# PHASE-3 PREPROCESSING

## 3.1 DATASET AND IT'S DETAIL EXPLAINATION IMPLEMENTATION

*The Dataset is taken from the kaggle , where Kaggle is a data science competition platform and online community of data scientists and machine learning practitioners under Google LLC. Kaggle enables users to find and publish datasets, explore and build models in a web-based data science environment, work with other data scientists and machine learning engineers, and enter competitions to solve data science challenges.*

*My Dataset Link: [https://www.kaggle.com/datasets/gpreda/covid-world-vaccination-progress](https://www.kaggle.com/datasets/gpreda/covid-world-vaccination-progress)*

*My Dataset about Covid-19 Vaccines Analysis. The Dataset contains many columns they are,*

- ***Country-*** *this is the country for which the vaccination information is provided;*
- ***Country ISO Code -*** *ISO code for the country;*
- ***Date -*** *date for the data entry; for some of the dates we have only the daily vaccinations, for others, only the (cumulative) total;*
- ***Total number of vaccinations -*** *this is the absolute number of total immunizations in the country;*
- ***Total number of people vaccinated -*** *a person, depending on the immunization scheme, will receive one or more (typically 2) vaccines; at a certain moment, the number of vaccination might be larger than the number of people;*
- ***Total number of people fully vaccinated -*** *this is the number of people that received the entire set of immunization according to the immunization scheme (typically 2); at a certain moment in time, there might be a certain number of people that received one vaccine and another number (smaller) of people that received all vaccines in the scheme;*
- ***Daily vaccinations (raw) -*** *for a certain data entry, the number of vaccination for that date/country;*
- ***Daily vaccinations -*** *for a certain data entry, the number of vaccination for that date/country;*
- ***Total vaccinations per hundred -*** *ratio (in percent) between vaccination number and total population up to the date in the country;*
- ***Total number of people vaccinated per hundred -*** *ratio (in percent) between population immunized and total population up to the date in the country;*

- ***Total number of people fully vaccinated per hundred*** - *ratio (in percent) between population fully immunized and total population up to the date in the country;*
- ***Number of vaccinations per day*** - *number of daily vaccination for that day and country;*
- ***Daily vaccinations per million*** - *ratio (in ppm) between vaccination number and total population for the current date in the country;*
- ***Vaccines used in the country*** - *total number of vaccines used in the country (up to date);*
- ***Source name*** - *source of the information (national authority, international organization, local organization etc.);*
- ***Source website*** - *website of the source of information;*

## 3.2 BEGIN BUILDING THE PROJECT BY LOAD THE DATASET

*To import the required libraries and read a CSV file,*

*data=pd.read_csv('D:\os\country_vaccinations.csv')*

*Print(data)*



## 3.3 PREPROCESS DATASET

**Import The Required Libraries:**

import numpy as np

import pandas as pd

import matplotlib.pyplot as plt

import plotly.express as px

*import plotly.graph_objects as go*

*import matplotlib.patches as mpatches*

*from sklearn.model_selection import train_test_split*

*from sklearn.impute import SimpleImputer*

*from plotly.subplots import make_subplots*

*from sklearn.preprocessing import StandardScaler*

*import seaborn as sns*

*import plotly.figure_factory as ff*

*from plotly.colors import n_colors*

**Importing the Dataset:**

- *Read Dataset,*

*data=pd.read_csv('D:\os\country_vaccinations.csv')*

*print(data)*

- *Create Matrix,*

*matrix=data.values*

*print(matrix)*



- *Other Imformation about dataset,*

## *Handling The Missing Data:*

 *Before Handling the Missing data, we use isnull () to show the null values and using isnull().sum to get total number of null values.*

