



# Major Assignment – Machine Learning

SUMIT SHAMLAL CHAURE

---

Data Science with Python by SkillAcademy

Batch 10

23/10/2023

Report for Assignment

[Google Colab](#) (Live Demo)

[Drive Folder](#) (Contains the report -pdf/word, Notebook, dataset etc)

[Assignment PDF](#) (Tasks & Details)

[Dataset](#) (Oil Spill Dataset)

**Note :-** *In the assignment I have attached the snippet picture as typing or copy-pasting everything and formatting is time consuming (You can view the colab link or notebook in drive to check the whole code for reference).*

*At some places I have skipped giving the results output screenshot as its long also few random code might not be added as it will make the report long. Do look into the colab link or inside zip to check for actual results.*

## INDEX

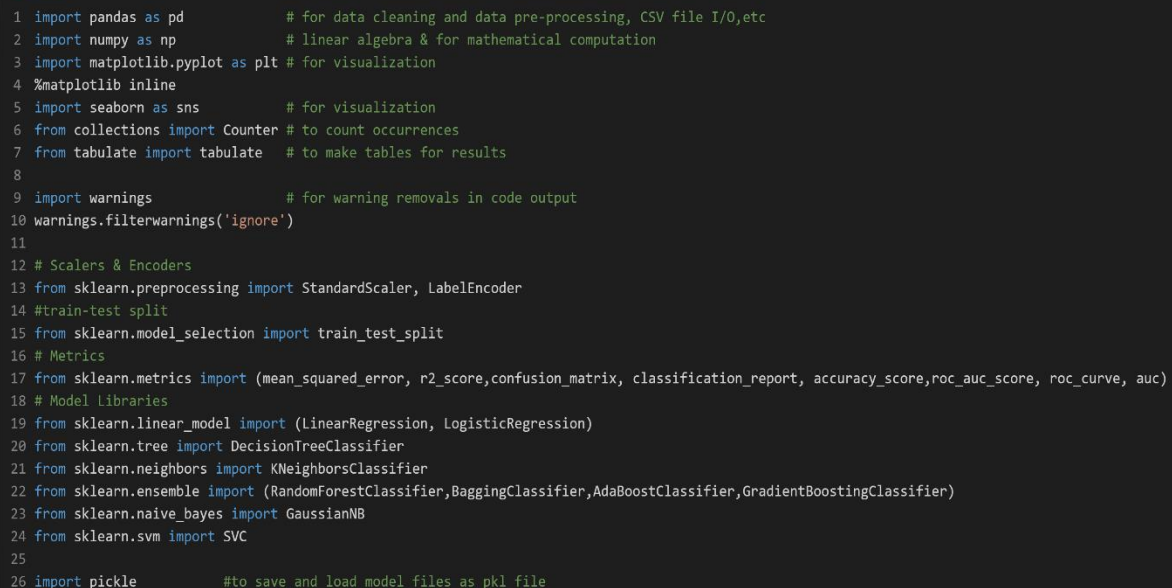
<i>Introduction</i> .....	1
<i>Overview</i> .....	2
<i>Approach</i> .....	2
<b>Question &amp; Answers</b> .....	3
Q1. <b>Download the Oil Spill Dataset and perform Data cleaning and Data Pre-Processing if Necessary.</b> *(Data cleaning & Processing Answer Continued in Q2).....	3
Q2. <b>Use various methods such as Handling null values, One-Hot Encoding, Imputation, and Scaling of Data Pre-Processing where necessary.</b> .....	4
Q3. <b>Derive some insights from the dataset.</b> .....	6
Q4. <b>Apply various Machine Learning techniques to predict the output in the target column, make use of Bagging and Ensemble as required, and find the best model by evaluating the model using Model evaluation techniques.</b> .....	9
Q5. <b>Save the best model and Load the model.</b> .....	12
Q6. <b>Take the original data set and make another dataset by randomly picking 20 data points from the oil spill dataset and applying the saved model to the same.</b> .....	14
<i>End</i> .....	22

## Overview

To Solve questions to test knowledge about **Machine Learning project life-cycle i.e. Starting from Data cleaning & pre-processing, Handling imputations, Data analysis & EDA, Model selection and evaluation techniques, model training & finding best fit models, optimizing the models and saving the best one, loading the saved model and testing against new set of data and checking the actual metrics such as recall, precision & accuracy of the trained model, re-iteration.**

## Approach:

### 1 - Library



```

1 import pandas as pd          # for data cleaning and data pre-processing, CSV file I/O, etc
2 import numpy as np           # linear algebra & for mathematical computation
3 import matplotlib.pyplot as plt # for visualization
4 %matplotlib inline
5 import seaborn as sns        # for visualization
6 from collections import Counter # to count occurrences
7 from tabulate import tabulate # to make tables for results
8
9 import warnings              # for warning removals in code output
10 warnings.filterwarnings('ignore')
11
12 # Scalers & Encoders
13 from sklearn.preprocessing import StandardScaler, LabelEncoder
14 # train-test split
15 from sklearn.model_selection import train_test_split
16 # Metrics
17 from sklearn.metrics import (mean_squared_error, r2_score, confusion_matrix, classification_report, accuracy_score, roc_auc_score, roc_curve, auc)
18 # Model Libraries
19 from sklearn.linear_model import (LinearRegression, LogisticRegression)
20 from sklearn.tree import DecisionTreeClassifier
21 from sklearn.neighbors import KNeighborsClassifier
22 from sklearn.ensemble import (RandomForestClassifier, BaggingClassifier, AdaBoostClassifier, GradientBoostingClassifier)
23 from sklearn.naive_bayes import GaussianNB
24 from sklearn.svm import SVC
25
26 import pickle                # to save and load model files as pkl file
  
```

### 1- Library Imports (all)

2 - Data Operations (Main Query) - Will paste the screenshot of things from notebook file as typing or copying pasting it will create white spaces and unnecessary time in formatting)

#### a. Question

b. **Code** (Screenshot - since long lines will take space and create whitespace) & **Result** (Screenshot of plots/graphs - some questions might have multiple depending on the approach I take)

c. **Insights & Approach** (Screenshot of the derived data insights & approach to solve if any test involved)

