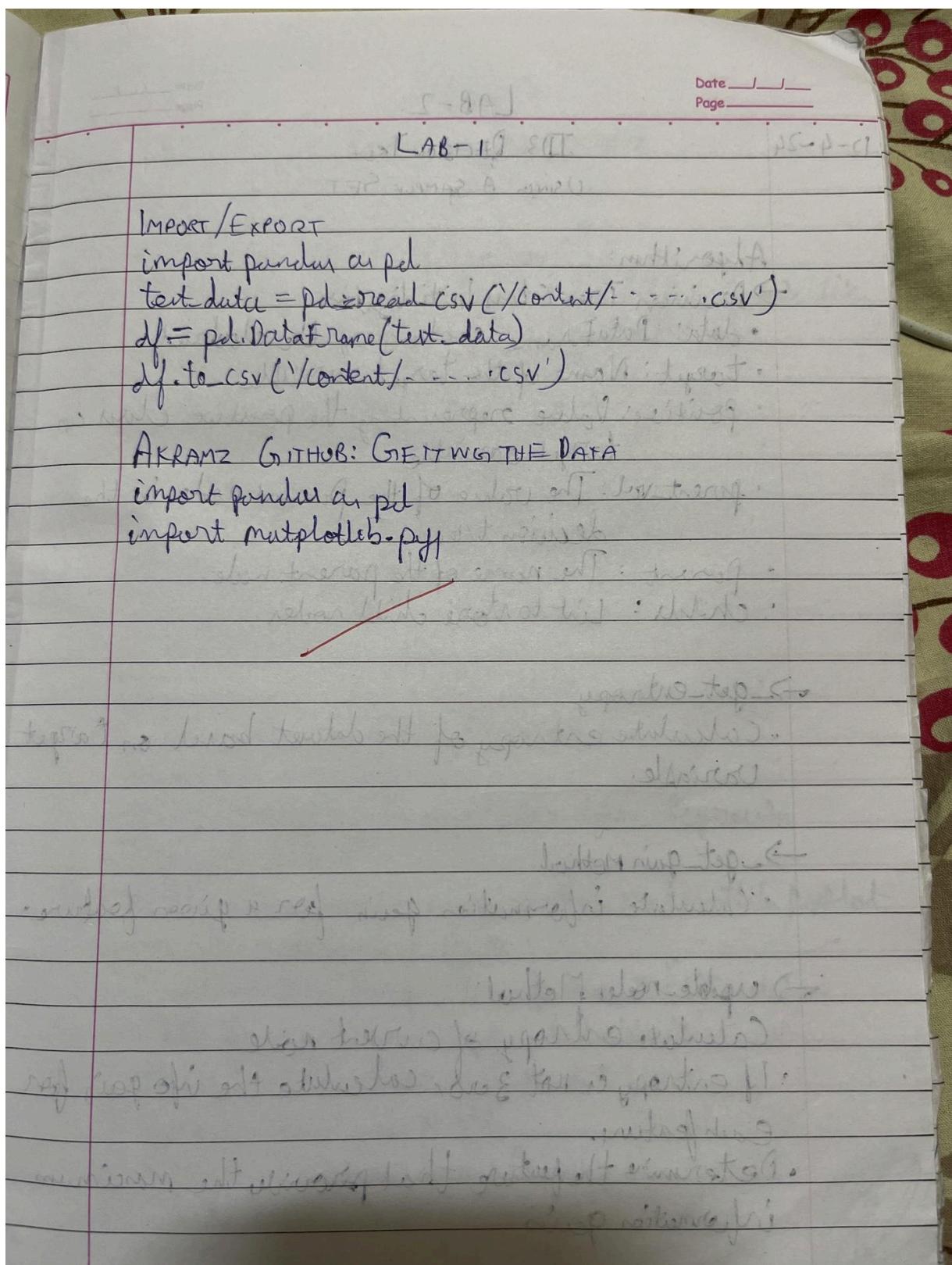
# Index

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

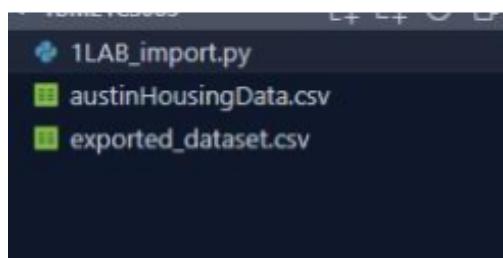


```
import pandas as pd df=pd.read_csv("austinHousingData.csv") print(df.head())
```

```
df.to_csv("exported_dataset.csv")
```

Output:

```
PS C:\Users\student\Documents\1BM21CS083> python -u "c:\Users\student\Documents\1BM21CS083\1LAB_import.py"
      zpid      city    streetAddress zipcode ... numOfBathrooms numOfBedrooms numOfStories      homeImage
0  111373431  pflugerville  14424 Lake Victor Dr  78668 ...          3.0             4               2  111373431_ffcc626843283d3365c11d81b8e65dcdf-p_f...
1  128998438  pflugerville  1184 Strickling Dr  78668 ...          2.0             4               1  128998438_8255c127be8dcf#Ba1a18b75639837088-p_f...
2  2084491383  pflugerville  1488 Fort Dessau Rd  78668 ...          2.0             3               1  2084491383_a2ad649e1a7e090111deca044a11c055-p...
3  128991374  pflugerville  1825 Strickling Dr  78668 ...          2.0             3               1  128991374_b469367a619da85b1f5ceb69b675d88e-p_f...
4  60134862  pflugerville  15005 Donna Jane Loop  78668 ...          3.0             3               2  60134862_btavtta3dF9f111e0092a691873e98ce2-p_F.jpg
[5 rows x 47 columns]
PS C:\Users\student\Documents\1BM21CS083>
```



```
import pandas as pd
url = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"

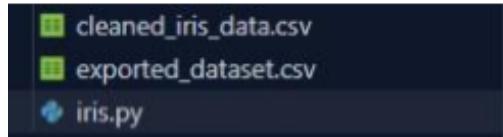
col_names = ["sepal_length_in_cm",
             "sepal_width_in_cm",
             "petal_length_in_cm",
             "petal_width_in_cm",
             "class"]
iris_data = pd.read_csv(url, names=col_names)

print(iris_data.head())

iris_data.to_csv("cleaned_iris_data.csv")
```

Output:

```
PS C:\Users\student\Documents\1BM21CS083> python -u "c:\Users\student\Documents\1BM21CS083\1lab\iris.py"
      sepal_length_in_cm  sepal_width_in_cm  petal_length_in_cm  petal_width_in_cm      class
0              5.1            3.5                1.4            0.2  Iris-setosa
1              4.9            3.0                1.4            0.2  Iris-setosa
2              4.7            3.2                1.3            0.2  Iris-setosa
3              4.6            3.1                1.5            0.2  Iris-setosa
4              5.0            3.6                1.4            0.2  Iris-setosa
PS C:\Users\student\Documents\1BM21CS083>
```



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

② Demonstrate various data preprocessing 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 columns having high no. of null values.
- 2) Use Fillna to replace a NULL values with a specified value.