*The below code for all the data cleaning,*



```python
In [18]: data.fillna(value=0,inplace=True)
data.total_vaccinations=data.total_vaccinations.astype(int)
data.people_vaccinated=data.people_vaccinated.astype(int)
data.people_fully_vaccinated=data.people_fully_vaccinated.astype(int)
data.daily_vaccinations_raw=data.daily_vaccinations_raw.astype(int)
data.daily_vaccinations=data.daily_vaccinations.astype(int)
data.total_vaccinations_per_hundred=data.total_vaccinations_per_hundred.astype(int)
data.people_fully_vaccinated_per_hundred=data.people_fully_vaccinated_per_hundred.astype(int)
print(data)

         country iso_code        date  total_vaccinations  \
0    Afghanistan      AFG  2021-02-22                   0
1    Afghanistan      AFG  2021-02-23                   0
2    Afghanistan      AFG  2021-02-24                   0
3    Afghanistan      AFG  2021-02-25                   0
4    Afghanistan      AFG  2021-02-26                   0
...          ...      ...         ...                 ...
86507    Zimbabwe      ZWE  2022-03-25             8691642
86508    Zimbabwe      ZWE  2022-03-26             8791728
86509    Zimbabwe      ZWE  2022-03-27             8845039
86510    Zimbabwe      ZWE  2022-03-28             8934360
86511    Zimbabwe      ZWE  2022-03-29             9039729
```



```python
         [86512 rows x 15 columns]

In [19]: date=data.date.str.split('-',expand=True)
print(date)

          0   1   2
0      2021  02  22
1      2021  02  23
2      2021  02  24
3      2021  02  25
4      2021  02  26
...     ...  ..  ..
86507  2022  03  25
86508  2022  03  26
86509  2022  03  27
86510  2022  03  28
86511  2022  03  29

[86512 rows x 3 columns]

In [ ]:
```



```python
In [22]: data['year']=date[0]
data['month']=date[1]
data['day']=date[2]
data.year=pd.to_numeric(data.year)
data.month=pd.to_numeric(data.month)
data.day=pd.to_numeric(data.day)
data.date=pd.to_datetime(data.date)
data.head()
```

Out[22]:

| | country | iso_code | date | total_vaccinations | people_vaccinated | people_fully_vaccinated | daily_vaccinations_raw | daily_vaccinations | total_vaccinations_per... |
|---|---|---|---|---|---|---|---|---|---|
| 0 | Afghanistan | AFG | 2021-02-22 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | Afghanistan | AFG | 2021-02-23 | 0 | 0 | 0 | 0 | 1367 | |
| 2 | Afghanistan | AFG | 2021-02-24 | 0 | 0 | 0 | 0 | 1367 | |
| 3 | Afghanistan | AFG | 2021-02-25 | 0 | 0 | 0 | 0 | 1367 | |

Some detailed features to specify the details using the below code,



```python
print('data point starts from',data.date.min())
print('data point ends at',data.date.max())
print('total number of countries in the data set',len(data.country.unique()))
print('total number of unique vaccines in the data set',len(data.vaccines.unique()))
```

```
data point starts from 2020-12-02 00:00:00
data point ends at 2022-03-29 00:00:00
total number of countries in the data set 223
total number of unique vaccines in the data set 84
```



```python
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 86512 entries, 0 to 86511
Data columns (total 18 columns):
 #   Column                               Non-Null Count  Dtype
---  ------                               --------------  -----
 0   country                              86512 non-null  object
 1   iso_code                             86512 non-null  object
 2   date                                 86512 non-null  datetime64[ns]
 3   total_vaccinations                   86512 non-null  int32
 4   people_vaccinated                    86512 non-null  int32
 5   people_fully_vaccinated              86512 non-null  int32
 6   daily_vaccinations_raw               86512 non-null  int32
 7   daily_vaccinations                   86512 non-null  int32
 8   total_vaccinations_per_hundred       86512 non-null  int32
 9   people_vaccinated_per_hundred        86512 non-null  float64
 10  people_fully_vaccinated_per_hundred  86512 non-null  int32
 11  daily_vaccinations_per_million       86512 non-null  float64
 12  vaccines                             86512 non-null  object
 13  source_name                          86512 non-null  object
 14  source_website                       86512 non-null  object
 15  year                                 86512 non-null  int64
```

*Using Data visualization we are going to draw some visuals to get insights from dataset,*

*Describe () function is used to get the statistics of each feature in dataset to get count, min, max, standard deviation, median, etc.,*