# Question & Answers

**Q1. Download the Oil Spill Dataset and perform Data cleaning and Data Pre-Processing if Necessary. \***(Data cleaning & Processing Answer Continued in Q2)

## Code Snippet:

```
1 # 2.1) Importing the dataset (With error handling)
2 # If you want to upload the dataset directly (Since on Google Colab it will be lost on re-run) - uncomment the below 2 line code and run
3 # from google.colab import files
4 # uploaded = files.upload()
5
6 file_path = "oil_spill.csv"
7 file_name = file_path.split("/")[-1]
8
9 try:
10     # Reading the CSV file into a Pandas DataFrame
11     df = pd.read_csv(file_path)
12     # Store the filename as an attribute in the DataFrame
13     df.file_name = file_name
14     print(f"\n '{df.file_name}' loaded successfully.")
15
16 # Exception to check if the file has some error like no file at the path, etc.
17 except FileNotFoundError:
18     print(f"Error: '{file_name}' not found at the specified location {
19         file_path}.")
20 except Exception as e:
21     print(f"An unexpected error occurred: {e}")
```

## Results & Insights:

41print(f"An unexpected error occurred: {e}")

[35]✓0.0s

...

'oil\_spill.csv' loaded successfully.

+ Code+ Markdown

Click here to ask Blackbox to help you code faster

1df

[3]✓0.2s 展开 'df'

...

		f_1	f_2	f_3	f_4	f_5	f_6	f_7	f_8	f_9	f_10	...	f_41	f_42	f_43	f_44	f_45	f_46	f_47	f_48	f_49	target
0	1	2558	1506.09	456.63	90	6395000	40.88	7.89	29780.0	0.19	...	2850.00	1000.00	763.16	135.46	3.73	0	33243.19	65.74	7.95	1	
1	2	22325	79.11	841.03	180	55812500	51.11	1.21	61900.0	0.02	...	5750.00	11500.00	9593.48	1648.80	0.60	0	51572.04	65.73	6.26	0	
2	3	115	1449.85	608.43	88	287500	40.42	7.34	3340.0	0.18	...	1400.00	250.00	150.00	45.13	9.33	1	31692.84	65.81	7.84	1	
3	4	1201	1562.53	295.65	66	3002500	42.40	7.97	18030.0	0.19	...	6041.52	761.58	453.21	144.97	13.33	1	37696.21	65.67	8.07	1	
4	5	312	950.27	440.86	37	780000	41.43	7.03	3350.0	0.17	...	1320.04	710.63	512.54	109.16	2.58	0	29038.17	65.66	7.35	0	
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
932	200	12	92.42	364.42	135	97200	59.42	10.34	884.0	0.17	...	381.84	254.56	84.85	146.97	4.50	0	2593.50	65.85	6.39	0	
933	201	11	98.82	248.64	159	89100	59.64	10.18	831.0	0.17	...	284.60	180.00	150.00	51.96	1.90	0	4361.25	65.70	6.53	0	
934	202	14	25.14	428.86	24	113400	60.14	17.94	847.0	0.30	...	402.49	180.00	180.00	0.00	2.24	0	2153.05	65.91	6.12	0	
935	203	10	96.00	451.30	68	81000	59.90	15.01	831.0	0.25	...	402.49	180.00	90.00	73.48	4.47	0	2421.43	65.97	6.32	0	
936	204	11	7.73	235.73	135	89100	61.82	12.24	831.0	0.20	...	254.56	254.56	127.28	180.00	2.00	0	3782.68	65.65	6.26	0	

937 rows x 50 columns

Sumit S. C

**Insights:** As the assignment is part of major submission i tried to add few exception handling steps to verify things like file not found error. (If we try to import any empty csv in the dataframe it will prompt an exception suggesting no file found at the specified location.)

- I have added an extra commented code for google colab imports directly to upload the csv file from local drive and then start the further process.
- To use on google colab uncomment the start line of code to import the file on google colab drive instance.
- I have not added extraneous code to get the filename from uploaded files and save them directly instead hardcoded/kept it so that even if code is run offline it will work perfectly fine.

**Note:** As Q.2 needs to do data pre-processing i have added the data cleaning steps in that section directly.

**Q2. Use various methods such as Handling null values, One-Hot Encoding, Imputation, and Scaling of Data Pre-Processing where necessary.**

**Code Snippet:**

```

2.2) Showing Basic Dataset Information

1 # Check the shape of the DataFrame
2 print("\nShape of the DataFrame:")
3 print(df.shape)
4 print(df.size)
5 num_rows, num_columns = df.shape
6 print(f"Rows: {num_rows}, Columns: {num_columns}")
7
8 # Display information about the dataset
9 print(f"\nDataset information for {df.file_name}:")
10 df.head(3)
11 df.tail(3)
12 print("\nDataset information:")
13 print(df.info())

[111] ✓ 0.0s

Shape of the DataFrame:
(937, 50)
46850
Rows: 937, Columns: 50

Dataset information for oil_spill.csv:

Dataset information:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 937 entries, 0 to 936
Data columns (total 50 columns):
#   Column  Non-Null Count  Dtype
---  -
0    f_1     937 non-null    int64
1    f_2     937 non-null    int64
2    f_3     937 non-null    float64

```

```

Columns of the dataset

1 # Display the columns & rows of dataset
2 print(f"The columns of our {file_name} dataframe\n")
3 print(df.columns)

[112] ✓ 0.0s

''' The columns of our oil_spill.csv dataframe

Index(['f_1', 'f_2', 'f_3', 'f_4', 'f_5', 'f_6', 'f_7', 'f_8', 'f_9', 'f_10',
       'f_11', 'f_12', 'f_13', 'f_14', 'f_15', 'f_16', 'f_17', 'f_18', 'f_19',
       'f_20', 'f_21', 'f_22', 'f_23', 'f_24', 'f_25', 'f_26', 'f_27', 'f_28',
       'f_29', 'f_30', 'f_31', 'f_32', 'f_33', 'f_34', 'f_35', 'f_36', 'f_37',
       'f_38', 'f_39', 'f_40', 'f_41', 'f_42', 'f_43', 'f_44', 'f_45', 'f_46',
       'f_47', 'f_48', 'f_49', 'target'],
      dtype='object')
'''