Encoding categorical data using pd.get\_dummies() which converts categorical data into dummy or indicator variables

```

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

[1]: dataset = pd.read_csv("./content/Data.csv")

[2]: df = pd.DataFrame(dataset)

[3]: df

[4]: Outputs are collapsed --

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

[6]: print(X)

[7]: Outputs are collapsed --

[8]: print(y)

[9]: ... ['No' 'Yes' 'No' 'No' 'Yes' 'Yes' 'No' 'Yes' 'No' 'Yes']

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

[11]: Outputs are collapsed --

```

## Dropna

```

[1]: df1 = df.copy()

[2]: # 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)

[3]: Outputs are collapsed --

[4]: # 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)

[5]: Outputs are collapsed --

```

## Scikit-Learn

```

[1]: X

[2]: Outputs are collapsed --

[3]: 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])

[4]: print(X)

[5]: Outputs are collapsed --

```

## Using Dummies

```

[1]: df2

[2]: Outputs are collapsed --

[3]: pd.get_dummies(df2) # Add to CodeSPT Chat (Ctrl + Shift + E).      Add to chat (Ctrl+L) | Edit highlighted code (Ctrl+I).

[4]: Outputs are collapsed --

```

```
: df2
```

	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	38.0	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	83000.0	No
9	France	37.0	67000.0	Yes

```
: pd.get_dummies(df2)
```

	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	38.0	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	83000.0	False	True	False	True	False
9	37.0	67000.0	True	False	False	False	True

12/04/2024

Decision Tree Algorithm:

12-4-24

## IT'S DECISION TREE

USING A SAMPLE SET

Algorithm:

→ Decision Tree Class Initialization

- data: DataFrame containing the dataset.

- target: Name of the target variable.

- positive: Value representing the positive class in target variable

- parent\_val: The value of the parent node in the decision tree.

- parent: The name of the parent node.

- children: List to store child nodes.

→ get\_entropy

- Calculate entropy of the dataset based on target variable.

→ get\_gain Method

- Calculate information gain for a given feature.

→ update\_node Method:

- Calculate entropy of current node

- If entropy is not zero, calculate the info gain for each feature.

- Determine the feature that provides the maximum information gain

OUTPUT:

Highest information gain =  $0.246 = \text{Outlook}$

12/4/2024 ~~Highest information gain =  $0.971 = \text{rainy}$~~   
~~attribute is wind.~~

## Code and output:

```
[1]: print(f'Entropy of the entire dataset: {find_entropy(data.play)}')
...
Entropy of the entire dataset: 0.94

[2]: entropy_and_infogain(data, 'temp')
...
Outputs are collapsed ...

[3]: entropy_and_infogain(data, 'humidity')
...
Outputs are collapsed ...

[4]: entropy_and_infogain(data, 'windy')
...
Outputs are collapsed ...

#Rainy -outlook

[5]: 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')], 'color': 'black'}, {'selector': 'tr:hover', 'props': [('background-color', 'white'), ('border', '1.5px solid black')]}])
...
Outputs are collapsed ...

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

[7]: entropy_and_infogain(rainy, 'temp')
...
Outputs are collapsed ...

[8]: entropy_and_infogain(rainy, 'humidity')
...
Outputs are collapsed ...

[9]: entropy_and_infogain(rainy, 'windy')
...
Outputs are collapsed ...

[10]: entropy_and_infogain(rainy, 'temp')
...
Outputs are collapsed ...

[11]: entropy_and_infogain(rainy, 'humidity')
...
Outputs are collapsed ...

[12]: entropy_and_infogain(rainy, 'windy')
...
Outputs are collapsed ...

[13]: wind has highest information gain
```

## Output:

```
[1]: 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')], 'color': 'black'}, {'selector': 'tr:hover', 'props': [('background-color', 'white'), ('border', '1.5px solid black')]}])
...

```

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

```
[1] print("Entropy of the entire dataset: {Find_entropy(data.play)}")
...
... Entropy of the entire dataset: 0.94

[2] entropy_and_infogain(data, 'temp')
...

...
... Entropy of temp - hot = 1.0
...

...
... Entropy of temp - mild = 0.918
...

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

[3] entropy_and_infogain(data, 'humidity')
...

...
... Entropy of humidity - high = 0.985
...

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

[4] entropy_and_infogain(data, 'windy')
...

...
... Entropy of windy - False = 0.811
...

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

#Rainy -outlook
```

```

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

```

	rainy	mild	high	False
rainy				
rainy		cool	normal	True
rainy		cool	normal	True
rainy		mild	normal	False
rainy		mild	high	True

```

print("Entropy of the Rainy dataset: " + find_entropy(rainy.play))

```

Entropy of the Rainy dataset: 0.971

```

entropy_and_infogain(rainy, 'temp')

```

	mild	False
mild		
mild		
mild		

Entropy of temp - mild = 0.918

	cool	False
cool		
cool		

Entropy of temp - cool = 1.0  
Information Gain for temp = 0.02

```

entropy_and_infogain(rainy, 'humidity')

```

	high	False
high		
high		

Entropy of humidity - high = 1.0

	normal	False
normal		
normal		
normal		

Entropy of humidity - normal = 0.918  
Information Gain for humidity = 0.082

```

entropy_and_infogain(rainy, 'windy')

```

	False	False
False		
False		
False		

Entropy of windy - False = 0.0

	True	False
True		
True		

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.:

LAB - 4  
LOGISTIC REGRESSION

→ Import necessary libraries  
import pandas as pd  
import numpy as np  
from sklearn.linear\_model import LogisticRegression

→ Load the dataset  
iris = load\_iris()  
X = iris.data  
y = iris.target

→ Create model and train  
model = LogisticRegression()  
model.fit(df.drop(['Outcome']), df['Outcome'])

→ Predict values  
model.predict(df.drop(['Outcome']))

K-NEAREST NEIGHBOUR

→ Import libraries necessary for KNN  
import pandas as pd  
import numpy as np

→ Read dataset  
df = pd.read\_csv('iris.csv')  
test = df[0:1]  
train = df[1:, :].drop(test.index)

→ Create model and predict  
clf = KNeighborsClassifier(n\_neighbors=5, weights='uniform')  
clf.fit(x\_train, y\_train)  
clf.predict(test[1:-1])

Steps:  
→ Import  
→ Convert column  
→ Bind  
→ a real  
→ Input  
→ Split  
→ Create  
→ Kernel  
→ Predict

Output:  
Accuracy

Steps:  
→ Load the  
→ Initialize  
→ Visualize  
→ Plot

Output:

3
2
1
2

## Code and output:

```
[1]: for var in df.columns:  
    print(df[var].value_counts())  
  
[2] Outputs are collapsed ...  
  
[3]: df['Bare_Nuclei'] = pd.to_numeric(df['Bare_Nuclei'], errors='coerce')  
  
[4] df.dtypes  
  
[5] Outputs are collapsed ...  
  
[6]: df.isnull().sum()  
  
[7] Outputs are collapsed ...  
  
[8]: # check 'na' values in the DataFrame  
df.isna().sum()  
  
[9] Outputs are collapsed ...  
  
[10]: # check frequency distribution of 'Bare_Nuclei' column  
df['Bare_Nuclei'].value_counts()  
  
[11] Outputs are collapsed ...  
  
[12]: # check unique values in 'Bare_Nuclei' column  
df['Bare_Nuclei'].unique()  
... array([ 1., 10., 2., 4., 3., 9., 7., nan, 5., 8., 6.])  
  
# check for nan values in "Bare_Nuclei" column ...  
... 16  
  
[13]: # view frequency distribution of values in 'Class' variable  
df['Class'].value_counts()  
  
[14] Outputs are collapsed ...  
  
[15]: import numpy as np  
  
# view summary statistics in numerical variables  
print(round(df.describe(),2))  
  
[16] Outputs are collapsed ...  
  
[17]: X = df.drop(['Class'], axis=1)  
y = df['Class']  
  
[18]: 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)  
  
[19]: X_train.shape, X_test.shape  
  
[20] Outputs are collapsed ...  
  
[21]: for col in X_train.columns:  
    if X_train[col].isnull().mean()>0:  
        print(col, round(X_train[col].isnull().mean(),4))  
  
... Bare_Nuclei 0.0239  
  
...  
  
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	2.02083	0.399506	0.289573	1.119077	-0.546543	1.858937	-0.577774	0.041241	-0.324258
1	1.669451	2.357680	2.304569	-0.624271	3.106879	1.297589	-0.159953	0.041241	-0.324258
2	-1.020205	-0.679581	-0.717925	0.074148	-1.003220	-0.104329	-0.095595	-0.608165	-0.324258
3	-0.175209	-0.026856	-0.046260	-0.624271	-0.546543	-0.665096	-0.159953	0.041241	-0.324258
4	0.239723	-0.353219	-0.382092	-0.274161	-0.546543	-0.665096	-0.577774	-0.283462	-0.324258