*Unique () function helps to get unique values,*

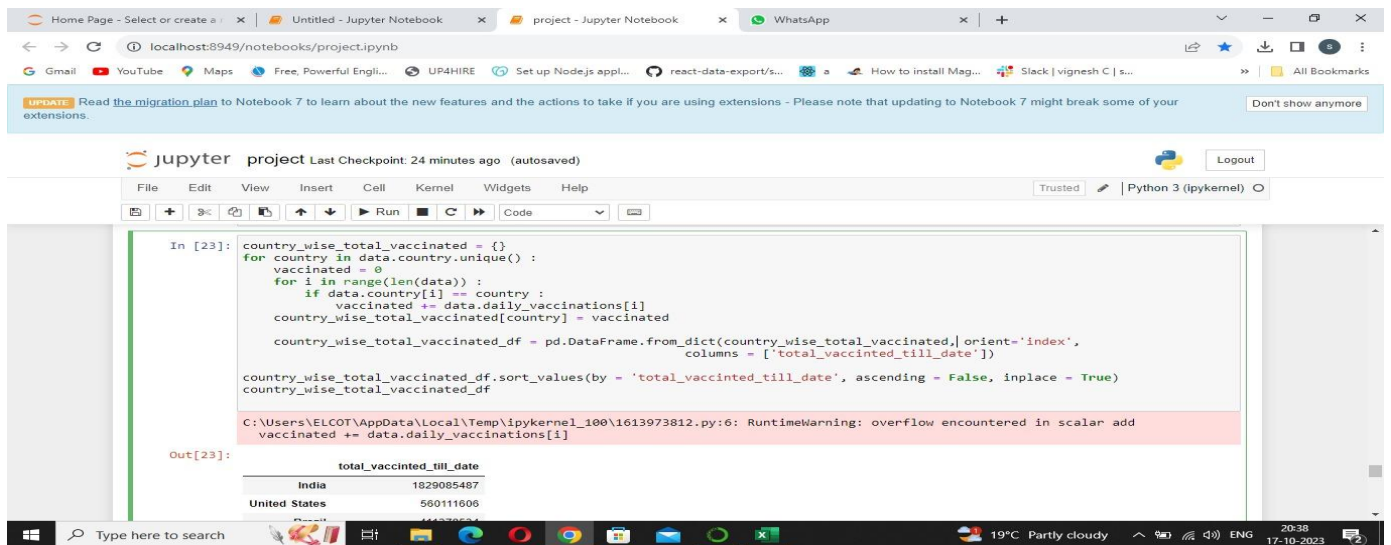*To see how many total vaccines have been used in each country using the code below,*
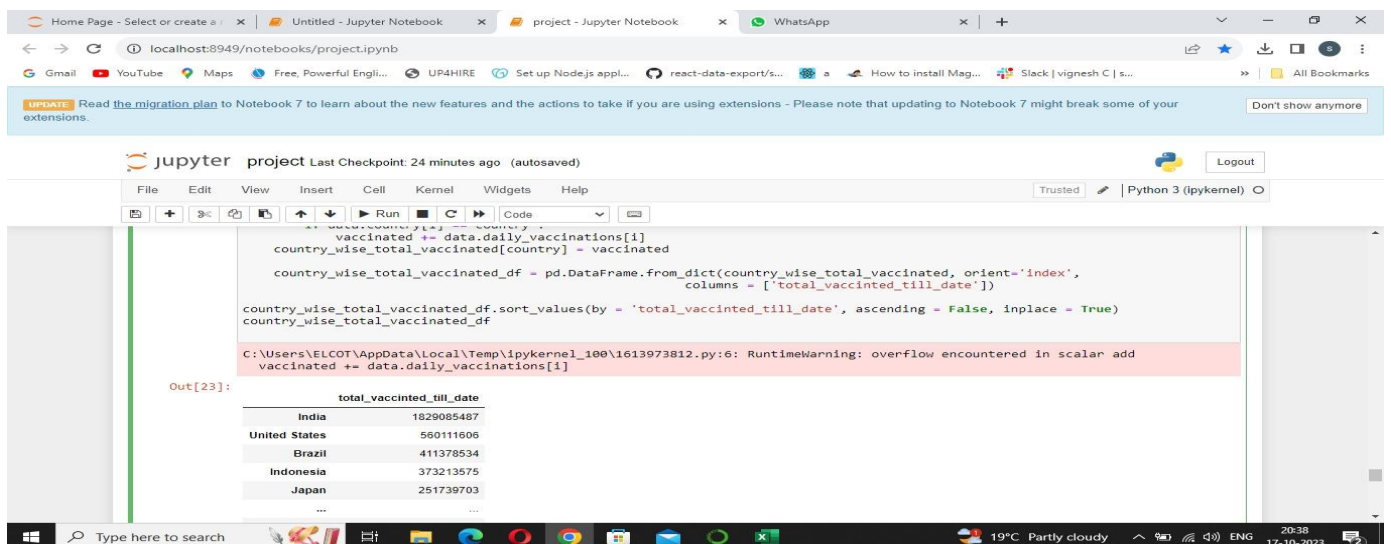


```
In [23]: country_wise_total_vaccinated = {}
         for country in data.country.unique() :
             vaccinated = 0
             for i in range(len(data)) :
                 if data.country[i] == country :
                     vaccinated += data.daily_vaccinations[i]
             country_wise_total_vaccinated[country] = vaccinated

             country_wise_total_vaccinated_df = pd.DataFrame.from_dict(country_wise_total_vaccinated, orient='index',
                                                          columns = ['total_vaccinted_till_date'])

         country_wise_total_vaccinated_df.sort_values(by = 'total_vaccinted_till_date', ascending = False, inplace = True)
         country_wise_total_vaccinated_df
```
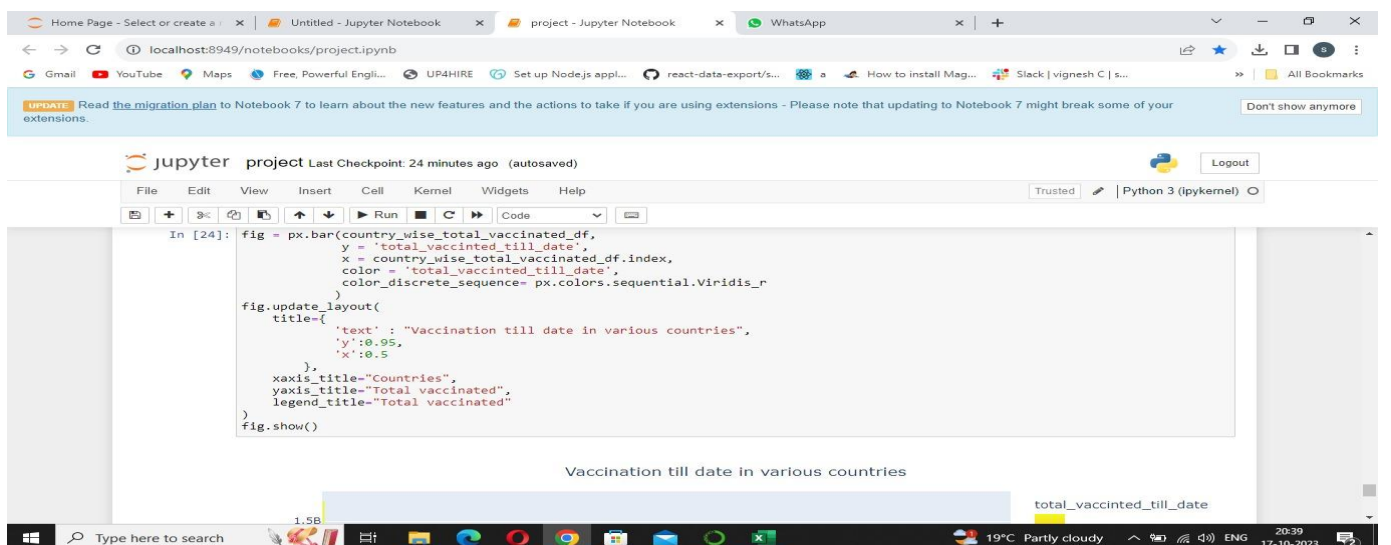
C:\Users\ELCOT\AppData\Local\Temp\ipykernel_100\1613973812.py:6: RuntimeWarning: overflow encountered in scalar add
  vaccinated += data.daily_vaccinations[i]