```

*I will paste the code part of the data cleaning as the operation gives long dataframe values so every output will be really long and most part of the data is cleaned – the Dataset does not contain any null values, duplicates, NA also no categorical column is present so we don't need to perform one-hot encoding to label them or impute the value also for the part of scaling as we need to predict the target data we let the imbalance for the training instead of balancing the dataset.*

*As the dataset is cleaned and the output results are really long adding screenshot of each and every check will take time and pages so just adding the screenshots of code snippets for data pre-processing, cleaning and checking.*

```

3) Data Checks to Perform

• Check Missing Values.
• Check Duplicates
• Check Data Type
• Check Unique Values
• Check Data Statistics
• Check Categorical Columns & Data

```

Oil\_Spill\_Classification.ipynb

Code

Markdown

Run All

Clear All Outputs

Outline

Detecting Kernels

- Check Duplicates
- Check Data Type
- Check Unique Values
- Check Data Statistics
- Check Categorical Columns & Data

Missing Values/Null Check

Click here to ask Blackbox to help you code faster

```

1 print("Missing values in the dataset:\n")
2 print(df.isnull().sum())
3 # NA value calculation
4 nullval = df.isna().sum() "nullval": Unknown word.
5 nullval = nullval[nullval > 0] "nullval": Unknown word.
6 print("\nSum of Missing values:\n", nullval) "nullval": Unknown word.

```

Python

No Null values or NA are present in our dataset.

Duplicate Values Check

Click here to ask Blackbox to help you code faster

```

1 print("\nChecking for duplicated values:\n")
2 print(df.duplicated())
3 print("\nSum of Duplicated Values in Dataframe :", df.duplicated().sum()) "Dataframe"

```

Python

The sum of duplicated value comes to be zero indicating that no duplicates are present in the dataset.

Oil\_Spill\_Classification.ipynb

Code

Markdown

Run All

Clear All Outputs

Outline

Detecting Kernels

```

1 # the number of missing values or null values in df
2 total_missing_values = df.isnull().sum().sum()
3 print("The number of missing values/NA in dataframe :", total_missing_values) "dataframe"
4
5 # Calculate the total number of values in df (excluding the missing values)
6 total_values = df.size
7 print("Total number of values in dataframe :", total_values) "dataframe": Unknown word
8
9 # percentage of missing values or null values in df
10 percentage_missing_values = (total_missing_values / total_values) * 100
11 print("Percentage of missing values or null values in df :",
    percentage_missing_values)

```

Python

No Missing values are present in the dataset.

Check Value Counts

Click here to ask Blackbox to help you code faster

```

1 print("Value counts of the dataset by datatypes")
2 df.dtypes.value_counts() "dtypes": Unknown word.

```

Python

Unique Values Check

Click here to ask Blackbox to help you code faster

```

1 print("Unique Value counts inside each columns")
2 df.nunique() "nunique": Unknown word.

```

Python

0 0 0 116

0

Connect

Live Share

1 hr 19 mins

Share Code Link

Explain Code

Comment Code

Code Chat

Discovering Python Interpreters

Cell 24 of 131

Blackbox

✓ Spoil

Sumit S. C

Oil\_Spill\_Classification.ipynb

Code

Markdown

Run All

Clear All Outputs

Outline

Detecting Kernels

Datatype Check

Click here to ask Blackbox to help you code faster

```

1 # Check data types of each column
2 print("Data types of each column:")
3 df.dtypes "dtypes": Unknown word.

```

Python

Statistics Summary

Click here to ask Blackbox to help you code faster

```

1 print("\nSummary statistics:\n")
2 print(df.describe())

```

Python

Categorical Column Check

Click here to ask Blackbox to help you code faster

```

1 df.select_dtypes(include="category") "dtypes": Unknown word.
2 categorical_columns = df.select_dtypes(include=["object"]).columns "dtypes": Unknown word.
3 print("\nCategorical columns:")
4 print(categorical_columns)

```

Python

Checking Descriptive Statistics

Click here to ask Blackbox to help you code faster

```

1 print(f"The descriptive Stats for the {file_name} dataset:")
2 df.describe()

```

Python

Click here to ask Blackbox to help you code faster

```

1 print("Complete Stats of every column")
2 print(df.describe())

```

Python

0 0 0 116

0

Connect

Live Share

1 hr 19 mins

Share Code Link

Explain Code

Comment Code

Code Chat

Discovering Python Interpreters

Sumit S. C

### Q3. Derive some insights from the dataset.

#### Code Snippet:

Q3) Derive some insights from the dataset.

Class Distribution (Target column)

```
1 # Basic Class summary
2 print("\nClass distribution:\n")
3 print(df["target"].value_counts())
4
5 # summarize the class distribution
6 target = df.values[:, -1]
7 counter = Counter(target)
8 print("\nClass Distribution Summary:\n")
9 for k, v in counter.items():
10     per = v / len(target) * 100
11     print("Class=%d, Count=%d, Percentage=%.3f%%" % (k, v, per))
```

Class distribution:

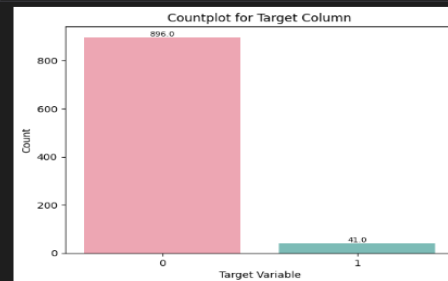
```
target
0    896
1     41
Name: count, dtype: int64
```

Class Distribution Summary:

```
Class=1, Count=41, Percentage=4.376%
Class=0, Count=896, Percentage=95.624%
```

Sumit S. C

```
1 # Countplot for Target Variable
2 ax = sns.countplot(x=df["target"], palette="husl", alpha=0.7)
3 plt.title("Countplot for Target Column")
4 plt.xlabel("Target Variable")
5 plt.ylabel("Count")
6 # Loop for annotation
7 for p in ax.patches:
8     ax.text(
9         p.get_x() + p.get_width() / 2.0,
10        p.get_height(),
11        f'{p.get_height()}',
12        ha="center",
13        va="bottom",
14        fontsize=8,
15        color="black",
16    )
17 plt.show()
```



The target column on which we need to work for our classification shows that the dataset indicates:

- 896 Non oil-spill data/regions or roughly 95.6% of dataset
- 41 Oil-spill data/regions i.e. around 4.4% of dataset.
- The above information will help us to design and test our model to check the predictions and its accuracy.

