

INDEX

Name of the Laboratory	WIT-JAVA AND PIG
Sub. Code	S022ABD
Name of the Staff In-Charge	R.D.SARASWATI, M.Tech

Sl No.	Date of Experiment	Name of the Experiment	Page No.	Date of Submission	Marks	Staff Signature with Date
7	14/08/23	Map Reduce using hadoop	124-126	01/09/23	10	EXPT 11
8	14/08/23	HDFS Connectivity and working with HDFS	27-28	01/09/23	10	EXPT 12
9a	20/08/23	Wordcount in Pig	29	07/09/23	10	EXPT 13
9b	20/08/23	Wordcount in Hive	30	07/09/23	10	EXPT 14
10a	20/08/23	YARN Connectivity	31-32	07/09/23	10	EXPT 15
10b	04/09/23	YARN Program	33-34	07/09/23	10	EXPT 16
11	04/09/23	Working with tableau	35-36	07/09/23	10	EXPT 17
		Explain	31-32	07/09/23	10	EXPT 18
		MLA for ML	33-34	07/09/23	10	EXPT 19
		MLA for ML	35-36	07/09/23	10	EXPT 20
		MLA for ML	37-38	07/09/23	10	EXPT 21
		MLA for ML	39-40	07/09/23	10	EXPT 22
		MLA for ML	41-42	07/09/23	10	EXPT 23
		MLA for ML	43-44	07/09/23	10	EXPT 24
		MLA for ML	45-46	07/09/23	10	EXPT 25
		MLA for ML	47-48	07/09/23	10	EXPT 26
		MLA for ML	49-50	07/09/23	10	EXPT 27
		MLA for ML	51-52	07/09/23	10	EXPT 28
		MLA for ML	53-54	07/09/23	10	EXPT 29
		MLA for ML	55-56	07/09/23	10	EXPT 30

Expt. No. 1 Page No. 01
 Expl. Name Data Processing-Building Good Datasets Date: 19/09/23

- Aim :
 To implement data preprocessing using python.
- Procedure :
1. The pandas library is used to load the dataset from the CSV file into a data frame.
 2. Any necessary data cleaning steps can be performed. Additionally missing values in income column are handled by filling them with mean value.
 3. The data set is split into features (x) and the target variable y.
 4. The data set is split into training and testing sets using the functions from the scikit learn module.
 5. Finally the shapes of the training and testing sets are printed to verify the sizes of the sets.

using this command job transformation or

most probably at last it is possible to do this by using

business as per requirement of prediction like

billions of numbers we can use

now more than most

probably

but still

it is

possible

to do

but still

it is

Program :-

```

import csv
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split

data = [{"Region": "Asia", "Age": 49, "Income": 86400, "no": 1,
          "Brazil": 32, 57600, "yes": 1},
        {"USA": 35, 64800, "no": 1,
         "India": 40, 69600, "yes": 1},
        {"Brazil": 43, 73200, "no": 1,
         "USA": 45, 0, "yes": 1},
        {"Brazil": 0, 62400, "no": 1,
         "India": 53, 94800, "yes": 1},
        {"USA": 55, 99600, "no": 1,
         "India": 42, 80000, "yes": 1}]

f = open("sample.csv", "w")
writer = csv.writer(f)
writer.writerow(["Region", "Age", "Income", "Online Shopping"])
writer.writerows(data)
writer.close()

data = pd.read_csv("sample.csv")
array = np.array(data)

def fill_empty(array, col):
    for i, j in enumerate(array):
        if i == 0:
            array[i][col] = "Asia"
        elif i == 1:
            array[i][col] = "Brazil"
        elif i == 2:
            array[i][col] = "USA"
        elif i == 3:
            array[i][col] = "India"
        else:
            array[i][col] = "Brazil"

fill_empty(array, 0)
array

```

```

margin = 10
margin_Cat = 10
margin_Dog = 10
margin_Human = 10
margin_Leopard = 10
margin_Lion = 10
margin_Sheep = 10
margin_Tiger = 10
margin_Wolf = 10
margin_Zebra = 10

def test():
    print("Margin for Cat is ", margin_Cat)
    print("Margin for Dog is ", margin_Dog)
    print("Margin for Human is ", margin_Human)
    print("Margin for Leopard is ", margin_Leopard)
    print("Margin for Lion is ", margin_Lion)
    print("Margin for Sheep is ", margin_Sheep)
    print("Margin for Tiger is ", margin_Tiger)
    print("Margin for Wolf is ", margin_Wolf)
    print("Margin for Zebra is ", margin_Zebra)

test()

```

```

if j[col] == 0:
    array[i][col] = "NAN"
else:
    tot, count = 0, 0
    for i in array[:, col]:
        if i == "NAN":
            tot += i
        else:
            mean = tot / count
            array[i][col] = mean
            count += 1
    print("Mean for column ", col, " is ", mean)

def NANTOMean(array, col):
    for i, j in enumerate(array):
        if j[col] == "NAN":
            array[i][col] = mean
    print("Mean for column ", col, " is ", mean)

def PrintArray():
    print("Cat", "Dog", "Human", "Leopard", "Lion", "Sheep", "Tiger", "Wolf", "Zebra")
    print("0", "0", "0", "0", "0", "0", "0", "0", "0")
    print("1", "1", "1", "1", "1", "1", "1", "1", "1")
    print("2", "2", "2", "2", "2", "2", "2", "2", "2")
    print("3", "3", "3", "3", "3", "3", "3", "3", "3")
    print("4", "4", "4", "4", "4", "4", "4", "4", "4")
    print("5", "5", "5", "5", "5", "5", "5", "5", "5")
    print("6", "6", "6", "6", "6", "6", "6", "6", "6")
    print("7", "7", "7", "7", "7", "7", "7", "7", "7")
    print("8", "8", "8", "8", "8", "8", "8", "8", "8")
    print("9", "9", "9", "9", "9", "9", "9", "9", "9")

PrintArray()

fillEmptyArray(1)
NANTOMeanArray(1)
print("Cat", "Dog", "Human", "Leopard", "Lion", "Sheep", "Tiger", "Wolf", "Zebra")
print("0", "0", "0", "0", "0", "0", "0", "0", "0")
print("1", "1", "1", "1", "1", "1", "1", "1", "1")
print("2", "2", "2", "2", "2", "2", "2", "2", "2")
print("3", "3", "3", "3", "3", "3", "3", "3", "3")
print("4", "4", "4", "4", "4", "4", "4", "4", "4")
print("5", "5", "5", "5", "5", "5", "5", "5", "5")
print("6", "6", "6", "6", "6", "6", "6", "6", "6")
print("7", "7", "7", "7", "7", "7", "7", "7", "7")
print("8", "8", "8", "8", "8", "8", "8", "8", "8")
print("9", "9", "9", "9", "9", "9", "9", "9", "9")

fillEmptyArray(2)
NANTOMeanArray(2)
print("Cat", "Dog", "Human", "Leopard", "Lion", "Sheep", "Tiger", "Wolf", "Zebra")
print("0", "0", "0", "0", "0", "0", "0", "0", "0")
print("1", "1", "1", "1", "1", "1", "1", "1", "1")
print("2", "2", "2", "2", "2", "2", "2", "2", "2")
print("3", "3", "3", "3", "3", "3", "3", "3", "3")
print("4", "4", "4", "4", "4", "4", "4", "4", "4")
print("5", "5", "5", "5", "5", "5", "5", "5", "5")
print("6", "6", "6", "6", "6", "6", "6", "6", "6")
print("7", "7", "7", "7", "7", "7", "7", "7", "7")
print("8", "8", "8", "8", "8", "8", "8", "8", "8")
print("9", "9", "9", "9", "9", "9", "9", "9", "9")

PrintArray()

x = array[:, :-1]
y = array[:, -1]

X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.8, random_state=1)

print("X shape: ", X.shape)
print("Y shape: ", Y.shape)