Out[23]:

| | total_vaccinted_till_date |
|---|---|
| India | 1829085487 |
| United States | 560111606 |



Out[23]:

| | total_vaccinted_till_date |
|---|---|
| India | 1829085487 |
| United States | 560111606 |
| Brazil | 411378534 |
| Indonesia | 373213575 |
| Japan | 251739703 |
| ... | ... |



```
In [24]: fig = px.bar(country_wise_total_vaccinated_df,
                      y = 'total_vaccinted_till_date',
                      x = country_wise_total_vaccinated_df.index,
                      color = 'total_vaccinted_till_date',
                      color_discrete_sequence= px.colors.sequential.Viridis_r
                      )
         fig.update_layout(
             title={
                 'text' : "Vaccination till date in various countries",
                 'y':0.95,
                 'x':0.5
             },
             xaxis_title="Countries",
             yaxis_title="Total vaccinated",
             legend_title="Total vaccinated"
         )
         fig.show()
```

Vaccination till date in various countries

*To draw a line plot where x-axis is Date and the y-axis is daily_vaccination using the in the image,*



```python
fig = px.line(data, x = 'date', y ='daily_vaccinations', color = 'country')
fig.update_layout(
    title={
        'text' : "Daily vaccination trend",
        'y':0.95,
        'x':0.5
    },
    xaxis_title="Date",
    yaxis_title="Daily Vaccinations"
)
fig.show()
```