[More info on color palette](#)

1. GFG
2. Seaborn Color

Sumit S. C

Histogram to see the data analysis

```
1 fig = plt.figure(figsize=(25, 35))
2 ax = fig.gca()
3
4 _ = df.hist(ax=ax, color="green", edgecolor="black")
5 # Add a title at the top of the subplots
6 plt.suptitle("Histograms of DataFrame Columns", y=0.90, fontsize=24)
7
8 plt.show()
```

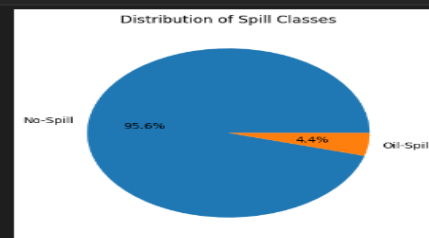
Boxplots for the outlier checking of dataset

```
1 # Set the number of subplots per row
2 subplots_per_row = 5
3 # Calculate the number of rows needed based on the number of columns and subplots per row
4 num_rows = (len(df.columns) - 1) // subplots_per_row + 1
5 # Set up the subplots
6 fig, axes = plt.subplots(
7     nrows=num_rows, ncols=subplots_per_row, figsize=(25, 25))
8 # Flatten the axes array for easy iteration
9 axes = axes.flatten()
10 # Iterate over each column and create boxplots
11 for ax, column in zip(axes, df.columns):
12     sns.boxplot(x=df[column], ax=ax, color="orange", width=0.5)
13     ax.set_title(column, fontsize=14)
14     ax.set_xlabel("Count")
15     ax.set_ylabel("Values")
16 # Adjust layout for better spacing between subplots
17 plt.tight_layout()
18 # Add a common title at the top of the subplots
19 fig.suptitle("Boxplots of DataFrame Columns (To Display Outliers in Dataset)",
20             y=1.02,
21             fontsize=24)
22
23 plt.show()
```

Sumit S. C

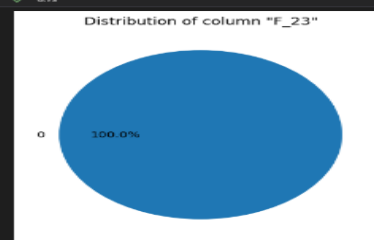
Simple Pie-chart to see target column distribution

```
1 piechart = df["target"].value_counts()
2 # Create a pie chart with labels, numbers, and percentages
3 plt.pie(piechart, labels=["No-Spill", "Oil-Spill"],
4         autopct="%0.1f%%", radius=1)
5 plt.title("Distribution of Spill Classes")
6 plt.show()
```



Pie-chart to show the single value in column F\_23

```
1 f_23_distribution = df["F_23"].value_counts()
2 labels = f_23_distribution.index.map(str)
3 values = f_23_distribution.values
4 plt.pie(values, labels=labels, autopct="%0.1f%%", radius=1)
5 plt.title("Distribution of column 'F_23'")
6 plt.show()
```