```

Output :-

```
[[ 'India' , 49 , 86400 , 'no' ] , [ 'Brazil' , 32 , 57600 , 'yes' ] , [ 'USA' , 35 , 64800 , 'no' ] , [ 'India' , 43 , 73200 , 'no' ] , [ 'USA' , 45 , 0 , 'yes' ] , [ 'India' , 40 , 67600 , 'no' ] , [ 'Brazil' , 0 , 62400 , 'yes' ] , [ 'India' , 53 , 94800 , 'no' ] , [ 'USA' , 55 , 99600 , 'no' ] , [ 'India' , 52 , 88800 , 'yes' ] , [ 'India' , 49 , 86400 , 'no' ] , [ 'Brazil' , 32 , 57600 , 'yes' ] , [ 'USA' , 35 , 64800 , 'no' ] , [ 'Brazil' , 43 , 73200 , 'no' ] , [ 'USA' , 45 , 0 , 'yes' ] , [ 'India' , 40 , 67600 , 'yes' ] , [ 'Brazil' , 0 , 62400 , 'no' ] , [ 'India' , 53 , 94800 , 'yes' ] , [ 'USA' , 55 , 99600 , 'no' ] , [ 'India' , 42 , 80000 , 'yes' ] ]
```

X train shape : (21, 3)

X test shape : (8, 3)

Input :-
1) - 49 - Possess = 1
2) - 32 - Possess = 1
3) - 35 - Possess = 1
4) - 43 - Possess = 1
5) - 45 - Possess = 1
6) - 40 - Possess = 1
7) - 53 - Possess = 1
8) - 55 - Possess = 1
9) - 42 - Possess = 1

Result :-
The data preprocessing for building good training sets is executed successfully.

Expt. Name _____

Date: _____

Output:

	Name	Twitter_Username	Account_Start_Dtm
0	A. Donald McEachin	Ruf_McEachin	2017-01-03T00:00:00.000Z
1	Aaron Richley	Ruf_Richley	2020-06-24T00:00:00.000Z
2	Aaron Perskin	Aaron_perkins	2010-11-15T00:00:00.000Z
3	Aaron Refra	Aaronperna	2007-10-31T00:00:00.000Z
4	Aaron Shick	Aaronshock	2009-03-12T00:00:00.000Z

Expt. No. 2
Expt. Name Manipulation of Twitter Dataset

Page No. 05
Date: 26/10/23

Algorithm:
1. Import the required modules, dataset and load the dataset.
2. Then drop a certain column from the dataset.
3. Drop a row from the dataset.
4. Then rename a column from the dataset.
5. Then group a data under a certain feature.
6. Slice the dataset to get required information.

7. Describe the dataset and print it.

Age	Instagram_username	Political_Party
59.0	rumpreakkin	Democratic party
42.0	NAN	Democratic party
56.0	aperstin52	Democratic party
61.0	NAN	Republican party
39.0	aaronshock	Republican party

Final output by printing the dataset in Jupyter notebook.

1. Name Twitter-Username Account-Start-time
 0 A. Donald McCalhin Repretachin 2017-01-03 00:00:00
 1 Aaron Michalutty Reprichelutty 2010-06-27 00:00:00
 2 Aaron Pekkin Aaronpekin 2010-11-13 00:00:00
 3 Aaron pena Aaron pena 2007-10-31 00:00:00
 4 Aaron Shook Aaron Shook 2009-03-12 00:00:00

~~df-dropped = data.drop ("Birthplace", axis=1)~~
~~df-dropped.head(1)~~

~~df : df.drop("Birthplace", axis=1)~~

Age Instagram_Username Political_Party
 59.0 rumpcalhin Democatic party
 42.0 (nawjwudoh) Democratic party
 56.0 apetism5 (Democratic party)
 61.0 NAW Republican party
 39.0 aaronshak Republican party

~~df : df.drop("Political_Party", axis=1)~~

Political_Party

Sex Instagram_Username Political_Party
 female rumpcalhin Repretachin
 male (nawjwudoh) Democratic party
 transgender female aaronshak Republican party
 transgender male

~~df : df.drop("Political_Party", axis=1)~~

Expt. Name _____

```

0    ruprichalvin
1    NaN
2    apeshin52
3    NAN
4    aaronshock
      .
  
```

Name: Triogram username, dtype : object

```

count    2514
unique   1506
top      Ben Carson
freq     12
  
```

Name: name, dtype: object.

	Name	Twitter_Username	Accont_Start_time
0	Aaron richelvitz	RupeFabrin	2017-01-03T00:00:00Z
1	Aaron richelvitz	Ruprichalvin	2010-06-07T00:00:00Z
2	Aaron Rich	AaronRich	2007-10-31T00:00:00Z
3	Aaron Schatz	aaronshock	2009-03-12T00:00:00Z
4	AbbyFrikkenacci	AbbyFrikken	2013-03-18T00:00:00Z

Amount_ID	Sex	Birthplace	Birthday
816161091673464448	male	Germany	1961-10-10T00:00:00Z
160246973	male	USA	1978-01-01T00:00:00Z
964332	male	USA	1989-06-08T00:00:00Z
2395197	male	USA	1981-06-08T00:00:00Z
1262017122	female	USA	1988-12-21T00:00:00Z

Age	Instagram - Username	Political - Party
59.0	supreachin	Democratic party
42.0	NAN	Democratic party
61.0	NAN	Republican party
39.0	aaronshock	Republican party
31.0	abby4jawn	Democratic party

Name	Twitter - Username	Account - Start time
0 A Donald McEachin	RepEachin1950	2017-01-03T00:00:00Z
1 Aaron Nichols	RepNichols1973	2010-06-07T00:00:00Z
2 Aaron Peskin	Aaronpeskin	2010-11-13T00:00:00Z
3 Aaron Perna	Aaronperna	2007-10-31T00:00:00Z
4 Aaron School	AaronSchool	2009-03-18T00:00:00Z

ID	Sex	Birthplace	Birthday
016101001673448448	male	Germany	1961-10-10T00:00:00Z
160a46973	male	USA	1978-01-05T00:00:00Z
015369273	male	USA	1964-06-17T00:00:00Z
98432352	male	USA	1959-06-08T00:00:00Z
23951197	male	USA	1981-05-25T00:00:00Z

Age	Instagram - Username	Political - Party
59.0	supreachin	Democratic party
42.0	NAN	Democratic party
56.0	peskin52	Democratic party
61.0	NAN	Republican party
39.0	aaronshock	Republican party

data.columns
~~df_renamed = df.rename(columns = {i: mount_ID[i] for i in df})~~
~~df_renamed.head()~~

Result: ~~df.to_csv('dataset.csv')~~

Thus the Twitter dataset has been manipulated and result is successfully executed

Of Machine learning Algorithm

Aim:-
 To evaluate the performance of results of machine learning algorithm

Procedure:-

1. Import the required modules, dataset and load the dataset.
2. Import the necessary libraries and model.
3. Import the model from Sklearn
4. The dataset is split into training and testing sets using the train-test-split() function from Sklearn
5. Use the predict() method from the model to examine our test dataset.
6. In this program, 75% of the data is allocated to the training set and 25% to the testing set.
7. Finally, the performance metrics are printed.