*Now, using the sklearn.preprocessing library contains class called imputer, helps in missing data by using the below:*



## Encoding categorical data(one-hot encoding)
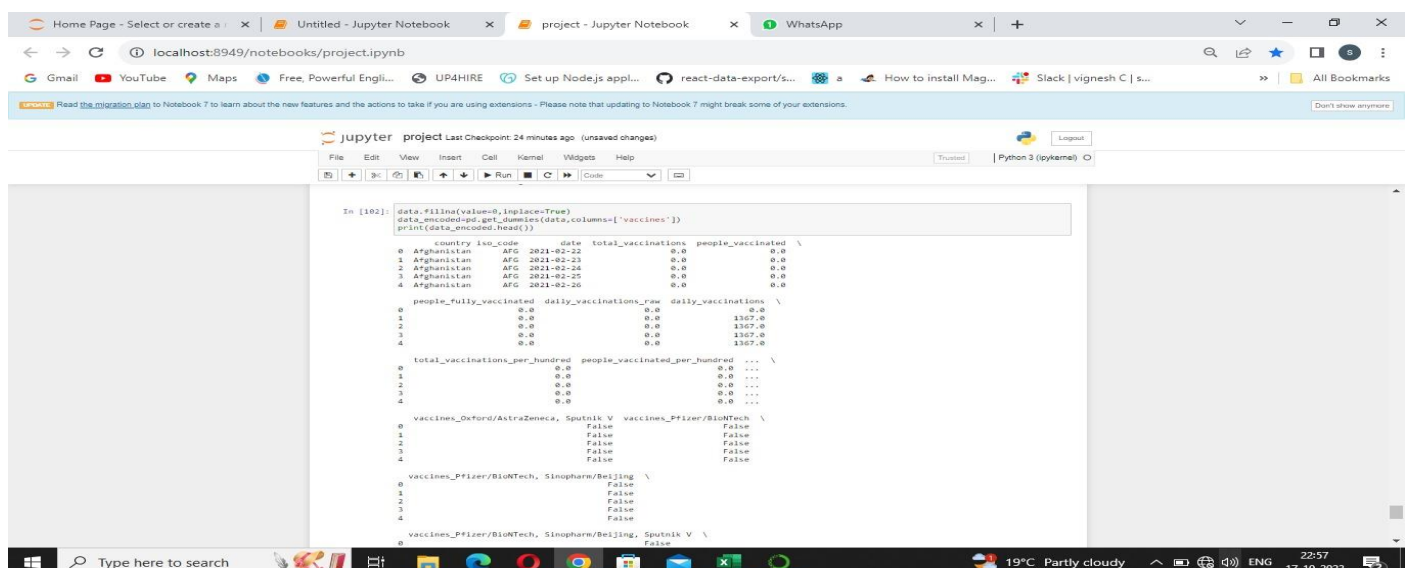
*One-hot encoding is a technique used to convert categorical data into a numerical format that machine learning algorithms can work with. Here's how you can perform one-hot encoding in Python, assuming you have a dataset with categorical variables:*

## Splitting the data set into test set and training set

By using,

Import train_test_split



## Feature Scaling

By using, import StandardScaler

## 3.4 PERFORMING DIFFERENT ANALYSIS

**1. Descriptive Analysis:**

```
import pandas as pd

data = pd.read_csv('your_dataset.csv')

print(data.describe())
```

**2. Time Series Analysis:**

```
import pandas as pd

data = pd.read_csv('time_series_data.csv')

# Perform time series analysis, e.g., using the statsmodels library
```

**3. Correlation Analysis:**

```
import pandas as pd

data = pd.read_csv('correlation_data.csv')

correlation_matrix = data.corr()

print(correlation_matrix)
```