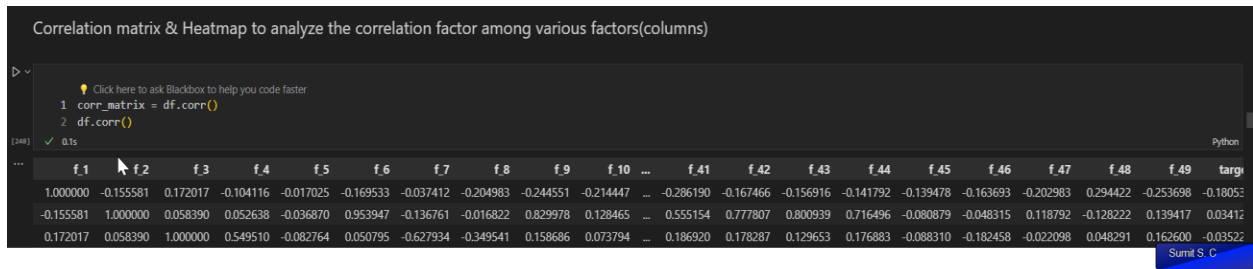
Sumit S. C



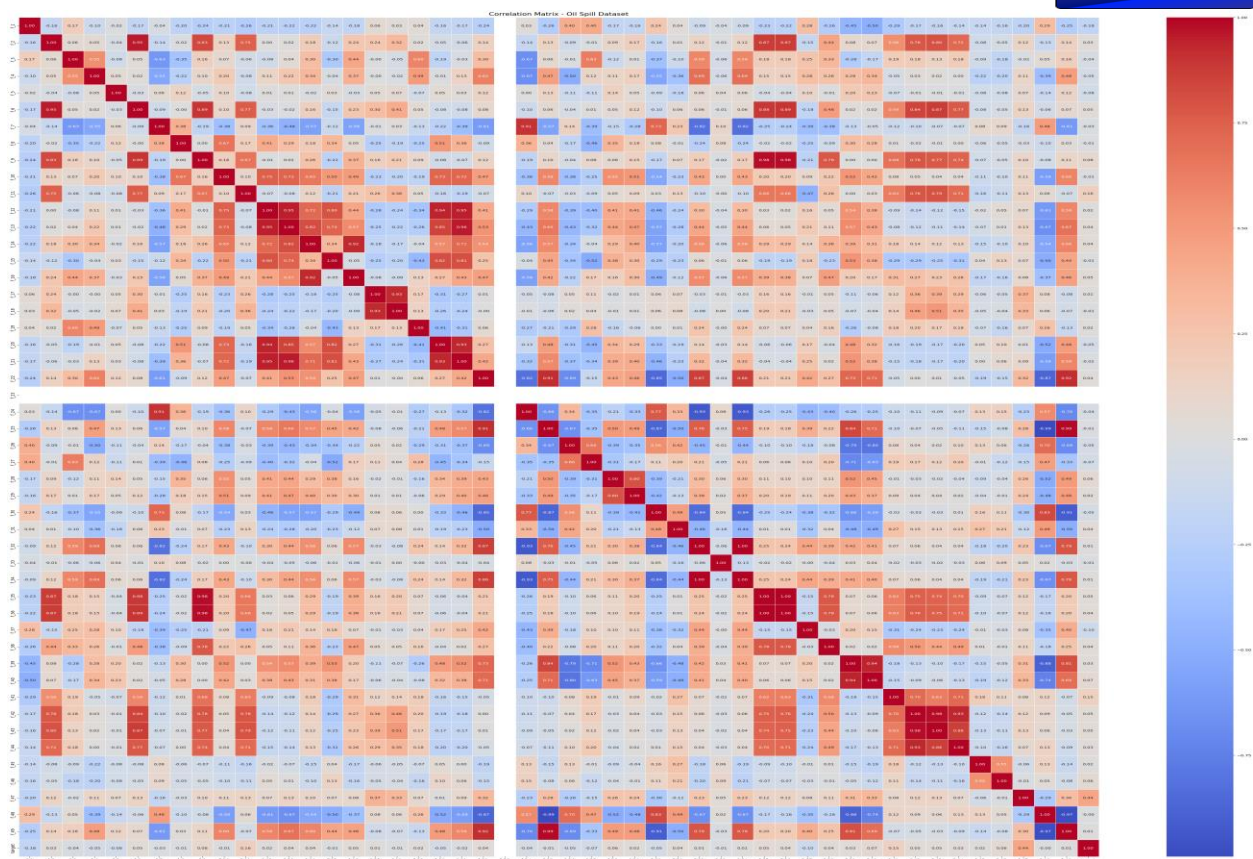
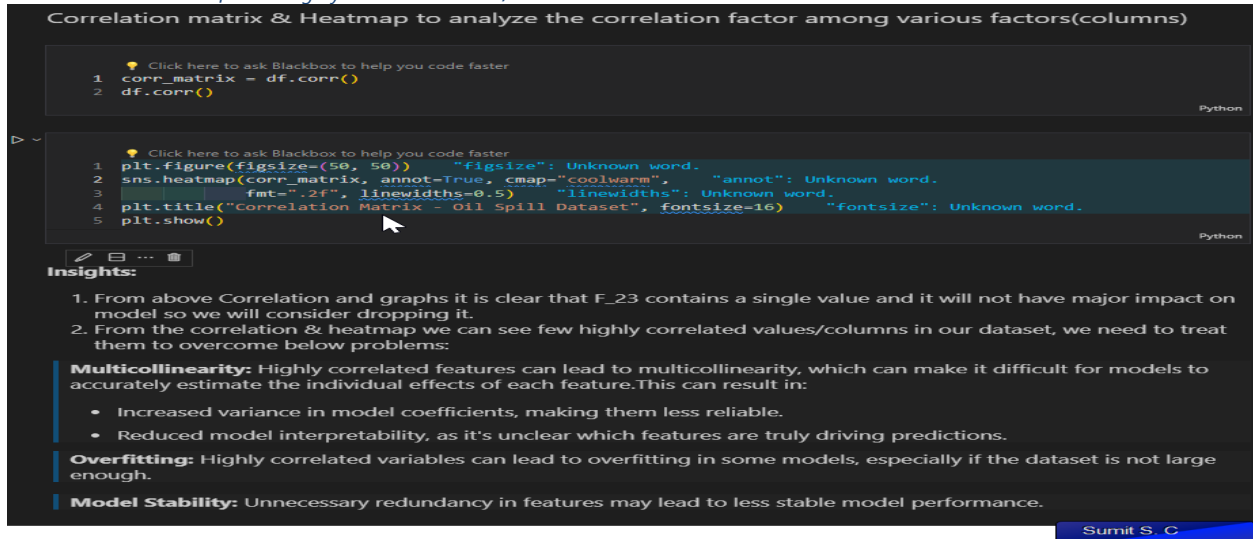


3 - Histogram, Boxplot & Piechart (prev. page) - to give the estimates pf data.





#### 4- correlation to check for the highly correlated values/columns



#### 5- Heatmap to show the correlated columns

### Removing Highly Correlated Columns

```

1 # Correlation matrix
2 corr_matrix = df.corr()
3
4 # Selecting upper triangle of correlation matrix
5 upper = corr_matrix.where(
6     np.triu(np.ones(corr_matrix.shape), k=1).astype(bool))
7
8 # Find features/columns with correlation greater than 0.90
9 to_drop = [column for column in upper.columns if any(upper[column] > 0.90)]
10 print(f"Columns with high correlation values : \n{to_drop} \n\nTotal Correlated columns : {len(to_drop)}")

```

Columns with high correlation values :  
['f\_6', 'f\_13', 'f\_16', 'f\_18', 'f\_20', 'f\_21', 'f\_24', 'f\_25', 'f\_34', 'f\_35', 'f\_36', 'f\_40', 'f\_43', 'f\_44', 'f\_49']

Total Correlated columns : 15

```

1 # Drop features/columns
2 df1 = df.copy()
3 df1.drop(to_drop, axis=1, inplace=True)
4 # dropping F_23 since it only has single value
5 f23 = "f_23"
6 df1.drop(f23, axis=1, inplace=True)
7 cleaned_df = df1.copy()
8 print("Original Dataframe:", df.shape)
9 print("Cleaned Dataframe (Highly correlated columns removed):", cleaned_df.shape)
10 print("\nRemoved Columns from dataset:\n", to_drop + [f23])

```

Original Dataframe: (937, 50)  
Cleaned Dataframe (Highly correlated columns removed): (937, 34)

Removed Columns from dataset:  
['f\_6', 'f\_13', 'f\_16', 'f\_18', 'f\_20', 'f\_21', 'f\_24', 'f\_25', 'f\_34', 'f\_35', 'f\_36', 'f\_40', 'f\_43', 'f\_44', 'f\_49', 'f\_23']

**Q4.** Apply various Machine Learning techniques to predict the output in the **target** column, make use of Bagging and Ensemble as required, and find the best model by evaluating the model using Model evaluation techniques.