Program:-

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, precision_score,
recall_score, f1_score
dataset = pd.read_csv('diabetes.csv')
```

import pandas as pd
 import numpy as np
 import matplotlib.pyplot as plt
 from sklearn.model_selection import train_test_split
 from sklearn.preprocessing import StandardScaler
 from sklearn.linear_model import LogisticRegression
 from sklearn.metrics import accuracy_score, precision_score,
 recall_score, f1_score
 dataset = pd.read_csv('diabetes.csv')

~~x = dataset. iloc [:, [4, 7]] values
 $y = \text{dataset. iloc [:, 8]. values}$~~

~~from sklearn.model_selection import train_test_split~~

~~x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.25)~~

~~from sklearn.linear_model import LogisticRegression~~

~~logisticRegr = LogisticRegression()~~

~~logisticRegr.fit(x_train, y_train)~~

~~Y_pred = logisticRegr.predict(x_test)~~

~~accuracy = accuracy_score(y_test, Y_pred)~~

~~precision = precision_score(y_test, Y_pred)~~

~~recall = recall_score(y_test, Y_pred)~~

~~f1 = f1_score(y_test, Y_pred)~~

~~print("Accuracy: ", accuracy)~~

~~print("Precision: ", precision)~~

~~print("Recall: ", recall)~~

~~print("F1 Score: ", f1)~~

~~Result: 0.7052~~

~~The evaluation of results of machine algorithm with diabetes dataset is executed successfully.~~

variables like blood sugar levels, number of pregnancies, age, etc. are present. When we consider all the variables, we can see that there is a positive correlation between blood sugar levels and blood pressure. We can also see that the number of pregnancies has a negative correlation with blood pressure. So, we can say that the number of pregnancies is a strong predictor of blood pressure. Now, we can split our data into training and testing sets. We can use the random forest classifier for this task. The code for this is as follows:

```
[1.35876151  0.404946103]
[0.81313372  0.20914106]
[-0.43506306 -0.80678687]
[-0.93010296  0.47622887]
[0.57101202  0.404946103]
[0.07171249  0.17757008]
[0.3546225  -0.57438803]
[-0.97010296 -0.67204882]
[-0.3546225  -0.28040867]
[0.57101202  0.47622887]
[0.07171249  0.17757008]
[0.3546225  -0.57438803]
[-0.97010296 -0.67204882]
[-0.3546225  -0.28040867]
```

Now, we can see that the data is balanced. We can use the random forest classifier for this task. The code for this is as follows:

```
accuracy = accuracy_score(y_test, Y_pred)
precision = precision_score(y_test, Y_pred)
recall = recall_score(y_test, Y_pred)
f1 = f1_score(y_test, Y_pred)

print("Accuracy: ", accuracy)
print("Precision: ", precision)
print("Recall: ", recall)
print("F1 Score: ", f1)
```

Result: 0.7052

The evaluation of results of machine algorithm with diabetes dataset is executed successfully.

Output :-

Salinity	T_deg_C	RTemp	Depth_m	Wt_m	Sta_ID	SL_out	Cat_out
0.20415626414174355	19.403CR	19.403CR	060	0930	05400560	05400560	0
1.0	19.403CR	19.403CR	060	0930	05400560	05400560	1
2.0	19.403CR	19.403CR	060	0930	05400560	05400560	2
3.0	19.403CR	19.403CR	060	0930	05400560	05400560	3
4.0	19.403CR	19.403CR	060	0930	05400560	05400560	4
..	19.403CR	19.403CR	060	0930	05400560	05400560	..
864858	20-1611SR	-MX-310-2839-09340264	-00004-	0934	0934	0934	0934
864859	20-1611SR	-MX-310-2839-09340264	-00004-	0934	0934	0934	0934
864860	20-1611SR	-MX-310-2839-09340264	-00004-	0934	0934	0934	0934
864861	20-1611SR	-MX-310-2839-09340264	-00004-	0934	0934	0934	0934
864862	20-1611SR	-MX-310-2839-09340264	-00004-	0934	0934	0934	0934

Expt. No. — 4.Cav
Expt. Name — Predicting Water Temperature based on Salinity Using Regression.

Salinity Using Regression

Aim :-

1. Import the libraries, dataset and load the dataset.
2. Read the dataset.
3. Remove the unwanted data from the dataset.
4. Splitting the data into training and testing data.
5. Dropping any rows with NaN values.
6. Data scatter of predicted values Exploring our results.

Program :-