```

# import KNeighborsClassifier 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: {:.4f}'.format(accuracy_score(y_test, y_pred)))

Outputs are collapsed ...

```

```

y_pred_train = knn.predict(X_train)

print('Training-set accuracy score: {:.4f}'.format(accuracy_score(y_train, y_pred_train)))

Outputs are collapsed ...

```

## Output:

```

df.shape

```

	Id	Clump_thickness	Uniformity_Cell_Size	Uniformity_Cell_Shape	Marginal_Adhesion	Single_Epithelial_Cell_Size	Bare_Nuclei	Bland_Chromatin	Normal_Nucleoli	Mitoses	Class
0	100025	5	1	1	1	2	1	3	1	1	2
1	100245	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

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

Clump_thickness	count
1	145
5	130
3	108
4	80
10	69
2	58
8	46
6	34
7	23
9	14

Uniformity_Cell_Size	count
1	360
10	67
3	52
2	45
4	40
5	38
8	29
6	27
7	18
9	6

Uniformity_Cell_Shape	count
1	458
4	241

Class	count
2	458
4	241

```

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

[10]: df.dtypes
      
```

Clump_thickness	Dtype
int64	
float64	
int64	
object	

```

[11]: df.isnull().sum()
      
```

Clump_thickness	count
0	0
Uniformity_Cell_Size	0
Uniformity_Cell_Shape	0
Marginal_Adhesion	0
Single_Epithelial_Cell_Size	0
Bare_Nuclei	0
Bland_Chromatin	0
Normal_Nucleoli	0
Mitoses	0
Class	0

Dtype	count
int64	0

```

[12]: # check 'na' values in the DataFrame
      df.isna().sum()
      
```

Clump_thickness	count
0	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	count
int64	0

```

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

[10]
...
Bare_Nuclei
1.0    402
10.0   132
2.0    50
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

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

[11]
...
array([ 1., 10., 2., 4., 3., 9., 7., nan, 5., 8., 6.])

# check for nan values in "Bare_Nuclei" column...
16

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

[12]
...
Class
2    458
4    241
Name: count, dtype: int64

```

```

import numpy as np

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

[13]
...
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            1.00            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
mean     3.44            2.87            2.69
std      2.44            3.05            1.72
min     1.00            1.00            2.00
25%    2.00            1.00            2.00
50%    3.00            1.00            2.00
75%    5.00            4.00            4.00
max    10.00           10.00           10.00

```

```

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

[14]
...
X_train.head()

[15]
...
Clump_Thickness Uniformity_Cell_Size Uniformity_Cell_Shape Marginal_Adhesion Single_Epithelial_Cell_Size Bare_Nuclei Bland_Chromatin Normal_Nucleoli Mitoses
0     2.098983         0.299506       0.289573     1.119077     -0.546543     1.858557     -0.577774     0.041241     -0.324598
1     1.669451         2.257600       2.304569    -0.624271     3.106679     1.297589     -0.159953     0.041241     -0.324598
2     -1.262005        -0.679581     -0.717925     0.074148    -1.003220     -0.104329     -0.995595     -0.608165     -0.324598
3     -0.115209        -0.026956     -0.046260    -0.624271     -0.546543     -0.665096     -0.159953     0.041241     -0.324598
4     0.233723        -0.353219     -0.382092    -0.274161     -0.546543     -0.665096     -0.577774     -0.280462     -0.324598

# import KNeighborsClassifier 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)

[16]
...
KNeighborsClassifier()
KNeighborsClassifier(n_neighbors=3)

y_pred = knn.predict(X_test)
y_pred

[17]
...
array([2, 2, 4, 2, 4, 2, 4, 2, 4, 2, 2, 2, 4, 4, 4, 2, 2, 2, 4, 4, 4,
       2, 2, 4, 2, 2, 4, 4, 2, 2, 2, 2, 2, 2, 2, 4, 2, 2, 2, 2, 2, 2,
       4, 4, 2, 4, 2, 4, 4, 2, 2, 4, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
       4, 2, 2, 4, 2, 2, 2, 4, 4, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
       4, 4, 2, 2, 2, 2, 4, 2, 2, 2, 2, 4, 4, 2, 2, 2, 2, 2, 2, 2, 2, 2,
       4, 4, 2, 2, 2, 2, 4, 4, 4, 4, 2, 2, 2, 2, 4, 4, 4, 4, 4, 4, 4, 2,
       2, 4, 4, 2, 2, 4, 4, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2])

```



03/05/2024

## Linear Regression

LAB-3  
LINEAR REGRESSION

**ALGORITHM:**

- Import data & libraries using `train`
- Visualize of dataset using different plots like heatmap distribution plot scatterplot etc.
- Preprocess and clean the data
- $X_{train}, X_{test}, Y_{train}, Y_{test}$  = `train test split(0.7)`
- Build model  
`lin_mod = LinearRegression()`
- `lin_mod.fit(X_train, Y_train)`

**MULTI REGRESSION**

**ALGORITHM:**

- Import necessary libraries like `LinearRegression`
- Import dataset
- Visualize dataset using `Matplotlib.pyplot`
- Encode categorical data  
`ct = ColumnTransformer([("encoder", OneHotEncoder(), [0])])`
- Split dataset to training and test dataset
- We can see multiple independent variables
- Create regression model  
`regressor = linearRegression()`
- Fit the training set
- Test the model using test set
- Compare the actual value and predicted value.

## Code and output:

```
y_pred = model.predict(X_test)

[45] y_pred
Python

array([ 90555.15, 59516.62, 106544.7 , 64219.45, 68922.24, 123474.81,
       84511.78, 65278.67, 65159.99, 61397.74, 57863.7 , 50111. ])

[46] from sklearn.metrics import mean_absolute_error, mean_absolute_percentage_error, mean_squared_error
Python

[47] import numpy as np
Python

[48] mean_absolute_error(y_test,y_pred)
Python

4005.9263101681768
```

### Multiple Regression

```
house = pd.read_csv('https://github.com/yGIFoundation/Dataset/raw/main/Boston.csv')

[49] house.head()
Python

      CRIM    ZN  INDUS CHAS   NX   RM   AGE   DIS   RAD   TAX PTRATIO     B   LSTAT MEDV
0  0.00632  18.0     2.31  0  0.538  6.575  65.2  4.09000  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  0.02327  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  5.23  36.2
```

```
house.describe()

[50] house.describe()
Python

      CRIM    ZN  INDUS CHAS   NX   RM   AGE   DIS   RAD   TAX PTRATIO     B   LSTAT MEDV
count 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000
mean  3.617524 11.363696 11.16779  0.069170  0.554695  6.204634 68.574901 3.79049  9.549407 408.237154 18.45534 56.674032 12.85963 22.512866
std   8.601545 23.324253 6.660553 0.253994 0.115879 0.702617 28.148861 2.105710 8.707259 168.537116 2.164946 91.294864 7.141062 9.197104
min   0.006320 0.000000 0.460000 0.000000 0.385000 3.561000 2.900000 1.129600 1.000000 167.000000 12.600000 0.320000 1.730000 5.000000
25%  0.02045  0.000000 5.190000 0.449000 5.885500 45.025000 2.100175 4.000000 279.000000 174.000000 375.377500 6.950000 17.015000
50%  0.236510 0.000000 9.690000 0.000000 0.538000 6.208500 77.500000 3.207450 5.000000 390.000000 18.050000 391.440000 11.360000 21.200000
75%  3.677063 12.500000 18.100000 0.000000 0.624000 6.623500 94.075000 5.188425 24.000000 666.000000 20.200000 396.225000 16.955000 25.000000
max  88.976200 100.000000 27.740000 1.000000 0.871000 8.780000 100.000000 12.126500 24.000000 711.000000 22.000000 396.900000 37.970000 50.000000
```

