# Unemployment analysis with python

Unemployment is measured by the unemployment rate which is the number of people who are unemployed as a percentage of the total labour force. We have seen a sharp increase in the unemployment rate during Covid-19,So analysing the unemployment rate can be a good data science project.

## Steps

**1. Importing requried libraries

2. Loading of dataset

3. Display Summary Statistics

4. Data Visualization

Used plots are pair plots, histograms, scatter plots, line plots 5. Correlation matrix 6. Analysis with Sunburst Plot

Python Plotly Library is an open-source library that can be used for data visualization and understanding data simply and easily. Plotly supports various types of plots like line charts, scatter plots, histograms, cox plots, etc.

In [60]:
```
!pip install plotly.express
```

```
Collecting plotly.express
  Downloading plotly_express-0.4.1-py2.py3-none-any.whl (2.9 kB)
Requirement already satisfied: scipy>=0.18 in c:\users\meghana\anaconda3\lib\site-
packages (from plotly.express) (1.7.1)
Requirement already satisfied: patsy>=0.5 in c:\users\meghana\anaconda3\lib\site-p
ackages (from plotly.express) (0.5.2)
Requirement already satisfied: statsmodels>=0.9.0 in c:\users\meghana\anaconda3\li
b\site-packages (from plotly.express) (0.12.2)
Requirement already satisfied: pandas>=0.20.0 in c:\users\meghana\anaconda3\lib\si
te-packages (from plotly.express) (1.3.4)
Requirement already satisfied: numpy>=1.11 in c:\users\meghana\anaconda3\lib\site-
packages (from plotly.express) (1.20.3)
Requirement already satisfied: plotly>=4.1.0 in c:\users\meghana\anaconda3\lib\sit
e-packages (from plotly.express) (5.12.0)
Requirement already satisfied: python-dateutil>=2.7.3 in c:\users\meghana\anaconda
3\lib\site-packages (from pandas>=0.20.0->plotly.express) (2.8.2)
Requirement already satisfied: pytz>=2017.3 in c:\users\meghana\anaconda3\lib\site
-packages (from pandas>=0.20.0->plotly.express) (2021.3)
Requirement already satisfied: six in c:\users\meghana\anaconda3\lib\site-packages
(from patsy>=0.5->plotly.express) (1.16.0)
Requirement already satisfied: tenacity>=6.2.0 in c:\users\meghana\anaconda3\lib\s
ite-packages (from plotly>=4.1.0->plotly.express) (8.1.0)
Installing collected packages: plotly.express
Successfully installed plotly.express-0.4.1
```

In [1]:
```
#importing libraries
```

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

In [2]:
```
#Loading Unemployment in India dataset
df = pd.read_csv('Unemployment in India.csv')
df
```

Out[2]:

| | Region | Date | Frequency | Estimated Unemployment Rate (%) | Estimated Employed | Estimated Labour Participation Rate (%) | Area |
|---|---|---|---|---|---|---|---|
| 0 | Andhra Pradesh | 31-05-2019 | Monthly | 3.65 | 11999139.0 | 43.24 | Rural |
| 1 | Andhra Pradesh | 30-06-2019 | Monthly | 3.05 | 11755881.0 | 42.05 | Rural |
| 2 | Andhra Pradesh | 31-07-2019 | Monthly | 3.75 | 12086707.0 | 43.50 | Rural |
| 3 | Andhra Pradesh | 31-08-2019 | Monthly | 3.32 | 12285693.0 | 43.97 | Rural |
| 4 | Andhra Pradesh | 30-09-2019 | Monthly | 5.17 | 12256762.0 | 44.68 | Rural |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 763 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| 764 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| 765 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| 766 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| 767 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |

768 rows × 7 columns

In [3]:
```
#Loading Unemployment rate upto 2020 dataset
df1 = pd.read_csv('Unemployment_Rate_upto_11_2020.csv')
df1
```