**4. Regression Analysis:**

```
import pandas as pd

from sklearn.linear_model import LinearRegression

data = pd.read_csv('regression_data.csv')

X = data[['independent_variable']]

y = data['dependent_variable']

model = LinearRegression()

model.fit(X, y)

# You can make predictions using the model
```
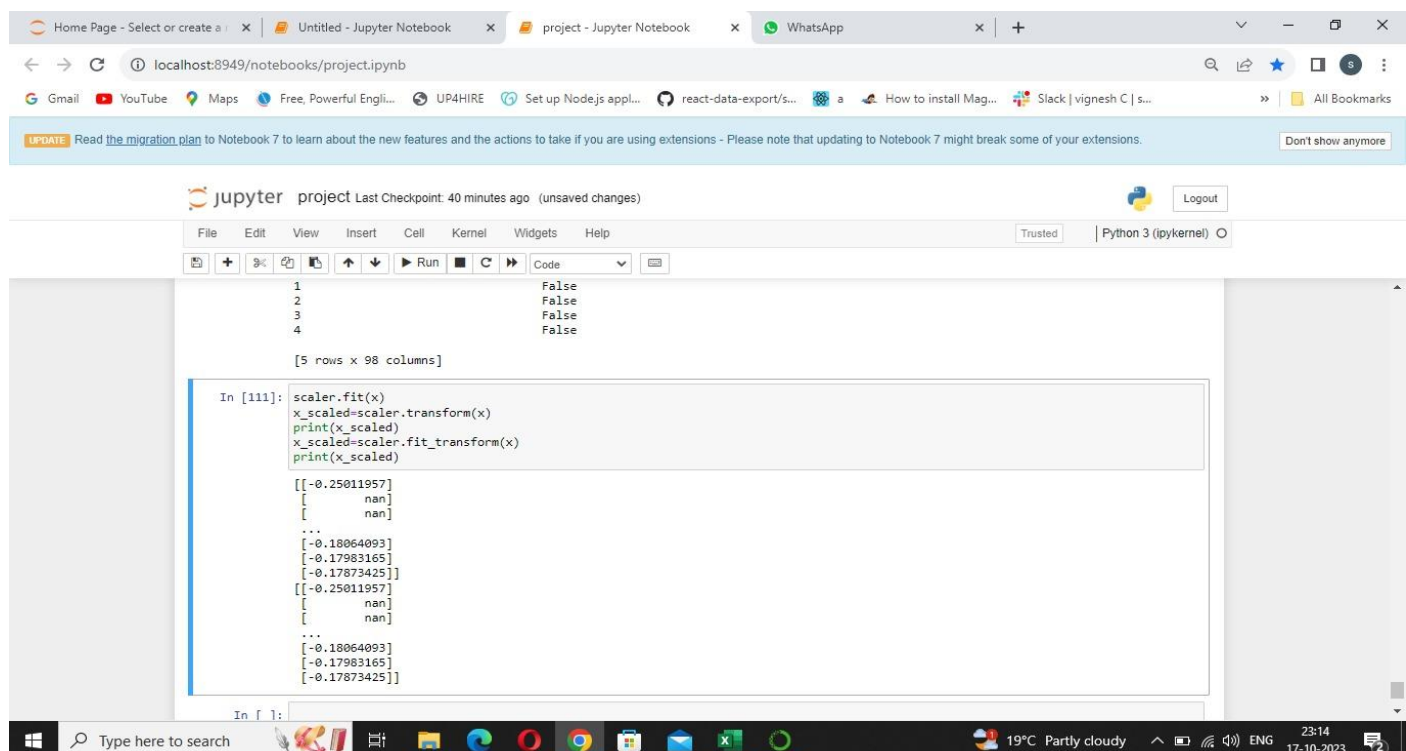
**5. Classification Analysis:**

```
import pandas as pd
```

```python
from sklearn.model_selection import train_test_split

from sklearn.ensemble import RandomForestClassifier

data = pd.read_csv('classification_data.csv')

X = data.drop('target_class', axis=1)

y = data['target_class']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)

model = RandomForestClassifier()

model.fit(X_train, y_train)

# Evaluate the model's performance
```

### 6. Text Analysis:

```python
from textblob import TextBlob

text = "Your text data here."

blob = TextBlob(text)

# Perform sentiment analysis, extract keywords, etc.
```

### 7. Cluster Analysis:

```python
from sklearn.cluster import KMeans

data = pd.read_csv('cluster_data.csv')

model = KMeans(n_clusters=3)

model.fit(data)

# Explore cluster assignments and centroids
```

### 8. Financial Ratio Analysis:

```python
import pandas as pd

data = pd.read_csv('financial_data.csv')

# Calculate financial ratios, e.g., profit margin, debt-to-equity ratio, etc.
```

### 9. Statistical Hypothesis Testing:

```python
import scipy.stats as stats

group1 = [1, 2, 3, 4, 5]

group2 = [6, 7, 8, 9, 10]

t_stat, p_value = stats.ttest_ind(group1, group2)

# Check if the p-value is significant
```