```
house.columns

[51] house.columns
Python

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

y = house['MEDV']

[52] y = house['MEDV']
Python

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

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

[54] X_train.shape, X_test.shape, y_train.shape, y_test.shape
Python

((354, 13), (152, 13), (354,), (152,))

from sklearn.linear_model import LinearRegression
model = LinearRegression()

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

```
# Step 6 : train or fit model
model.fit(X_train,y_train)

...
# LinearRegression
LinearRegression()

Python
```

```
model.intercept_
34.21916368862993

Python
```

```
model.coef_
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,
      -5.45e-01])

Python
```

```
# Step 7 : predict model
y_pred = model.predict(X_test)

Python
```

```
y_pred
array([31.72, 22.02, 21.17, 39.78, 20.1, 22.06, 18.36, 14.79, 22.56,
       21.55, 18.38, 27.57, 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,
      20.57, 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, 21.85,
     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,
   31.66, 28.97, 32.46, 13.93, 35.49, 19.36, 19.6 , 1.44, 24.1 ,
   33.67, 20.62, 26.89, 19.29, 30.24, 29.74, 13.6 , 13.82, 19.76,
  21.52, 23.5 , 23.54, 23.54, 23.54, 23.54, 23.54, 23.54, 23.54,
  16.77, 25.41, 14.95, 3.72, 15.03, 10.91, 23.45, 30.45, 30.75,
  23.73, 22.19, 13.76, 18.47, 18.15, 36.6 , 27.49, 11. , 17.26,
  22.49, 16.53, 29.49, 22.89, 24.68, 20.38, 19.69, 22.95, 27.32,
  24.86, 20.2 , 29.14, 7.43, 5.85, 25.38, 38.75, 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, 18.15, 31.24, 20.85, 27.36, 30.53])
```

```
# Step 8 : model accuracy
from sklearn.metrics import mean_absolute_error, mean_absolute_percentage_error, mean_squared_error

Python
```

```
mean_absolute_error(y_test,y_pred)

Python
```

```
3.155030527602485
```

03/05/2024

## Logistic Regression

LAB - 4  
LOGISTIC REGRESSION

→ Import necessary libraries  
import pandas as pd  
import numpy as np  
from sklearn.linear\_model import LogisticRegression

→ Load the dataset  
iris = load\_iris()  
X = iris.data  
y = iris.target

→ Create model and train  
model = LogisticRegression()  
model.fit(X.drop(['Outcome']), y)

→ Predict values  
model.predict(X.drop(['Outcome']))

→ Load dataset  
df = pd.read\_csv('iris.csv')  
test = df[0:1]  
train = df[1:,:]

→ Create model and predict  
clf = KNeighborsClassifier(n\_neighbors=5, weights='uniform')  
clf.fit(x\_train, y\_train)  
clf.predict(test[1:-1])

Steps:  
→ Import  
→ Convert column  
→ Bind  
→ a set  
→ split  
→ Create  
→ Predict

Output:  
Accuracy

Steps:  
→ Load the dataset  
→ Initialize  
→ Visualize  
→ Plot

Output:  
3  
2  
1  
0

```

D def initialize_weights_and_bias(dimension):
    w = np.full((dimension,1),0.01)
    b = 0.0
    return w, b
[6] Python

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

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
[14] Python

def update(w, b, x_train, y_train, learning_rate,number_of_iteration):
    cost_list = []
    cost_list2 = []
    index = []
    # updating(learning) parameters is number_of_iteration times
    for i in range(number_of_iteration):
        # 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 %: %f" %(i, cost))
        # we update(learn) parameters weights and bias
        parameters = {"weight": w,"bias": b}
        plt.plot(index,cost_list)
        plt.plot(index,cost_list2)
        plt.xlabel("Number of iteration")
        plt.ylabel("Cost")
        plt.show()
    return parameters, gradients, cost_list
[14] Python

D def predict(w,b,x_test):
    # x_test is a input for forward propagation
    z = sigmoid(np.dot(w.T,x_test))
    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
[14] Python

!pip install mymodule
[14] Python

-- 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=2cdcfad015fdc631a533b7e720b40c29061c95d1a6a50b070d71cb04054f3f
    Stored in directory: /root/.cache/pip/wheels/f4/64/ed/4ac2e7514f3e85b8aea509dff4c236817b10c4db005e511d8
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):
    n = X.shape[0]
    A = sigmoid(np.dot(w.T, X) + b)
    cost = -1 / n * 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)
    cost = -1 / m * np.sum(y * np.log(A) + (1 - y) * np.log(1 - A))
    db = 1 / m * np.sum(A - y)
    dw = 1 / m * np.sum(X * (A - y).T)
    return dw, db
...
def logistic_regression(X_train, y_train, X_test, y_test, learning_rate, num_iterations):
    # Initialize
    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 = w - learning_rate * dw
        b = b - learning_rate * db

        # Record the costs
        if i % 100 == 0:
            cost = compute_cost(w, b, X_train)
            costs.append(cost)
            print("Cost after iteration {} : {}".format(i, cost))

    # Predictions
    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: {} %".format(train_accuracy))
    print("Test accuracy: {} %".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=1, num_iterations=100)

```

(a)

```

...
Cost after iteration 0: 0.6782740160052536
Train accuracy: 29.74534161490683 %
Test accuracy: 21.3953488372093 %

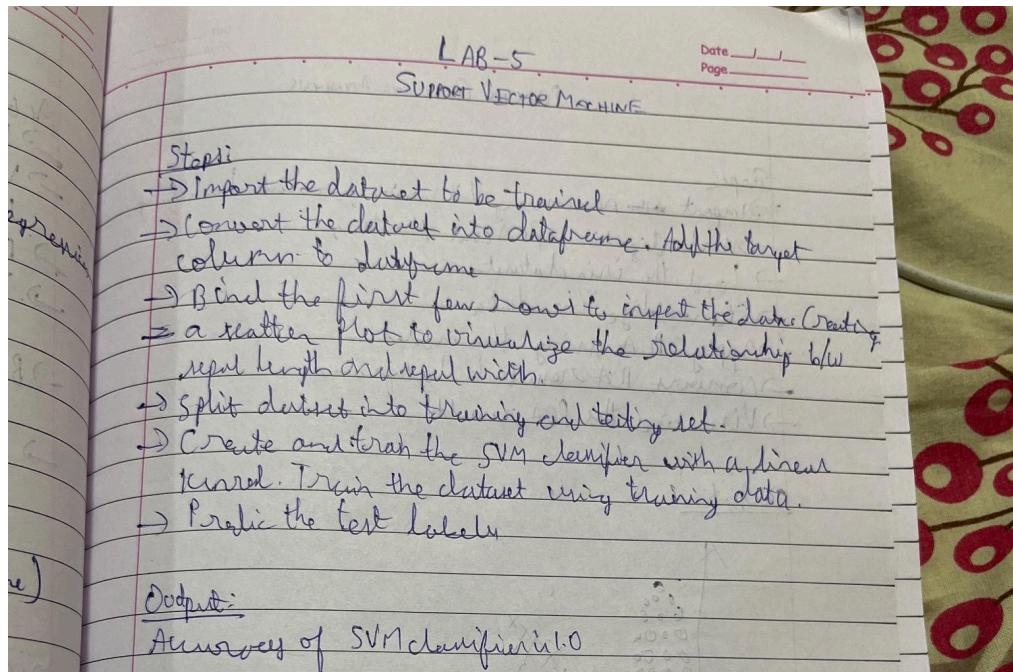
...
(array([[ 1.77806654e-02],
       [ 1.0168388e-02],
       [ 1.27806976e-01],
       [ 1.25749649e+00],
       [ 1.25931875e-05],
       [ 2.68863405e-04],
       [ 4.29020848e-04],
       [ 2.63106003e-04],
       [ 3.49357933e-05],
       [ 2.02145931e-05],
       [ 1.25698784e-03],
       [-3.98285024e-04],
       [ 8.96937014e-03],
       [ 2.02426962e-01],
       [ 3.00718647e-06],
       [ 4.19158446e-05],
       [ 6.03411729e-05],
       [ 2.00740406e-05],
       [-6.24383672e-06],
       [ 6.24944788e-07],
       [ 2.79506973e-02],
       [ 1.29326368e-02],
       [ 1.28774929e-01],
       [ 3.39189908e+00],
       [ 5.79135019e-05],
       ...
       [ 1.25862280e-03],
       [ 4.68695564e-04],
       [ 1.09671301e-04],
       [ 3.52490835e-05]]),
-1.5161875221606185)