```

import numpy as np
import pandas as pd
import Seaborn as sns
import matplotlib.pyplot as plt
from sklearn import preprocessing, svm
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
df = pd.read_csv('bottle.csv')
dt_binary = df[['Salinity', 'T_degC']]
dt_binary.columns = ['Sal', 'Temp']
dt_binary.head()

```

निम्नलिखित विषयों परिचय

प्रारूप बिल्डर से बुना प्रतिक्रिया करने वाली एक अभियान विस्तृत विस्तृत

प्रारूप विस्तृत

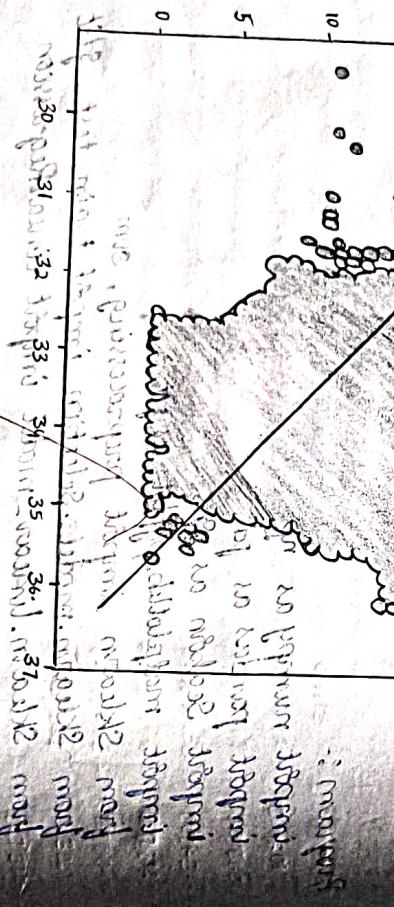
~~df = binary_fillna(method='ffill', inplace=True)~~

~~classf = np.array(df['binary'].values).reshape(-1, 1)~~

~~preds = np.array(df['binary'].values).reshape(-1, 1)~~

~~preds = np.array(df['binary'].values).reshape(-1, 1)~~

~~preds = np.array(df['binary'].values).reshape(-1, 1)~~



[Temp :- 15.0 17.5 19.0 20.5 22.0 23.5 25.0 26.5 28.0 29.5]
[Salinity :- 10.0 12.5 15.0 17.5 20.0 22.5 25.0 27.5 30.0 32.5]

[Temp :- 15.0 17.5 19.0 20.5 22.0 23.5 25.0 26.5 28.0 29.5]
[Salinity :- 10.0 12.5 15.0 17.5 20.0 22.5 25.0 27.5 30.0 32.5]

Thus ~~the water temperature can be predicted on Salinity~~
~~using Regression and the result is successfully predicted.~~

(1) ~~first~~ - 10

Output:
 Correlation matrix for Iris dataset is as follows:

	Sepal length(cm)	Sepal width(cm)	Petal length(cm)	Petal width(cm)
Sepal length(cm)	1.000000	-0.11570550	-0.02871739	
Sepal width(cm)	0.11757050	1.00000000	-0.428440	
Petal length(cm)	-0.02871739	-0.428440	1.000000	
Petal width(cm)			(0.817945) \rightarrow 0.846126 (second column after 285)	
Petal width(cm)				petal width(cm)

Sepal length(cm) \rightarrow 0.817945 (row 1 - 4, column 1, second X)

Sepal width(cm) \rightarrow -0.11570550 (row 2 - 4, column 2, second X)

Petal length(cm) \rightarrow -0.02871739 (row 3 - 4, column 3, second X)

Petal width(cm) \rightarrow (0.817945) \rightarrow 0.846126 (second column after 285)

Correlation matrix for Boston Housing dataset is as follows:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PRATIO	B	STAT
CRIM	1.000000	-0.200469	0.406583	-0.058842	0.420122								
ZN	-0.200469	1.000000	-0.533888	-0.062038	-0.516604								
INDUS	0.406583	-0.533888	1.000000	-0.042897	-0.509920								
CHAS	-0.058842	-0.062038	-0.042897	1.000000	-0.509920								
NOX	0.420122	-0.516604	-0.509920	-0.509920	1.000000								
RM		-0.219447	0.311991	-0.391676	0.091251	-0.30215							
AGE		0.352734	-0.569537	-0.644779	0.086518	-0.73407							
DIS			-0.379670	0.664408	-0.708027	-0.099176	-0.769250						
RAD			0.625505	-0.311948	0.595129	-0.007368	0.611441						
TAX			0.582264	-0.314563	0.720760	-0.035587	0.668023						
PRATIO			0.289946	-0.391679	0.383248	-0.121515	0.08933						
B			-0.385064	0.175520	-0.356977	-0.048788	0.50238005						
STAT			0.455621	-0.4122995	-0.6038000	-0.0539897	0.50080001	1.0000000	0.0000000	0.0000000	0.0000000	0.0000000	0.0000000

- Print the correlation matrix for the Iris dataset.
- Load the Boston Housing dataset using function from Sklearn datasets.
- Convert the Iris dataset into a pandas DataFrame.
- Calculate the correlation matrix for the Iris dataset.
- Print the correlation matrix for the Iris dataset.
- Load the Boston Housing dataset using function from Sklearn datasets.
- Convert the Boston Housing dataset to a pandas DataFrame.
- Calculate the correlation matrix for the Boston Housing dataset.
- Print the correlation matrix for the Boston Housing dataset.
- Calculate the pairwise correlation between the selected features in the Boston housing dataset using the `corr()` function on the respective columns.
- Print the calculated correlation value.

RH	AGE	DIS	RAD	CPTM
-0.219247	0.352734	-0.379670	0.625505	
0.311991	-0.369537	0.6649063	-0.2311946	
-0.3916765	0.644777	-0.700027	0.595129	INDIUS
0.091251	0.066518	-0.091176	0.001266	CHASE
-0.0802168	0.731470	0.769230	0.611441	DIABETES
1.000000	-0.240245	0.205246	-0.209847	RM
-0.240245	1.000000	-0.747861	0.456022	AGE
0.205246	-0.747861	1.000000	-0.494588	DIS
-0.209847	0.456022	-0.494588	1.000000	RAD
-0.298048	0.506456	-0.524432	0.910220	TAX
-0.355501	0.721515	-0.232471	0.464761	PRIMRO
0.128069	-0.273524	0.291512	-0.444413	BINR
-0.613808	0.602339	0.496796	-0.408676	LSTAT
TAX	PRIMRO	BINR	LSTAT	LSIM
CPTM	0.582764	0.265946	-0.365064	0.455010
ZN	-0.314563	-0.316779	0.175520	-0.412995
INDUS	0.720160	0.363248	-0.356977	0.608800
CHAS	-0.035587	-0.121515	0.048788	-0.058828
NOX	0.668023	0.188933	-0.280051	0.590899
RH	-0.395048	-0.355501	0.126069	-0.613908
AGE	0.506456	0.261515	-0.273524	0.602339
DIS	-0.534423	0.232471	0.291512	0.496796
RAD	0.910220	0.464761	-0.444413	0.408676
TAX	1.000000	0.460853	-0.444808	0.545995
PRIMRO	0.406456	1.000000	-0.171783	0.374049
E	-0.441100	-0.171783	1.000000	-0.366087
LSTAT	0.543995	-0.374049	-0.366087	1.000000

Program:-

import numpy as np
import pandas as pd

from sklearn.datasets import load_iris, load_boston

iris = load_iris()

iris.data = iris.data

iris.feature_names = iris.feature_names

iris_df = pd.DataFrame(iris.data, columns=iris.feature_names)

iris_corr_matrix = iris_df.corr()

print("Correlation matrix for iris dataset: ")

print(iris_corr_matrix)

boston = load_boston()

boston.data = boston.data

boston.feature_names = boston.feature_names

boston_df = pd.DataFrame(boston.data, columns=boston.feature_names)

boston_corr_matrix = boston_df.corr()

print("Correlation matrix for Boston Housing dataset: ")

print(iris_corr_matrix)

Print Boston - corr-matrix

Print()

Corr. val = boston_dt['CRIM'].corr(boston_dt['TAX'])
Housing dataset : "Crim-Corr"

() house - boat = house - boat
stob - house = stob - house

House - Correlation = house - house - house

House - value - house - house = house - house - house

House - tax - house = house - tax - house

() house - tax = house - house - tax - tax

() house

() house - boat = house - boat

stob - house = stob - house

House - house = house - house - house

House - value = house - value - house

Result is

Thus the correlation matrix of both iris and boston housing datasets are calculated using the correlation function.

() house - boat = house - boat - house
House - boat - house = house - house - house

(Random - 100 - record) Using1) import

Aim:-
To write a classification algorithm in Support vector machine
and execute it.

Procedure:-

1. Import the libraries
2. We read csv() function of pandas library which is used to read a csv file and perform various operations on it.
3. Use iloc[] method of pandas library, used to extract the required rows and columns from the dataset.
4. To create the best line or decision boundary that can segregate n dimensional space into classes so that we easily put the new data point in the correct category in the future.
5. The SVM algorithm helps to find the best line & decision boundary, this best boundary or region is called as a hyperplane.
6. The dimensions of hyperplane dependent on the features present in the dataset which means if they are 2 features then hyperplane will be a straight line.
7. If there are 3 features, then hyperplane will be a 2 dimensional plane.
8. To create the SVM classifier, import SVC class from sklearn.svm library.
9. Predict the output for test set. for this, will create a new vector y-pred.
10. From the dataset, the SVM classifier has divided the users into two regions (purchased & not purchased)
11. Users who purchased the SUV are in red region with red scattered points. Find who did not purchase the SUV are in the green region with green scatter points.

Expt. No. _____
 Expt. Name _____
 Date: _____

which shows separate two multiple classifications & show of
 linear basis despite no multiple classifications.

Q. The hyperplane has divided the two classes into
 purchased and not purchased variable.

Program:

```
import matplotlib.pyplot as plt
import pandas as pd
import numpy as np
df = pd.read_csv('social-network-ad.csv')
x = df.iloc[:, [2, 3]].values
y = df.iloc[:, 4].values
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, ytest = train_test_split(x, y, test_size=0.25,
random_state=0)
```

from sklearn.preprocessing import StandardScaler

sc = StandardScaler()

```
x_train = sc.fit_transform(x_train)
x_test = sc.transform(x_test)
from sklearn.svm import SVC
classifier = SVC(kernel='rbf', random_state=0)
classifier.fit(x_train, y_train)
y_pred = classifier.predict(x_test)
```

Result will be shown with different no. of features like mean
 no. of hours spent by user in buying company and their
 no. of purchases made by user like hours spent in buying
 hours above user's purchase made etc.

and which job will be given work depends on the
whether building job has been building

from sklearn.metrics import confusion_matrix, accuracy_score
cm = confusion_matrix(y_true, y_pred)

print(cm)