### Code Snippet:

### Dependant (y) & Independent (x) Features

#### 1. Dropping dependant feature from dataset

```

1 x = cleaned_df.drop("target", axis=1)
2 y = cleaned_df["target"]
3
4 print(type(x))
5 print(type(y))
6 print(x.shape)
7 print(y.shape)

```

<class 'pandas.core.frame.DataFrame'>  
<class 'pandas.core.series.Series'>  
(937, 33)  
(937,)

#### 2. Splitting The dataset

```

1 # splitting the dataset into 70% training data and 30% test data
2 X_train, X_test, y_train, y_test = train_test_split(
3     x, y, test_size=0.3, random_state=42
4 )
5 print(f"Split Check Test values : {937 * 0.3} & Train values : {937 * 0.7}")
6 # rows , columns
7 print(X_train.shape)
8 print(X_test.shape)
9 print(y_train.shape)
10 print(y_test.shape)

```

Split Check Test values : 281.09999999999997 & Train values : 655.9  
(655, 33)  
(282, 33)  
(655,)  
(282,)

### 3. Display Splitted data

```

1 X_train, X_test

```

(655 rows x 33 columns),

	f_1	f_2	f_3	f_4	f_5	f_7	f_8	f_9	f_10	f_11	...
321	29	105	881.92	1128.79	83	38.90	8.51	2710.0	0.22	96.9	...
70	60	111	1153.32	1283.44	41	41.25	5.98	1760.0	0.14	157.7	...
209	17	867	1059.49	581.31	46	31.08	8.26	15780.0	0.27	137.4	...
656	9	85	71.06	469.47	140	70.85	11.28	4626.0	0.16	148.8	...
685	38	15	32.47	582.13	156	73.27	12.11	1880.0	0.17	112.5	...
...	...	...	...	...	...	...	...	...	...	...	...
430	183	51	1340.16	898.61	64	42.45	7.88	1430.0	0.19	89.2	...
292	317	117	1269.88	917.89	123	29.16	8.85	2440.0	0.30	119.9	...
412	151	64	991.70	1018.53	175	37.52	9.27	1400.0	0.25	114.3	...
557	63	59	1253.20	1192.53	76	29.51	7.32	1664.0	0.25	49.9	...
133	123	72	1606.14	1110.06	99	36.50	6.89	1760.0	0.19	102.3	...

(282 rows x 33 columns)

```

1 y_train, y_test

```

(757, 0)  
(693, 0)  
(854, 0)  
(501, 0)  
(664, 1)  
...  
(106, 0)  
(270, 0)  
(860, 0)  
(435, 0)  
(102, 0)

Name: target, Length: 655, dtype: int64

4. Standardizing the dataset

```

1 Click here to ask Blackbox to help you code faster
2 sc = StandardScaler()
3 sc.fit_transform(X_train)
4 sc.transform(X_test)

```

array([[ -0.84347997, -0.12493331, 0.30561268, ..., -0.38616422,
 -0.69720058, 0.41600293],
 [-0.36292032, -0.12189254, 0.75321954, ..., -0.38616422,
 -0.30840223, 0.43029588],
 [-1.02949997, 0.26124516, 0.59847027, ..., -0.38616422,
 -0.34532796, 0.39408706],
 ...,
 [ 1.0477233 , -0.14571195, 0.48666751, ..., -0.38616422,
 -0.77709157, 0.41314433],
 [-0.31642332, -0.14824593, 0.91794678, ..., -0.38616422,
 2.82886418, -2.35873678],
 [ 0.61367665, -0.14165758, 1.50003362, ..., -0.38616422,
 -0.53625066, 0.42839016]])

Sumit S. C

## Functions to – get the model scores, model parameters & train it, plot the graphs

```

1 # Function to evaluate and store results in a dictionary
2 def calculate_scores(model, X_train, y_train, X_test, y_test):
3     train_score = accuracy_score(
4         y_train, model.predict(X_train)
5     ) # Calculate train score
6     test_score = accuracy_score(
7         y_test, model.predict(X_test)) # Calculate test score
8
9     return train_score, test_score
10 # Function to evaluate the model to do the training on split data and calculate various params
11 def evaluate_model(model, model_name, X_test, y_test):
12     y_pred = model.predict(X_test)
13     y_prob = model.predict_proba(X_test)[ :, 1]
14
15     # Mean Squared Error and R-squared Score
16     mse = mean_squared_error(y_test, y_pred)
17     r2 = r2_score(y_test, y_pred)
18
19     # Confusion Matrix and classification report
20     cm = confusion_matrix(y_test, y_pred)
21     acc = accuracy_score(y_test, y_pred)
22     auc_score = roc_auc_score(y_test, y_prob)
23     fpr, tpr, _ = roc_curve(y_test, y_prob)
24     auc_value = auc(fpr, tpr)
25
26     cls_report = classification_report(y_test, y_pred, zero_division=0)
27
28     # Display results in tabular format
29     results_table = [
30         ["Model", model_name],
31         ["Mean Squared Error", mse],
32         ["R-squared Score", r2],
33         ["Confusion Matrix", f"({cm})"],
34         ["True Positive", cm[0, 0]],
35         ["False Negative", cm[0, 1]],
36         ["False Positive", cm[1, 0]],
37         ["True Negative", cm[1, 1]],
38         ["Accuracy", acc],
39         ["AUC", auc_score],
40         ["Train Score", train_score],
41         ["Test Score", test_score],
42     ]
43
44     print(tabulate(results_table, headers=[
45         "Metric", "Value"], tablefmt="heavy_grid"))
46
47     # Display Classification Report
48     print("\nClassification Report:\n")
49     print(cls_report)
50
51     # Plot Confusion Matrix
52     plt.matshow(cm, cmap=plt.cm.Blues)
53     plt.title(f"Confusion Matrix for {model_name}")
54     plt.colorbar()
55     plt.xlabel("Predicted")
56     plt.ylabel("True")
57
58     # Add annotations to matrix
59     for i in range(cm.shape[0]):
60         for j in range(cm.shape[1]):
61             plt.text(j, i, str(cm[i, j]), ha="center",
62                     va="center", color="black")
63     plt.show()
64
65     # Store results in the dictionary
66     return {
67         "Model": model_name,
68         "Mean Squared Error": mse,
69         "R-squared Score": r2,
70         "True Positive": cm[0, 0],
71         "False Negative": cm[0, 1],
72         "False Positive": cm[1, 0],
73         "True Negative": cm[1, 1],
74         "Accuracy": acc,
75         "AUC": auc_score,
76         "ROC Curve FPR": fpr,
77         "ROC Curve TPR": tpr,
78         "AUC Value": auc_value,
79         "Confusion Matrix": cm,
80         "Train Score": (train_score),
81         "Test Score": (test_score),
82     }