```

24/05/2024

## Build Support vector machine model for a given dataset

Algorithm:



Code and Output:

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

```
df.describe().T
```

	count	mean	std	min	25%	50%	75%	max
radius_mean	560.0	14.127392	3.534069	659.0000	11.700000	13.370000	15.790000	28.11000
texture_mean	560.0	19.289469	4.301036	97.01000	16.170000	18.840000	21.800000	39.28000
perimeter_mean	560.0	91.969033	24.09981	437.90000	75.170000	68.340000	104.10000	198.50000
area_mean	560.0	654.889104	3519.4159	143.00000	420.30000	551.10000	782.70000	2510.00000
smoothness_mean	560.0	0.09636	0.014054	0.02539	0.086270	0.095870	0.105300	0.16340
compactness_mean	560.0	0.10341	0.052613	0.01959	0.064925	0.092300	0.130400	0.34540
concavity_mean	560.0	0.087959	0.09720	0.00000	0.02956	0.071940	0.130700	0.42860
concave points_mean	560.0	0.046919	0.03800	0.00000	0.02019	0.03590	0.074000	0.21120
symmetry_mean	560.0	0.91165	0.07414	0.10690	0.16900	0.175200	0.195700	0.30400
fractal_dimension_mean	560.0	0.05298	0.07060	0.046940	0.05700	0.061940	0.06120	0.09744
radius_se	560.0	0.05175	0.07731	0.111500	0.23400	0.38400	0.47900	2.97100
texture_se	560.0	0.11685	0.051948	0.09020	0.093000	0.100000	0.107000	0.20000
perimeter_se	560.0	0.06629	0.07025	0.02000	0.02500	0.03200	0.03700	0.09000
area_se	560.0	0.33709	45.42006	0.02000	17.83000	24.53000	45.92000	541.20000
smoothness_se	560.0	0.007041	0.00002	0.00711	0.00116	0.003930	0.00145	0.07113
compactness_se	560.0	0.035479	0.017908	0.00232	0.01390	0.02489	0.023445	0.12940
concavity_se	560.0	0.031894	0.030196	0.00000	0.01590	0.02930	0.02309	0.36900
concave points_se	560.0	0.011759	0.006170	0.00000	0.00738	0.01930	0.014710	0.05729
symmetry_se	560.0	0.030542	0.002656	0.007982	0.015160	0.016731	0.023400	0.07955
fractal_dimension_se	560.0	0.003755	0.002656	0.000085	0.00238	0.001357	0.004558	0.02384
radius_worst	560.0	15.26910	4.833842	7.92000	13.01000	14.97000	18.79000	36.04000
texture_worst	560.0	25.677223	6.140538	12.02000	21.08000	25.41000	27.72000	49.54000
perimeter_worst	560.0	107.261215	33.02552	50.41000	84.11000	97.66000	125.40000	251.20000
area_worst	560.0	880.981318	562.35960	185.20000	515.30000	684.50000	1084.00000	4254.00000
smoothness_worst	560.0	0.132380	0.022832	0.07170	0.11660	0.131300	0.146000	0.22290
compactness_worst	560.0	0.254265	0.15733	0.02739	0.147200	0.21900	0.39100	1.05800
concavity_worst	560.0	0.272188	0.20652	0.00000	0.14500	0.22870	0.38300	1.25200
concave points_worst	560.0	0.114000	0.095732	0.00000	0.066930	0.115400	0.21900	0.39100
symmetry_worst	560.0	0.280076	0.061867	0.159500	0.250400	0.282200	0.317900	0.66300
fractal_dimension_worst	560.0	0.083946	0.016061	0.059340	0.071460	0.080040	0.090080	0.20750

```
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 = 0.7)
relevant_features = cor_target[cor_target >= 0.7]

# 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')

# Consider the results
print(names)

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

```

Python

```

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=42, 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 traing 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:
    -----
    initialize_parameters()
        Initializes the weights and bias.
    gradient_descent(X, y)
        Updates the weights and bias using gradient descent.
    update_weights(X, y)
        Updates the weights and bias.
    fit(X, y)
        Fits the SVM to the data.
    predict()
        Predicts the labels for the given data.
    ...

    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.
        ...
    
```

```

n = x.shape[0]
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.
    """
    y_ = np.where(y <= 0, -1, 1)
    for i in range(len(X)):
        if y_[i] * (np.dot(X[i], self.w) + self.b) >= 1:
            dw += 2 * self.lambdas * self.w
            db += 0
        else:
            dw += 2 * self.lambdas * self.w - np.dot(X[i], y_[i])
            db += y_[i]
    self.update_parameters(dw, db)

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

    Parameters:
    -----------
    dw : numpy array
        The change in weights.
    db : float
        The change in bias.
    """
    self.w -= self.lr * dw
    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()
    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 outputs
    output = np.dot(X, self.w) + self.b
    # get the signs of the labels depending on if it's greater/less than zero
    labels = np.sign(output)
    # predictions set them to -1 if they are less than or equal to -1 else set them to 1
    predictions = np.where(labels <= -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 = SVC()
model.fit(X_train,y_train)
predictions = model.predict(X_test)
accuracy(y_test, predictions)

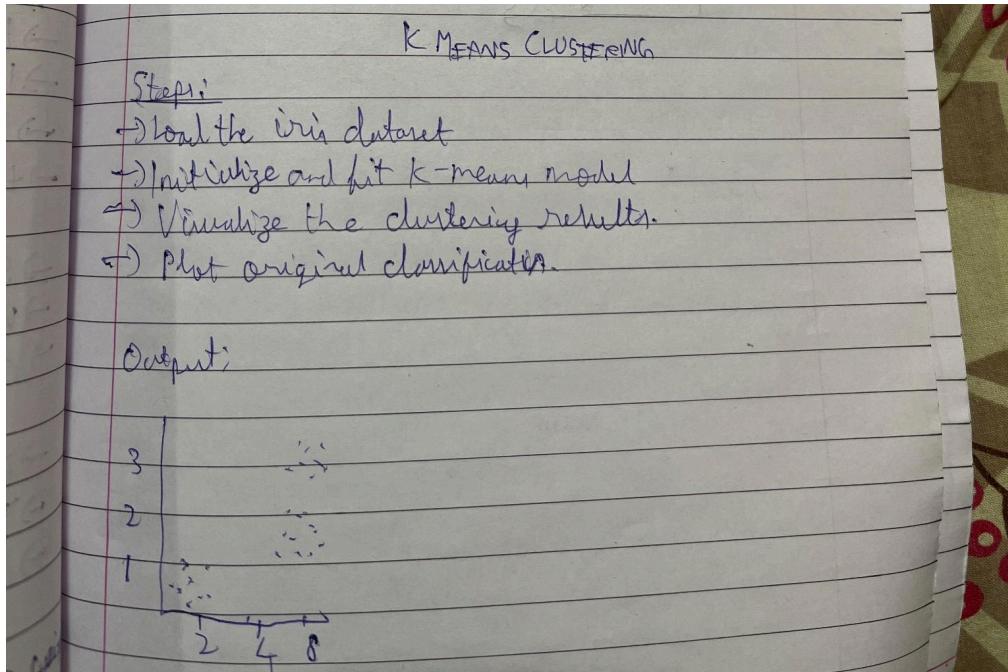
0.882880804957522