```
[[ 10  10  10]
 [ 10  10  10]
 [ 10  10  10]]
```

accuracy_score(y_true, y_pred)

x_set, y_set = x_test, y_test

x1, x2 = np.meshgrid(np.arange(x_set[:, 0].min() - 1,

stop = x_set[:, 0].max() + 1, step = 0.01),

y1, y2 = np.arange(x_set[:, 1].min() - 1,

stop = x_set[:, 1].max() + 1, step = 0.01))

if i == 0:
 y_set = np.argmax(y_set, axis=1)
else:
 y_set = np.argmax(y_set, axis=1) + 1

for i in range(0, 3):
 plt.contourf(x1, x2, classifier.predict(np.array([x1.ravel(), y1.ravel()]).T).reshape(2, 2), alpha=0.75, cmap='viridis')

plt.xlim(x_set[:, 0].min(), x_set[:, 0].max()),

plt.ylim(x_set[:, 1].min(), x_set[:, 1].max())

for i, j in enumerate(np.unique(y_set)):
 plt.scatter(x_set[y_set == i, 0], x_set[y_set == i, 1], c = ListedColormap(['pink', 'green'])(i), label=j)

plt.title('SVM (Test set)')

plt.xlabel('Age')

plt.ylabel('Estimated Salary')

for i in range(0, 3):
 plt.legend()
 plt.show()

(0: idle, 1: working, 2: working)
 (idle, working, working) = 0, 1, 2
 (idle, working, working) = 0, 1, 2

Result: Thus, the classification algorithm of SVM is executed successfully.

Output :
 Sodium 0.00065 ion density measure near
 occurs Chai = 0.00065
 emp local 0.00065

~~Suitable to School~~

(b) (c) 0.00065

~~for: if~~
 Acknowledged true
 insertedId "objectID" "64db49843810075b847c1"
 1.0 0.00065

~~{
 "name": "Sangeetha",
 "age": 30,
 "gra": 3.2}~~



~~name: Sangeetha
 age: 30
 gra: 3.2~~

Program:-

test > Show db Students.distinct("name") by
 test > Use School
 School > db.createCollection("Students")
 School > db.Students.insertOne({name: "Sangeetha", age: 30, gra: 3.2})
 School > db.Students.find()

Algorithm :-
 To Create, update and delete No SQL documents using MongoDB

1. Install MongoDB (Server from official link) & Run
2. Install .msi file.
3. Install mongoDB shell
4. Set it in c:
5. Set the path
6. Show database using mongoDB command : "use"
7. Use database using mongoDB command : "show databases"
8. Create a School database
9. Insert many documents into School database
10. Display the documents in the database using "db"
11. Insert many documents in the collection using "insert"
12. Set the document in alphabetical order, 1-alphabetical order and -1 for reverse order.
13. Update one or many documents using "update" command
14. Delete one or many documents using "delete" command

```

{
    "acknowledged": true,
    "insertId": "2020-01-28T17:28:29Z"
}

Object ID ("64db23cc6f17a13e398f206")
Object ID ("64db23cc6f17a13e398f206") exists in wallet
Object ID ("64db23cc6f17a13e398f206"), and it is not in wallet

{
    "id": "objectID("64db23cc6f17a13e398f206")",
    "name": "Sandy", "age": 22, "gpa": 4.0, "lname": "Gandy", "gpa": 2.53}
}

School > db.Students.insertMany([{"name": "Patrice", "age": 38, "gpa": 15.3,
    "frame": "Sandy"}, {"name": "Gary", "age": 32, "gpa": 4.0}, {"name": "Sandy", "gpa": 2.53}])

{
    "id": "objectID("64db23cc6f17a13e398f206")",
    "name": "Gary", "age": 18, "gpa": 2.5}
}

School > db.Students.find().set({names: 1})
{
    "id": "objectID("64db23cc6f17a13e398f206")",
    "name": "Patrice", "age": 38, "gpa": 15.3, "frame": "Gary", "gpa": 2.53}
}

{
    "id": "objectID("64db23cc6f17a13e398f206")",
    "name": "Patrick", "age": 38, "gpa": 15.3
}

School > db.Students.find().set({names: 1})
{
    "id": "objectID("64db23cc6f17a13e398f206")",
    "name": "Sandy", "age": 22, "gpa": 4.0, "lname": "Gandy", "gpa": 2.53}
}

{
    "id": "objectID("64db23cc6f17a13e398f206")",
    "name": "Sandy", "age": 22, "gpa": 4.0}
}

{
    "id": "objectID("64db23cc6f17a13e398f206")",
    "name": "SpongeBob", "age": 30, "gpa": 3.2}
}

```

die Schwestern, die Kinder und die Eltern der Kinder sind. Sie sind alle sehr gut aufgezogen und sehr höflich.

"modified Kant": I

Id: ObjectID ("64db2c4bfef93e43011b9c"),
"name": " spongebob ", "age": 30, "gpa": 3.2 }

"acknowledged": true
"mathtt{student}": true,
"modifiedCount": 4

1

"_id": ObjectID("64db309cbf47bdb7323736a1"),
"name": "Songba", "age": 30, "gra": 32, "fulltim

`-id": objectID("64db30c6bf47bdb73a373b92"),
"name": "Patrick", "age": 38, "gpa": 1.5, "fulltime": false}`

Schad > db.studentik.UpdateOne("Name": "Spongetob", "Lppset": "fulltime":
true}) : result = p.get("Lppset", "fulltime") == "true" .
Schad > db.Student.findById (name: "Spongetob")

```
1
" - id": ObjectId("64db30c6fb7bde732373b931"), Jb > donde
  "name": "Sandy", "age": 27, "gpa": 4, "fulltime": false
}
3
3
2
" - id": ObjectId("64db30c6fb7bde732373b94"),
  "name": "Garry", "age": 18, "gpa": 25, "fulltime": true
}
3
3
1
"acknowledged": true
"deletetount": 1
}
1
" - id": ObjectId("64db3cd18d7a5e7d2fa559fc8"),
  "name": "Rabid", "age": 38, "gpa": 1.5, "fulltime": false
}
3
3
2
" - id": ObjectId("64db3cd18d7a5e7d2fa559fa"),
  "name": "Sandy", "age": 27, "gpa": 4, "fulltime": false
}
3
3
1
" - id": ObjectId("64db3cd18d7a5e7d2fa559fb"),
  "name": "Garry", "age": 18, "gpa": 25, "fulltime": false
}
3
3
1
"acknowledged": true, "deletetount": 3
}
```

School > db.students.deleteOne({ name: "Sandy" })

School > db.students.find()

School rdb.Students.tad()

(A) Insertion, Deletion, Update & Delete

1) insert, update & delete

(student: student_id: name) (teacher: teacher_id, teacher_name)

Result :-
Thus, the creation, updation and deletion commands
of no sql documents using MongoDB is executed successfully.

(1) Inst. Analysis: Map Reduce

Aim:-

To Write a python code to print word count in map reduce using cloudera.

Algorithm:-

1. Read lines from standard input ('stdin')
2. For each line, remove leading and trailing whitespace and split it into words.
3. For each word in the line, emit a key-value pair where 'the word' is the key and '1' is the value.
4. This process generates a list of key-value pairs for each word in the input.
5. Now, for the 'reduce.py' → Read lines from Standard input ('stdin')
6. Initialize variables 'current_word' and 'current_count' to None and 0 respectively.
7. For each line, split it into a word and a count.
8. If the current word is the same as the new word, add the count to the current count.
9. If the current word changes, print previous word and its total count
10. At the end of the input, make sure to print the last word and its total count if it exists.

Program:-

```
#!/usr/bin/python  
'''  
mapper.py'''  
  