Out[3]:

| | Region | Date | Frequency | Estimated Unemployment Rate (%) | Estimated Employed | Estimated Labour Participation Rate (%) | Region.1 | long |
|---|---|---|---|---|---|---|---|---|

| | Region | Date | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 0 | Andhra Pradesh | 31-01-2020 | M | 5.48 | 16635535 | 41.02 | South | 15 |
| 1 | Andhra Pradesh | 29-02-2020 | M | 5.83 | 16545652 | 40.90 | South | 15 |
| 2 | Andhra Pradesh | 31-03-2020 | M | 5.79 | 15881197 | 39.18 | South | 15 |
| 3 | Andhra Pradesh | 30-04-2020 | M | 20.51 | 11336911 | 33.10 | South | 15 |
| 4 | Andhra Pradesh | 31-05-2020 | M | 17.43 | 12988845 | 36.46 | South | 15 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 262 | West Bengal | 30-06-2020 | M | 7.29 | 30726310 | 40.39 | East | 22 |
| 263 | West Bengal | 31-07-2020 | M | 6.83 | 35372506 | 46.17 | East | 22 |
| 264 | West Bengal | 31-08-2020 | M | 14.87 | 33298644 | 47.48 | East | 22 |
| 265 | West Bengal | 30-09-2020 | M | 9.35 | 35707239 | 47.73 | East | 22 |
| 266 | West Bengal | 31-10-2020 | M | 9.98 | 33962549 | 45.63 | East | 22 |

267 rows × 9 columns

◄ ██████████████████████████ ►

In [4]:
```python
#First and last rows of Unemployment in India dataset
print("Rows from start are: ")
print(df.head(6))
print("\n")
print("Rows from bottom: ")
print(df.tail(8))
```

```
Rows from start are:
          Region        Date Frequency   Estimated Unemployment Rate (%)  \
0  Andhra Pradesh  31-05-2019   Monthly                              3.65
1  Andhra Pradesh  30-06-2019   Monthly                              3.05
2  Andhra Pradesh  31-07-2019   Monthly                              3.75
```

```
3  Andhra Pradesh  31-08-2019  Monthly                                    3.32
4  Andhra Pradesh  30-09-2019  Monthly                                    5.17
5  Andhra Pradesh  31-10-2019  Monthly                                    3.52
```

```
   Estimated Employed  Estimated Labour Participation Rate (%)   Area
0          11999139.0                                   43.24   Rural
1          11755881.0                                   42.05   Rural
2          12086707.0                                   43.50   Rural
3          12285693.0                                   43.97   Rural
4          12256762.0                                   44.68   Rural
5          12017412.0                                   43.01   Rural
```

```
Rows from bottom:
     Region  Date  Frequency   Estimated Unemployment Rate (%)  \
760     NaN   NaN        NaN                               NaN
761     NaN   NaN        NaN                               NaN
762     NaN   NaN        NaN                               NaN
763     NaN   NaN        NaN                               NaN
764     NaN   NaN        NaN                               NaN
765     NaN   NaN        NaN                               NaN
766     NaN   NaN        NaN                               NaN
767     NaN   NaN        NaN                               NaN
```

```
     Estimated Employed  Estimated Labour Participation Rate (%) Area
760                 NaN                                     NaN  NaN
761                 NaN                                     NaN  NaN
762                 NaN                                     NaN  NaN
763                 NaN                                     NaN  NaN
764                 NaN                                     NaN  NaN
765                 NaN                                     NaN  NaN
766                 NaN                                     NaN  NaN
767                 NaN                                     NaN  NaN
```

In [5]:
```python
#First and last rows in the dataset "umemployment rate till 2020"
print("Rows from start are: ")
print(df1.head(6))
print("\n")
print("Rows from bottom: ")
print(df1.tail(8))
```

```
Rows from start are:
          Region         Date  Frequency   Estimated Unemployment Rate (%)  \
0  Andhra Pradesh   31-01-2020          M                             5.48
1  Andhra Pradesh   29-02-2020          M                             5.83
2  Andhra Pradesh   31-03-2020          M                             5.79
3  Andhra Pradesh   30-04-2020          M                            20.51
4  Andhra Pradesh   31-05-2020          M                            17.43
5  Andhra Pradesh   30-06-2020          M                             3.31
```

```
   Estimated Employed  Estimated Labour Participation Rate (%) Region.1  \
0            16635535                                    41.02    South
1            16545652                                    40.90    South
2            15881197                                    39.18    South
3            11336911                                    33.10    South
4            12988845                                    36.46    South
5            19805400                                    47.41    South
```

```
   longitude  latitude
0    15.9129     79.74
```

```
1    15.9129    79.74
2    15.9129    79.74
3    15.9129    79.74
4    15.9129    79.74
5    15.9129    79.74
```

```
Rows from bottom:
          Region        Date  Frequency   Estimated Unemployment Rate (%)  \
259  West Bengal   31-03-2020          M                              6.92
260  West Bengal   30-04-2020          M                             17.41
261  West Bengal   31-05-2020          M                             17.41
262  West Bengal   30-06-2020          M                              7.29
263  West Bengal   31-07-2020          M                              6.83
264  West Bengal   31-08-2020          M                             14.87
265  West Bengal   30-09-2020          M                              9.35
266  West Bengal   31-10-2020          M                              9.98

     Estimated Employed   Estimated Labour Participation Rate (%) Region.1  \
259             35903917                                    47.27     East
260             26938836                                    39.90     East
261             28356675                                    41.92     East
262             30726310                                    40.39     East
263             35372506                                    46.17     East
264             33298644                                    47.48     East
265             35707239                                    47.73     East
266             33962549                                    45.63     East

     longitude  latitude
259    22.9868    87.855
260    22.9868    87.855
261    22.9868    87.855
262    22.9868    87.855
263    22.9868    87.855
264    22.9868    87.855
265    22.9868    87.855
266    22.9868    87.855
```