```

2. Plot for the graph of roc curve

```

1 Click here to ask Blackbox to help you code faster
2 # Function to plot ROC curve
3 def plot_roc_curve(model, X_test, y_test):
4     y_prob = model.predict_proba(X_test)[ :, 1] # "proba": Unknown word.
5     fpr, tpr, _ = roc_curve(y_test, y_prob)
6     auc_value = auc(fpr, tpr)
7
8     # Plot ROC curve
9     plt.plot(fpr, tpr, color="orange",
10             label=f"ROC Curve (AUC = {auc_value:.4f})")
11
12     plt.plot([0, 1], [0, 1], label="TPR=FPR", linestyle="--") # "linestyle": Unknown word.
13     plt.title(f"ROC Curve for {model_name}")
14     plt.xlabel("False Positive Rate (FPR)") # "xlabel": Unknown word.
15     plt.ylabel("True Positive Rate (TPR)") # "ylabel": Unknown word.
16     plt.grid()
17     plt.legend()
18     plt.show()

```

Model To look into

List of models to evaluate (just an example of parameters)

Models

```

models = [
    ("Logistic Regression", LogisticRegression(max_iter=1000, C=1.0, solver='lbfgs')),
    ("K-Nearest Neighbors", KNeighborsClassifier(n_neighbors=5, weights='uniform')),
    ("Decision Tree", DecisionTreeClassifier(max_depth=None, min_samples_split=2, min_samples_leaf=1)),
    ("Random Forest", RandomForestClassifier(n_estimators=100, max_depth=None, min_samples_split=2, min_samples_leaf=1)),
    ("AdaBoost", AdaBoostClassifier(n_estimators=50, learning_rate=1.0)),
    ("Bagging", BaggingClassifier(n_estimators=10, max_samples=1.0, max_features=1.0)),
    ("Gradient Boosting", GradientBoostingClassifier(n_estimators=100, learning_rate=0.1, max_depth=3)),
    ("Gaussian Naive Bayes", GaussianNB()),
    ("SVM", SVC(probability=True, C=1.0, kernel='rbf'))
]

```

3. Store the results of models built (by calling above evaluation code)

```

1 # Placeholder for results
2 evaluation_results = []

```

Sumit S. C

```

1 # 1. Passing the model name
2 model_name = "Logistic Regression"
3
4 # 2. model parameters
5 model = LogisticRegression(max_iter=1000, C=1.0, solver="lbfgs")
6
7 # 3.1: Fit the Model
8 model.fit(X_train, y_train)
9 print("\n", model, "\n")
10
11 # 3.2: Score Calculation
12 train_score, test_score = calculate_scores(
13     model, X_train, y_train, X_test, y_test)
14
15 # 3.3: Evaluate and Store Results
16 results = evaluate_model(model, model_name, X_test, y_test)
17 evaluation_results.append(results)
18
19 # Section 3.3: Plot ROC curve
20 print("==" * 75)
21 plot_roc_curve(model, X_test, y_test)
22 print("-" * 75)
23
24 # Model Detail
25 model

```

7 - in the above code need to just change the name of model and the model params and it will run the defined function and give the o/p so will not paste code for every (9) models

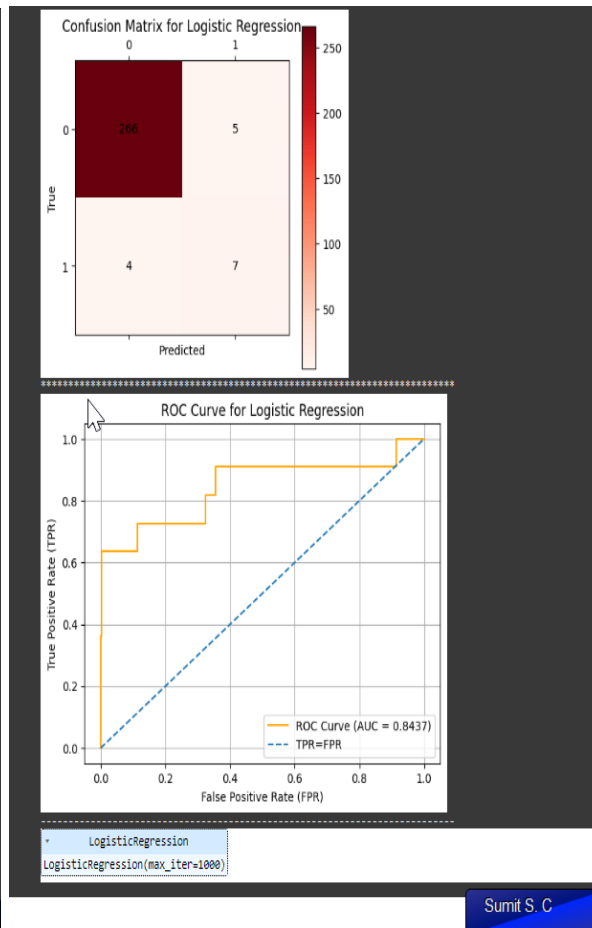
LogisticRegression(max\_iter=1000)

Metric	Value
Model	Logistic Regression
Mean Squared Error	0.031914893617021274
R-squared Score	0.14860784971486074
Confusion Matrix	[[266 5] [ 4 7]]
True Positive	266
False Negative	5
False Positive	4
True Negative	7
Accuracy	0.9680851063829787
AUC	0.8436766185843676
Train Score	0.9679389312977099
Test Score	0.9680851063829787