# importing modules, modules will work  
# with file system to support print statement  
import sys  
  
for line in sys.stdin:  
    print line
```

```
#!/usr/bin/env python
# reducer.py : "reduces" the word count
# input comes from STDIN
# line = line.strip()
# words = line.split()
# for word in words :
#     print '%s\t% s' % (word, 1)

line = line.strip()
words = line.split()
for word in words :
    print '%s\t% s' % (word, 1)

# !/usr/bin/env python
''' reducer.py : "reduces" the word count
# input comes from STDIN
# line = line.strip()
# words = line.split()
# for word in words :
#     print '%s\t% s' % (word, 1)
# current_word = None
# current_count = 0
# word = None

# Convert count (currently a String) to int
try :
    count = int(count)
except ValueError:
    continue

# If current_word is None, then we're still
# initializing it
if current_word != word :
    if current_word :
        print '%s\t% s' % (current_word, current_count)
    # reset
    current_word = word
    current_count = 0
else:
    current_count += 1

print '%s\t% s' % (current_word, current_count)
```

operatingsystem 2

1 pig

Expt. No. _____
Expt. Name _____
Date: _____

(pig.200) = null
whether null = whatever
not our best guess

1. (best) A 200's guess

multiple individual result
answer

if current_word == word:
current_count += count
else:
if current_word: *if current_word != current_count*:
current_count = count
current_word = word
print("%s %s (%s)" % (current_word, current_count))

Result: python3 wordcount.py ./shakespeare.txt
William Shakespeare wrote 1025505
to Henry V play was not released to public domain yet as present now

multiple answers possible because
of many different ways to read
private information from shakespeare

more more real input
several possibilities

multiple answers - no problem because
of many different ways to read
private information from shakespeare

multiple answers possible because
of many different ways to read
private information from shakespeare

multiple answers possible because
of many different ways to read
private information from shakespeare

multiple answers possible because
of many different ways to read
private information from shakespeare

multiple answers possible because
of many different ways to read
private information from shakespeare

multiple answers possible because
of many different ways to read
private information from shakespeare

multiple answers possible because
of many different ways to read
private information from shakespeare

multiple answers possible because
of many different ways to read
private information from shakespeare

multiple answers possible because
of many different ways to read
private information from shakespeare

multiple answers possible because
of many different ways to read
private information from shakespeare

multiple answers possible because
of many different ways to read
private information from shakespeare

multiple answers possible because
of many different ways to read
private information from shakespeare

multiple answers possible because
of many different ways to read
private information from shakespeare

multiple answers possible because
of many different ways to read
private information from shakespeare

multiple answers possible because
of many different ways to read
private information from shakespeare

multiple answers possible because
of many different ways to read
private information from shakespeare

multiple answers possible because
of many different ways to read
private information from shakespeare

multiple answers possible because
of many different ways to read
private information from shakespeare

Output :-

```

Re-format filesystem in storage. Directly root = /tmp/hadoop-
root/dfs/name; location = null? (y or n) Y
2023-09-01 04:16:59, 475 INFO namenode.FSImage: Allocated new
BlockPoolID: BP-150640189-17a.28.0.12-1693541819454
2023-09-01 04:16:59, 475 INFO common.Storage: will remove file:
[/tmp/hadoop-root/dfs/name/current/seen.trid, /tmp/hadoop-
root/dfs/name/current/fsimage.10000000000000000000, /tmp/hadoop-
root/dfs/name/current/version, /tmp/hadoop-
root/dfs/name/current/fsimage-0000000000000000.mds]
2023-09-01 04:16:59, 515 INFO common.Storage: storage directly
/tmp/hadoop-root/dfs/name has been successfully formatted.
2023-09-01 04:16:59, 556 INFO namenode.FSImageFormatButtلاف: giving
image file /tmp/hadoop-
root/dfs/name/Current/fsimage.ckpt_0000000000000000 using no
compression
2023-09-01 04:16:59, 732 INFO namenode.FSImageFormatButtلاف: fsync file
/tmp/hadoop-root/dfs/name/current/fsimage.ckpt_0000000000000000
Size 399 bytes saved in 0 seconds
2023-09-01 04:16:59, 745 INFO namenode.NamenodeRegistrationManager: trying
to retain 1 images with trid >= 0
2023-09-01 04:16:59, 746 INFO namenode.FSNamespaceSystem: stopping servers
Started for active state [0] 0
Started for active state [0] 0
Started for standby state [0]
2023-09-01 04:16:59, 745 INFO namenode.FSNamespaceSystem: stopping servers
checkpoint: trid = 0 when meet shutdown
  
```

Aim :-

To demonstrate the running status of HDFS file system.

Procedure :-

1. Install Tarball if not already installed.
2. Download Hadoop 3.3.1
3. Extract the downloaded Hadoop archive.
4. Set environment variables from tar and Hadoop paths.
5. Configure HDFS by updating core-site.xml and hdfs-site.xml
6. Format HDFS using hdfs namenode format.
7. Verify HDFS is running using jps command.

Code :-

```

! apt-get install openjdk-8-jdk-headless -q
! wget -q https://downloads.apache.org/hadoop/common/hadoop-3.3.1/
hadoop-3.3.1.tar.gz
! tar xf hadoop-3.3.1.tar.gz
import os
os.environ['JAVA_HOME'] = '/usr/lib/jvm/java-8-openjdk-amd64'
os.environ['HADOOP_HOME'] = 'content/hadoop-3.3.1'
os.environ['HADOOP_CONF_DIR'] = 'content/hadoop-3.3.1/etc/hadoop'
os.environ['PATH'] = os.environ['HADOOP_HOME'] + '/bin:' +
os.environ['PATH']