```

24/05/2024

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

Algorithm:



```
>Returns:  
numpy.ndarray: array containing the index of the centroid for each point  
  
X = np.expand_dims(X, axis=0) # 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):  
        centroids[i] = 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() # initialize the centroids  
    for i in range(iterations):  
        points = self.assign_points_to_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 max(points) < X.shape[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, 3000)
```

```
fig = go.Figure()
fig.add_trace(go.Scatter(
    x=[points == 0, 1], y=[points == 0, 1],
    mode='markers', marker_color="#00CED1", name='iris-setosa'
))

fig.add_trace(go.Scatter(
    x=[points == 1, 0], y=[points == 1, 0],
    mode='markers', marker_color="#D9E9F7", name='iris-versicolor'
))

fig.add_trace(go.Scatter(
    x=[points == 2, 1], y=[points == 2, 1],
    mode='markers', marker_color="#FF7F0E", name='iris-virginica'
))

fig.add_trace(go.Scatter(
    x=[centroids[0], 0], y=[centroids[1], 0],
    mode='markers', marker_color="#00CED1", marker_symbol="x", marker_size=15, name='centroids'
))

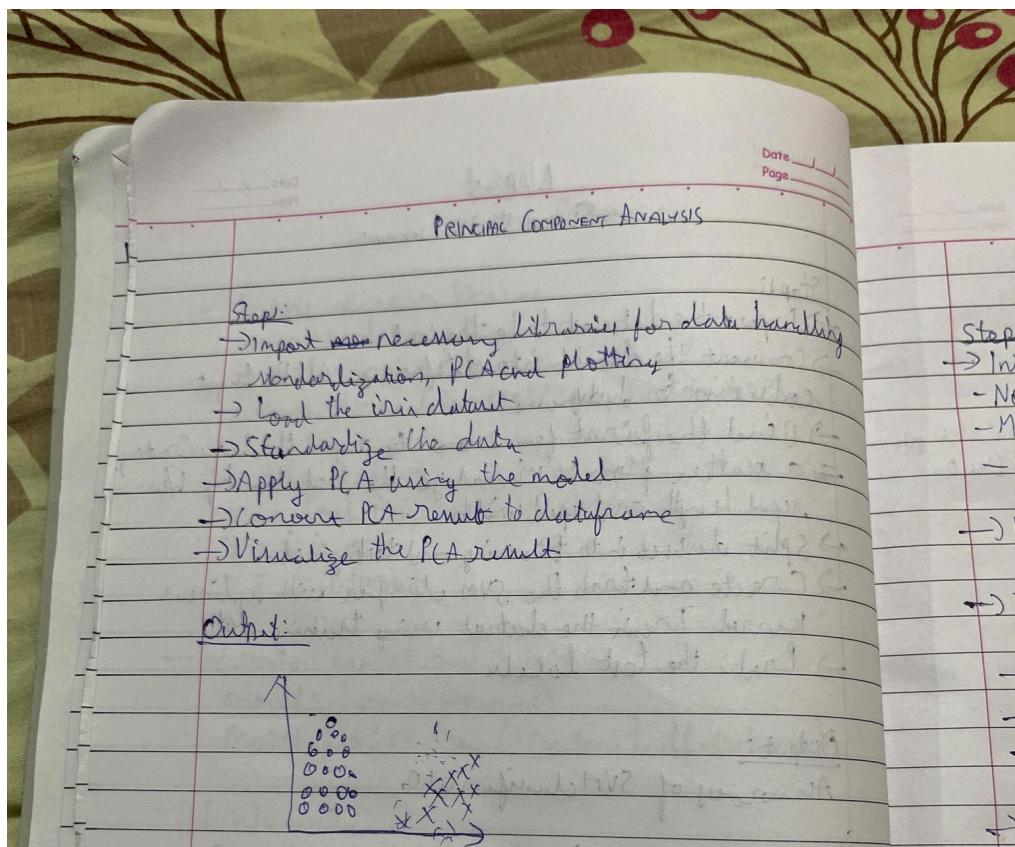
fig.update_layout(template='plotly_dark', width=3800, height=900)
```