Classification Report:

	precision	recall	f1-score	support
0	0.99	0.98	0.98	271
1	0.58	0.64	0.61	11
accuracy			0.97	282
macro avg	0.78	0.81	0.80	282
weighted avg	0.97	0.97	0.97	282

Sumit S. C



Different Models Testing, evaluation & result/graph

1. Logistic Regression  
1 cell hidden
2. KNeighbors Classifier  
1 cell hidden
3. Decision Tree Classifier  
1 cell hidden
4. Random Forest Classifier  
1 cell hidden
5. Ada Boost Classifier  
1 cell hidden
6. Bagging Classifier  
1 cell hidden
7. Gradient Boosting Classifier  
1 cell hidden
8. Gaussian Naive Bayes  
1 cell hidden
9. SVM  
1 cell hidden

Results Check

```

1 print(f"Total models used & Evaluated :")
2 print(f"len(evaluation_results) &")
3 print(f"Saved parameters in each results :{len(results)}")
4 for result in evaluation_results:
5     print("\n", result)

```

Total models used & Evaluated : 9 &  
Saved parameters in each results :15

Sumit S. C

8 - The score calculation & tabulated results of various model parameters after training & the ROC graph and confusion matrix and model name (Will not paste the code or results of every model as its time consuming so can look into colab link for more) & the various models trained on list and results - 9 models





- Q5) Save the best model and Load the model

- ✓ Saving Model as pickle file and dumping it to use later on

- ✓ Loading the saved model

Name of loaded Model : Random Forest

Best Model Name: Random Forest

Retrieved Model Instance: RandomForestClassifier()

```
▶ # Testing the imported model
```

```
print("Length of test data: ", len(load_model.predict(X_test)))
load_model.predict(X_test)
```

Length of test data: 282

Best Selected Model name : 'Random Forest' &

its parameters :

```
{'bootstrap': True, 'ccp_alpha': 0.0, 'class_weight': None, 'criterion': 'gini', 'max_depth': None, 'max_features': 'sqrt', 'max
```

Sumit S. C

Sumit S. C

**Q6.** Take the original data set and make another dataset by randomly picking 20 data points from the oil spill dataset and applying the saved model to the same.

### Code Snippet:

Generating sample data from cleaned df to test on the trained model.

```
[ ] random_datasample = cleaned_df.sample(20)
random_datasample_df = random_datasample.drop("target", axis=1)
print(random_datasample_df.shape)
random_datasample_df.head()
```

Resetting the index as the randomly generated data has no continuous index (will delete later just for understanding)

```
[ ] random_datasample_df.reset_index()
```

Show hidden output

Saving the random sample dataset and removing the index

```
[ ] random_datasample_df.to_csv("20_random_sample.csv", index=False)
```

Loading the sample data and checking basics

```
[ ] testsample_df = pd.read_csv("20_random_sample.csv")
print("Shape of loaded sample dataframe:", testsample_df.shape, "\n\nSample Dataframe contents")
testsample_df
```

Show hidden output

```
[ ] # making prediction on random data
predicted_data = load_model.predict(testsample_df)
print("The predicted data from {f_modelname} model:\n", predicted_data)
```

The predicted data from Random Forest model:  
[0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0]

Comparison of Actual and Predicted values by the model

```
# Compare the actual data and predicted data
prediction_data = random_datasample.copy()
prediction_data["predicted_target"] = predicted_data

# Print the actual and predicted data
print(f"Actual Data and Predicted Data Comparison based on {f_modelname} model:\n")
# print(prediction_data[["target", "predicted_target"]])
comparision = {"Actual Target": random_datasample["target"], "Predicted Target": predicted_data}
final_results = pd.DataFrame(comparision)
final_results
```

Actual Data and Predicted Data Comparison based on Random Forest model:

Actual Target	Predicted Target
855	0
148	0
728	0
902	0
434	0
473	0
409	0
96	0
235	0
362	1
808	0
705	0
558	0

13 - Generating 20 sample data and then predicting it from the loaded model and checking results

Calculating the correctness of model

```
# Calculate the number of correct predictions
correct_predictions = (
    prediction_data["predicted_target"] == prediction_data["target"]).sum()
# Calculate the percentage of correct predictions
percentage_correct_predictions = (
    correct_predictions / len(prediction_data)) * 100
# Print the result
print(f"\nPercentage of Correct Predictions: {percentage_correct_predictions:.2f}%")
if (percentage_correct_predictions >= 90):
    print(f"\nOur model based on '{f_modelname}' is well trained having prediction accuracy of {percentage_correct_predictions:.2f}%")
else:
    print(f"\nOur model based on '{f_modelname}' needs to be trained more to achieve atleast 95% prediction accuracy from our current results : {percentage_correct_predictions:.2f}%")
```

Percentage of Correct Predictions: 100.00%

Our model based on 'Random Forest' is well trained having prediction accuracy of 100.00%

Saving the results in a output file

```
# Saving the final results in a output file
final_results.to_csv('final_results.csv', index=False)
with open('final_results.txt', 'w') as f:
    f.write(final_results.to_string())
with open('final_results.txt', 'a') as f:
    f.write(f"\n\n-----\n\nPrinting the results of our {f_modelname} prediction on random 20 data samples.")
    f.write(f"\nNumber of correct predictions: {sum(final_results['Actual Target'] == final_results['Predicted Target'])}")
    f.write(f"\nPercentage of correct predictions: {sum(final_results['Actual Target'] == final_results['Predicted Target']) / len(final_results)}")
# Print the result in output file
if (percentage_correct_predictions >= 90):
    f.write(f"\nOur model based on '{f_modelname}' is well trained having prediction accuracy of {percentage_correct_predictions:.2f}%")
else:
    f.write(f"\nOur model based on '{f_modelname}' needs to be trained more to achieve atleast 95% prediction accuracy from our current results : {percentage_correct_predictions:.2f}%")
```

14 - Check the results of actual target value & predicted ( Extra step saving the results in csv and inside a txt file with the last output results)





---

***Thank You***

---