# Configure HDFS
! cp /content/hadoop-3.3.1/etc/hadoop/core-site.xml /content/hadoop-
3.3.1/etc/hadoop/core-site.xml.back
  
```

2023-09-01 04:16:59,801 INFO namenode.Namenode: SHUTDOWN_MSG: /*****

SHUTDOWN_MSG: Shutting down Namenode at ba73d57e1e43/172.28.0.12

2023-JPS

Re-format filesystem in Storage Directory root = /tmp/hadoop-root/hdfs/names-

root/dts/name; location = null? (y or n) N

Format aborted via storage Directory root = /tmp/hadoop-root/hdfs/names-

location=null? (y or n) N

2023-09-01 04:38:05,514 INFO namenode.FSNamesystem: Stopping Service

Started for active state

2023-09-01 04:38:05,515 INFO namenode.FSNamesystem: Stopping Service

Started for standby state.

2023-09-01 04:38:05,520 INFO util.ExitUtil: Exiting with status 1: [0]

ExitException thrown from org.apache.hadoop.util.ExitUtil - [0]

2023-09-01 04:38:05,521 INFO namenode.Namenode: SHUTDOWN_MSG:[/*****

SHUTDOWN_MSG: Shutting down Namenode at ba73d57e1e43/172.28.0.12

2023-JPS

Re-format filesystem in Storage Directory root = /tmp/hadoop-root/hdfs/names-

root/dts/name; location = null? (y or n) N

Format aborted via storage Directory root = /tmp/hadoop-root/hdfs/names-

location=null? (y or n) N

2023-09-01 04:38:05,514 INFO namenode.FSNamesystem: Stopping Service

Started for active state

2023-09-01 04:38:05,515 INFO namenode.FSNamesystem: Stopping Service

Started for standby state.

2023-09-01 04:38:05,520 INFO util.ExitUtil: Exiting with status 1: [0]

ExitException thrown from org.apache.hadoop.util.ExitUtil - [0]

2023-09-01 04:38:05,521 INFO namenode.Namenode: SHUTDOWN_MSG:[/*****

SHUTDOWN_MSG: Shutting down Namenode at ba73d57e1e43/172.28.0.12

2023-JPS

Re-format HDFS

!hdfs namenode -format

2023-09-01 04:38:05,521 INFO namenode.Namenode: SHUTDOWN_MSG:[/*****

SHUTDOWN_MSG: Shutting down Namenode at ba73d57e1e43/172.28.0.12

2023-JPS

verify HDFS is running

!jps

2023-09-01 04:38:05,521 INFO namenode.Namenode: SHUTDOWN_MSG:[/*****

SHUTDOWN_MSG: Shutting down Namenode at ba73d57e1e43/172.28.0.12

2023-JPS

Re-format HDFS

!hdfs namenode -format

2023-09-01 04:38:05,521 INFO namenode.Namenode: SHUTDOWN_MSG:[/*****

SHUTDOWN_MSG: Shutting down Namenode at ba73d57e1e43/172.28.0.12

2023-JPS

Re-format HDFS

!hdfs namenode -format

2023-09-01 04:38:05,521 INFO namenode.Namenode: SHUTDOWN_MSG:[/*****

SHUTDOWN_MSG: Shutting down Namenode at ba73d57e1e43/172.28.0.12

2023-JPS

Re-format HDFS

!hdfs namenode -format

2023-09-01 04:38:05,521 INFO namenode.Namenode: SHUTDOWN_MSG:[/*****

SHUTDOWN_MSG: Shutting down Namenode at ba73d57e1e43/172.28.0.12

2023-JPS

Re-format HDFS

!hdfs namenode -format

2023-09-01 04:38:05,521 INFO namenode.Namenode: SHUTDOWN_MSG:[/*****

SHUTDOWN_MSG: Shutting down Namenode at ba73d57e1e43/172.28.0.12

2023-JPS

Re-format HDFS

!hdfs namenode -format

2023-09-01 04:38:05,521 INFO namenode.Namenode: SHUTDOWN_MSG:[/*****

SHUTDOWN_MSG: Shutting down Namenode at ba73d57e1e43/172.28.0.12

2023-JPS

Re-format HDFS

!hdfs namenode -format

2023-09-01 04:38:05,521 INFO namenode.Namenode: SHUTDOWN_MSG:[/*****

SHUTDOWN_MSG: Shutting down Namenode at ba73d57e1e43/172.28.0.12

2023-JPS

Re-format HDFS

!hdfs namenode -format

2023-09-01 04:38:05,521 INFO namenode.Namenode: SHUTDOWN_MSG:[/*****

SHUTDOWN_MSG: Shutting down Namenode at ba73d57e1e43/172.28.0.12

2023-JPS

Re-format HDFS

!hdfs namenode -format

2023-09-01 04:38:05,521 INFO namenode.Namenode: SHUTDOWN_MSG:[/*****

!cp /content/hadoop-3.3.1/etc/hadoop/hdfs-site.xml /content/hadoop-3.3.1/etc/hadoop/hdfs-site.xml.bak
!echo -e 'dfs.default.name'<value></value></property><configuration>'> /content/hadoop-3.3.1/etc/hadoop-site.xml
!echo -e '<configuration><property><name>dfs.replication</name><value>'></value></property></configuration>'> /content/hadoop-3.3.1/etc/hadoop-site.xml
#Format HDFS
!hdfs namenode -format
verify HDFS is running
!jps

./wordcount -mrx-size-2tbn/pigball/stl.1.5.5 -ppbball/trivial sp
 wordcount
 >edit wordcount.pig

>edit testfile.txt

5- add

< pig -x local wordcount.pig > < wordcount > < wordcount >
 Input(s): < wordcount > < wordcount > < wordcount > < wordcount >
 successfully read records from: < /home/cloudera/documents/testfile.txt >

< output(s): < wordcount > < wordcount > < wordcount > < wordcount >

successfully stored records in: < file:///tmp/tmp1234539/tmp-405488 >

Job DAG:

job-local123453066_0001

(a, 2)

(bt, 1)

(in, 1)

(is, 3)

(to, 1)

(pig, 1)

(and, 1)

(bt, 1)

(fun, 1)

(pig, 1)

(thus, 2)

(thus, 2)

>>> (word,1) = wordcount
 (pig,1) ~~OP~~ ~~ON~~
 (text,1) ~~OP~~ ~~ON~~
 (text,1)

Aim:- To print the wordcount using pig in cloudera.

Algorithm:-

1. Load the input text file 'testfile.txt' into a relation called 'lines' and treat each line as a character array ('chararray')
2. Tokenize each line into words using 'split(')' as the delimiter and generate a relation called 'words' where each word is flattened.
3. Group the 'words' relation by the 'word' column
4. Calculate the count of words within each group and create a relation called 'wordCount' that contains the word and its corresponding count.
5. Dump the 'wordCount' relation to the console, displaying the word and its count.

Program:-

-wordcount.pig

lines = LOAD '/home/cloudera/Documents/testfile.txt' AS (line:chararray);

words = foreach lines GENERATE FLATTEN(TOKENIZE(line));

grouped = GROUP words BY word;

wordCount = foreach grouped GENERATE group, COUNT(words);

dump wordCount;

Result:-

Thus the execution of the wordcount using pig in cloudera is successful.

Ques

Aim:-
To find word count in file.

Output:-

word	Count
salt	2
soil	2
water	2
land	2
tree	2
sun	2
clouds	2
grass	1
red	1
green	1
yellow	1

Note :- The words which are not in the output are removed.

- ### Algorithm:-
- Start
 - Create table files
 - Read data from input files
 - Select word from files
 - Group and Order the words
 - Stop.

Program:-

```
CREATE TABLE FILES(CLINESTRING);
LOAD DATA IN PATH 'data' OVERWRITE INTO TABLE FILES;
CREATE TABLE WORDCOUNTS AS
SELECT word, count(*) AS COUNT FROM
SELECT replicate(SPLIT(line, ' '), AS FROM FILES)
GROUP BY word
ORDER BY word
```

Note :- The words which are not in the output are removed.

Result:- ✓
Thus the above program has been executed successfully.

Aim:-

To implement YARN connectivity using python.

Algorithm:-

Step 1: Download YARN V3.8.0 from Yarn site - <https://github.com/apache/yarn>
 Step 2: YARN config get and set
 Step 3: Create rpm authToken -> https://www.yarn.apache.org/docs/rpm.html

Program:-

yarn set version latest

yarn set version canary

yarn set version classic

yarn set version 3.x

yarn set version 3.0.0-rc.10

yarn set version 1.0.0

yarn config get

yarn config get yarn path

yarn config get package Extension

yarn config get 'npm scopes="my-company".npm register server'

yarn config get npm auth Token - no redundant

yarn config get package Extensions - json

yarn config set

yarn config set init scope my scope

Output:-

> latest/berry/stable → The most recent stable berry version is v1.0.0
(>= 2.0.0) relax

> cocony → The most recent cocony release (v1.0.0)

berry (>= 2.0.0) relax

> classic → The most recent classic (v0.11.11.1) release

> Yarnpath: ./scripts/yarn -v 2.0.0 -rc 001.js

> Success set "init - license" to "BSD - 2 clause" version is v2.0.0

> *-* done in 0.05s

npmAuth Token: "XXXXXXXXX-XXXX-XXXX-XXXX-XXXXXXXXXXXX" ready

Yarn config set 'npm Registry' ["https://npm.example.com"].npm

Yarn config set 'Auth Token' • "XXXXXXXXX-XXXX-XXXX-XXXX-XXXXXXXXXXXX"

Yarn config set unsafe Help whitelist ---.json ["*.example.com"]
Yarn config set package Extensions - json i " @label/prefix@* " :
dependencies": { "@label": "1.3.3" }

OK

Yarn config set unsafe Help whitelist ---.json ["*.example.com"]
Yarn config set package Extensions - json i " @label/prefix@* " :
dependencies": { "@label": "1.3.3" }

Result:-
Thus the above program has been executed successfully.

Issue: Suppose you want to send some data from Node.js to Python script. How do you do it?

Answer: We can use Python's `subprocess` module.

Process: Python program will make a process named `subprocess` which will run the command in the terminal.

Output: The output of the command will be sent back to the Python program.

Code: In node.js, we can use `child_process` module to run the command.

Example: `child_process.spawn('python myscript.py')`

Output: The output of the command will be sent back to the Python program.

Advantages: We can use this method to interact with Python programs.

Disadvantages: This method is not memory efficient and will result in memory leak.

Aim:-

To implement your program using Python.

Algorithm:-

1. Start
2. Import required libraries in Python.
3. Running the Python Script
4. Exchanging data between Node and Python
5. Stop.

Program:-