Python

24/05/2024

## Implement Dimensionality reduction using Principle Component Analysis (PCA)

Algorithm:



```
# get the absolute value of the correlation
cor_target = abs(corr[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']

x = df[names].values
```

```
class PCA:
    """Principal component analysis (PCA) class for dimensionality reduction.

    Attributes:
        n_components (int): Number of principal components to retain.

    Methods:
        __init__(self, n_components):
            Constructor method that initializes the PCA object with the number of components to retain.
            Args:
                - n_components (int): Number of principal components to retain.
            self.n_components = n_components
        fit(self, X):
            Fits the PCA model to the input data and computes the principal components.
            Args:
                - X (numpy.ndarray): Input data matrix with shape (n_samples, n_features).
            # 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
```

```

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

# Compute the eigenvalues and eigenvectors of the covariance matrix.
eigenvalues, eigenvectors = np.linalg.eigh(cov)
# Reverse the order of the eigenvalues and eigenvectors.
eigenvalues = eigenvalues[-1:-0:-1]
eigenvectors = eigenvectors[-1:-0:-1]

# Keep only the first n_components eigenvalues as the principal components.
self.components_ = eigenvectors[1:n_components]

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

# Compute the variance explained by each principal component
self.explained_variances_ = eigenvalues[1: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.

    Parameters
    ----------
    X : numpy.ndarray
        Input data matrix with shape (n_samples, n_features).

    Returns
    -------
    transformed_data : numpy.ndarray
        Transformed data matrix with shape (n_samples, n_components).
    """
    # Center the input data around zero using the mean computed during the fit step.
    X -= np.mean(X, axis=0)

    # Project the centered input data onto the principal components.
    transformed_data = np.dot(X, self.components_)

    return transformed_data

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

    Parameters
    ----------
    X : numpy.ndarray
        Input data matrix with shape (n_samples, n_features).
    """
    self.fit(X)
    transformed_data = self.transform(X)
    return transformed_data

pca = PCA(2)
pca.fit(X)

X_transformed = pca.transform(X)

X_transformed[1].shape
...
(562,)

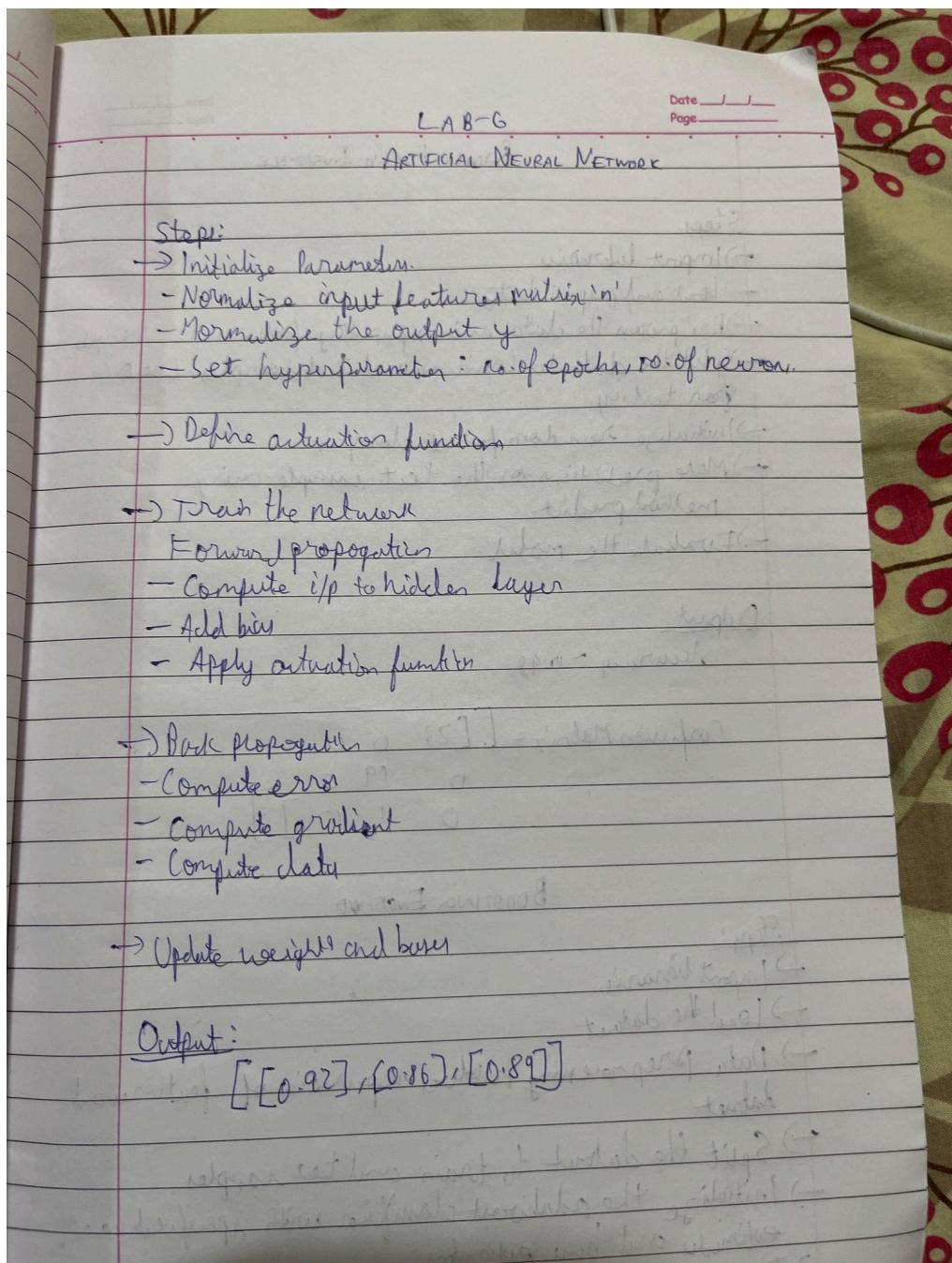
fig = plt.figure(figsize=(10, 6))
fig.suptitle("PCA transformed data for breast cancer dataset",
             fontsize=16)
ax0 = fig.add_subplot(111, title="PC1", aspect="equal")
ax1 = fig.add_subplot(111, title="PC2", aspect="equal")
for i in range(len(y0)):
    if y0[i] == 0:
        ax0.scatter(X_transformed[i][0], X_transformed[i][1], color="red")
    else:
        ax0.scatter(X_transformed[i][0], X_transformed[i][1], color="blue")
for i in range(len(y1)):
    if y1[i] == 0:
        ax1.scatter(X_transformed[i][0], X_transformed[i][1], color="red")
    else:
        ax1.scatter(X_transformed[i][0], X_transformed[i][1], color="blue")
fig.show()

```

31.05.2024

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

Algorithm:



```

✓ 0s  import numpy as np
    from sklearn.model_selection import train_test_split

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

    ↴ Database raw shape (86,7130)

✓ 0s  np.random.shuffle(db)
    y = db[:, 0]
    x = np.delete(db, [0], axis=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, 7129) (9, 7129)

✓ 0s  [4] 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))

✓ 0s  [5] 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

✓ 0s  [6] def activate_layer(layer, weights, x):
        for i in range(0, len(layer)):
            layer[i] = 1.7159 * np.tanh(2.0 * sum_function(weights, i, x) / 3.0)

✓ 0s  [7] def soft_max(layer):
        soft_max_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))
            soft_max_output_layer[i] = np.exp(layer[i] - np.max(layer)) / denominator
        return soft_max_output_layer

✓ 0s  [8] 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[j] * activation[i]) + weights[i][j]

✓ 0s  [9] 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[i][j]
            hidden_gradient[i] = sum_ * hidden_derivative[i]
        recalculate_weights(learning_rate, hidden_weights, output_gradient, hidden_layer)
        recalculate_weights(learning_rate, weights, hidden_gradient, x)

```

```
✓ [10] one_hot_encoding = np.zeros((2,2))
    for i in range(0, len(one_hot_encoding)):
        one_hot_encoding[i][i] = 1
    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 = soft_max(output_layer)
        training_correct_answers += 1 if y_train[i] == np.argmax(output_layer) else 0
    print("MLP Correct answers while learning: %s / %s (Accuracy = %s) on %s database." % (training_correct_answers, len(x_train),
                                                                                           training_correct_answers/len(x_train), "Duke breast cancer"))
```

→ MLP Correct answers while learning: 44 / 77 (Accuracy = 0.5714285714285714) on Duke breast cancer database.