```
npm install python-shell
import {PythonShell} from 'python-shell';
PythonShell.runString ('x=1+1;print(x)',null,function (err) {
if (err) throw err;
console.log('finished');
});
let pyshell = new PythonShell('my-script.py');
pyshell.send ("hello");
pyshell.on ('message', function (message) {
console.log(message);
});
pyshell.end (function (err,code,signal) {
if (err) throw err;
console.log ('The exit code was : '+code);
console.log ('The signal was : '+signal);
console.log ('finished');
});
```

Output:

Filename : foopy
Inside file : hello world

~~Output~~
writing to standard output and standard error
inside running file program
writing to standard output file foopy

```
pythonshell.default_options = { 'script_path' : ... /Scripts '3;
pythonshell.run('script.py',null, function(error, results) {
    console.log(`Hello world!`); file = open('foopy', "w");
    print("Hello world"); file = open('foopy', "w");
    import { PythonShell, NewlineTransformer, options } from 'python-shell';
    python_shell.setOptions({ options: {
        'stdio': [
            'pipe', 'pipe', 'pipe', 'pipe'
        ]
    }
});
```

```
const pyshell = new PythonShell('FOO.py', options);
const customPipe = pyshell.childProcess.stdio[3];
customPipe.pipe(process.stdout);
Result = Buffer.from(customPipe.read(1));
process.stdout.write(result.toString());
```

```
if (process.stdin === process.stdout) {
    process.stdin = process.stdout;
}
```

```
if (process.stdin === process.stderr) {
    process.stderr = process.stdout;
}
```

```
if (process.stdin === process.stdin) {
    process.stdin = process.stdin;
}
```

```
if (process.stdin === process.stderr) {
    process.stderr = process.stderr;
}
```

```
if (process.stdin === process.stdout) {
    process.stdout = process.stdout;
}
```

```
if (process.stdin === process.stderr) {
    process.stderr = process.stderr;
}
```

```
if (process.stdin === process.stdin) {
    process.stdin = process.stdin;
}
```

```
if (process.stdin === process.stderr) {
    process.stderr = process.stderr;
}
```

Result:-
~~Output~~
Thus the above program has been executed successfully.

↳ `df = pd.read_csv('C:/Users/Asus/Desktop/Unacademy.csv')` working
↳ `df`

Output: most worked company, working hours, working place

Marks.csv

NAME	AGE	YEAR	MARKS
BALU	19	1st	85
ANIL	20	2nd	90
MATESH	18	3rd	75
NADAV	15	2nd	74
SUNNY	20	3rd	98
NEETH	19	2nd	57
REDDY	20	2nd	65
SANATH	19	3rd	74
KAVAN	18	2nd	75
KUSHIK	19	2nd	70
HARR	20	3rd	84
ROHITH	19	3rd	98

↳ `df`

most worked company, working hours, working place

↳ `df = pd.read_csv('C:/Users/Asus/Desktop/Unacademy.csv')` working
↳ `df`

: 'df'

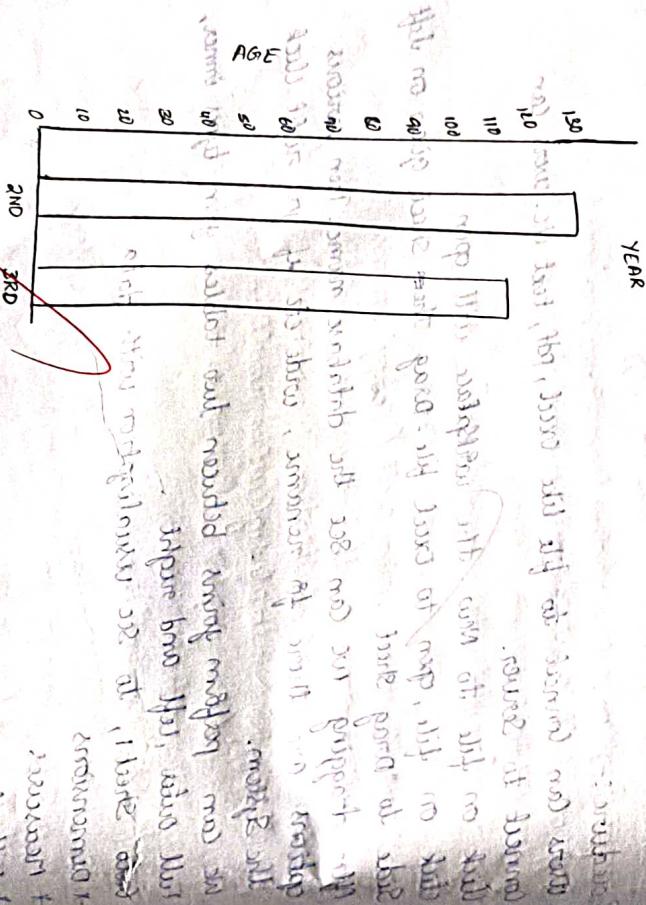
[85, 90, 75, 74]

1. we can connect to file like excel, pdf, text etc. users can connect to server.
2. click on file to New the workplace will open
3. click on file, open to excel file. drag these sheet given on left side to "drag sheet".
4. after dragging we can see the database name. Then various options are there for rename, width etc. if you right click the system.
5. we can perform joins between two tables join types inner, full outer, left and right.
6. goto sheet 1, to see visualization with data

- * dimensions
- * measures
- * filters
- * show me
- 7. we will drag values from dimensions a measures and make a visualization.
- 8. we can change sum to average, media.
- 9. we can change the size of the graph click on size → adjust size as per your wish
- 10. Drag region field into the filter, we can filter by region wise. The options are general, wild card, condition, Top.

Ques. P. missionaries mission have collected after 1800 of India.

11. Shortcut keys to play with data.
12. Tableau Dashboard contains more than one sheet, select color, size, rename as per the design we want.
13. The two sheets are dragged into dashboard. In two bar graphs, one is for sum of sales other is for number of records from Sheets. i.e.,



Result :- Thus the data has been visualized using tableau and the required information has been acquired.

After running the code, we get the following output:

```

SELECT COUNT(*) AS COUNT_OF_RECORDS FROM SHEET1;
+-----+
| COUNT_OF_RECORDS |
+-----+
|          400    |
+-----+

```

Next we run the following query, we get the following output:

```

SELECT AGE, COUNT(*) AS COUNT_OF_RECORDS FROM SHEET1 GROUP BY AGE;
+-----+-----+
| AGE | COUNT_OF_RECORDS |
+-----+-----+
|   0  |          120    |
| 10- |          100    |
| 20- |          80     |
| 30- |          60     |
+-----+-----+

```