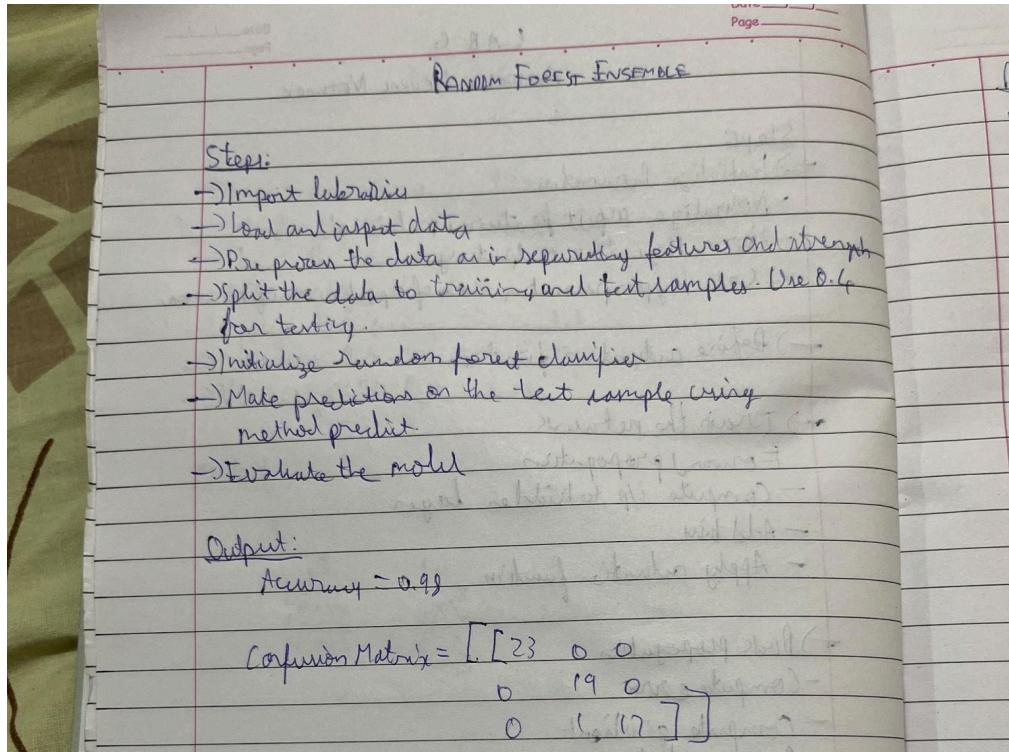
```
✓ [11] testing_correct_answers = 0
    for i in range(0, len(x_test)):
        activate_layer(hidden_layer, weights, x_test[i])
        activate_layer(output_layer, hidden_weights, hidden_layer)
        output_layer = soft_max(output_layer)
        testing_correct_answers += 1 if y_test[i] == np.argmax(output_layer) else 0
    print("MLP Correct answers while testing: %s / %s (Accuracy = %s) on %s database." % (testing_correct_answers, len(x_test),
                                                                                           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:



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

def train_test_split(X, y, random_state=42, test_size=0.2):
    """
    Splits the data into training and testing sets.

    Parameters:
        X (numpy.ndarray): Feature 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

X = iris.iloc[:, :-1].values
y = iris.iloc[:, -1].values.reshape(-1, 1)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

from sklearn.tree import DecisionTreeClassifier
n = DecisionTreeClassifier()
```

```

class RandomForest:
    """Random forest classifier.

    Parameters
    ----------
    n_estimators : int, default=5
        The number of trees in the random forest.
    max_depth : int, default=None
        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
        or each decision tree in the random forest.

    Attributes
    ----------
    n_estimators : 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
        or each decision tree in the random forest.
    trees : list of DecisionTreeClassifier
        The decision trees in the random forest.

    """
    def __init__(self, n_estimators=5, max_depth=None, min_samples=2):
        """Initialize the random forest classifier.

        Parameters
        ----------
        n_estimators : int, default=5
            The number of trees in the random forest.
        max_depth : int, default=None
            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
            or each decision tree in the random forest.
        trees : list of DecisionTreeClassifier
            The decision trees in the random forest.

    """
        self.n_estimators = n_estimators
        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
        -------
        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)
        # Create the trees.
        for _ in range(self.n_estimators):
            # 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_samples(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)
            # Append the tree in the list of trees.
            self.trees.append(tree)
        return self

    def bootstrap_samples(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_samples = dataset.shape[0]
        # Generate random indices to index into the dataset with replacement.
        np.random.seed()
        indices = np.random.choice(n_samples, n_samples, replace=True)
        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_occuring_value : int or float
            The most common label in the array.
        """
        y = list(y)
        # get the highest present class in the array
        most_occuring_value = max(y, key=y.count)
        return most_occuring_value

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

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

        Returns
        -------
        y : array-like of shape (n_samples,)
            Predicted classes for X.
        """
        # Predict the class for each tree.
        predictions = [tree.predict(X) for tree in self.trees]
        # Get the most common prediction for each sample.
        predictions = np.array(predictions).T
        most_common_predictions = np.apply_along_axis(lambda x: Counter(x).most_common(1)[0][0], 1, predictions)
        return most_common_predictions

```

```
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 samples from all trees for each sample in the test data
preds = np.stack(predictions, 0)
# iterate over each row in the preds and store it in the final predictions array
majority_predictions = np.array([self.most_common_label(pred) for pred in preds])
return majority_predictions
```

```
def accuracy(y_true, y_pred):
    """
    Compute 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()
y_train_encoded = label_encoder.fit_transform(y_train)
y_test_encoded = label_encoder.transform(y_test)
```

```
... /usr/local/lib/python3.6/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/python3.6/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().
```

```
from sklearn.preprocessing import LabelEncoder
label_encoder = LabelEncoder()
x_train_encoded = label_encoder.fit_transform(x_train)()
y_train_encoded = label_encoder.transform(y_train_encoded())
model = RandomForestClassifier(30, 30, 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)

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

# Calculate accuracy
accuracy(y_test_encoded, predictions)
```

```
... 0.9
```

31.05.2024

Implement Boosting ensemble method on a given dataset.

Algorithm:

Boosting Ensemble	
Step:	
→ Import libraries	
→ Load the dataset	
→ Data preprocessing involves separation of features and dataset	
→ Split the dataset to train and test samples	
→ Initialize the adaboost classifier with specified no. of estimator and base estimator.	
→ Train the model using the training dataset	
→ Make prediction for test sample using trained model	
→ Evaluate model	

```

[16] # Define AdaBoost class
class AdaBoost:

    def __init__(self):
        self.alphas = []
        self.G_M = []
        self.M = None
        self.training_errors = []
        self.prediction_errors = []

    def fit(self, X, y, M = 100):
        ...

        Fit model. Arguments:
        X: independent variables - array-like matrix
        y: target variable - array-like vector
        M: number of boosting rounds. Default is 100 - integer
        ...

        # Clear before calling
        self.alphas = []
        self.training_errors = []
        self.M = M

        # Iterate over M weak classifiers
        for m in range(0, M):

            # Set weights for current boosting iteration
            if m == 0:
                w_i = np.ones(len(y)) * 1 / len(y) # At m = 0, weights are all the same and equal to 1 / N
            else:
                # (d) Update w_i
                w_i = update_weights(w_i, alpha_m, y, y_pred)

            # (a) Fit weak classifier and predict labels
            G_m = DecisionTreeClassifier(max_depth = 1)      # Stump: Two terminal-node classification tree
            G_m.fit(X, y, sample_weight = w_i)
            y_pred = G_m.predict(X)

            self.G_M.append(G_m) # Save to list of weak classifiers

            # (b) Compute error
            error_m = compute_error(y, y_pred, w_i)
                w_i = update_weights(w_i, alpha_m, y, y_pred)

            # (a) Fit weak classifier and predict labels
            G_m = DecisionTreeClassifier(max_depth = 1)      # Stump: Two terminal-node classification tree
            G_m.fit(X, y, sample_weight = w_i)
            y_pred = G_m.predict(X)

            self.G_M.append(G_m) # Save to list of weak classifiers

            # (b) Compute error
            error_m = compute_error(y, y_pred, w_i)
            self.training_errors.append(error_m)

            # (c) Compute alpha
            alpha_m = compute_alpha(error_m)
            self.alphas.append(alpha_m)

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

    def predict(self, X):
        ...

        Predict using fitted model. Arguments:
        X: independent variables - array-like
        ...

        # Initialise dataframe with weak predictions for each observation
        weak_preds = pd.DataFrame(index = range(len(X)), columns = range(self.M))

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

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

        return y_pred

```

```
[47] import pandas as pd
     import numpy as np
     from sklearn.model_selection import train_test_split
     from sklearn.tree import DecisionTreeClassifier

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

# Column names
names = pd.read_csv('/content/spambase.names', sep = ':', skiprows=range(0, 33), header = None)
col_names = list(names[0])
col_names.append('Spam')

# Rename df columns
df.columns = col_names

# Convert classes in target variable to {-1, 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').values,
                                                    df['Spam'].values,
                                                    train_size = 3065,
                                                    random_state = 2)
```

```
# Fit model
ab = AdaBoost()
ab.fit(X_train, y_train, M = 400)

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

from sklearn.metrics import accuracy_score

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

```
Accuracy: 0.9440104166666666
```