In [6]:
```python
#Summary statistics of Unemployment in India dataset
print("Shape of the data set ",df.shape)
print("Size of the data set",df.size)
print("\n")
print("Info of the dataset \n",df.info)
print("\n")
print("Descriptive statistics of the dataset \n",df.describe)
```

```
Shape of the data set  (768, 7)
Size of the data set 5376


Info of the dataset
 <bound method DataFrame.info of               Region       Date Frequency   Est
imated Unemployment Rate (%)  \
0     Andhra Pradesh   31-05-2019   Monthly                             3.65
1     Andhra Pradesh   30-06-2019   Monthly                             3.05
2     Andhra Pradesh   31-07-2019   Monthly                             3.75
3     Andhra Pradesh   31-08-2019   Monthly                             3.32
4     Andhra Pradesh   30-09-2019   Monthly                             5.17
..               ...          ...       ...                              ...
763              NaN          NaN       NaN                              NaN
```

|     | | | |     |
| --- | --- | --- | --- |
| 764 | NaN | NaN | NaN |     | NaN |
| 765 | NaN | NaN | NaN |     | NaN |
| 766 | NaN | NaN | NaN |     | NaN |
| 767 | NaN | NaN | NaN |     | NaN |

|     | Estimated Employed | Estimated Labour Participation Rate (%) | Area |
| --- | --- | --- | --- |
| 0 | 11999139.0 | 43.24 | Rural |
| 1 | 11755881.0 | 42.05 | Rural |
| 2 | 12086707.0 | 43.50 | Rural |
| 3 | 12285693.0 | 43.97 | Rural |
| 4 | 12256762.0 | 44.68 | Rural |
| .. | ... | ... | ... |
| 763 | NaN | NaN | NaN |
| 764 | NaN | NaN | NaN |
| 765 | NaN | NaN | NaN |
| 766 | NaN | NaN | NaN |
| 767 | NaN | NaN | NaN |

[768 rows x 7 columns]>


Descriptive statistics of the dataset
<bound method NDFrame.describe of

|     | Region | Date | Frequency | E |
| --- | --- | --- | --- | --- |
| stimated Unemployment Rate (%)  \ | | | | |
| 0 | Andhra Pradesh | 31-05-2019 | Monthly | 3.65 |
| 1 | Andhra Pradesh | 30-06-2019 | Monthly | 3.05 |
| 2 | Andhra Pradesh | 31-07-2019 | Monthly | 3.75 |
| 3 | Andhra Pradesh | 31-08-2019 | Monthly | 3.32 |
| 4 | Andhra Pradesh | 30-09-2019 | Monthly | 5.17 |
| .. | ... | ... | ... | ... |
| 763 | NaN | NaN | NaN | NaN |
| 764 | NaN | NaN | NaN | NaN |
| 765 | NaN | NaN | NaN | NaN |
| 766 | NaN | NaN | NaN | NaN |
| 767 | NaN | NaN | NaN | NaN |

|     | Estimated Employed | Estimated Labour Participation Rate (%) | Area |
| --- | --- | --- | --- |
| 0 | 11999139.0 | 43.24 | Rural |
| 1 | 11755881.0 | 42.05 | Rural |
| 2 | 12086707.0 | 43.50 | Rural |
| 3 | 12285693.0 | 43.97 | Rural |
| 4 | 12256762.0 | 44.68 | Rural |
| .. | ... | ... | ... |
| 763 | NaN | NaN | NaN |
| 764 | NaN | NaN | NaN |
| 765 | NaN | NaN | NaN |
| 766 | NaN | NaN | NaN |
| 767 | NaN | NaN | NaN |

[768 rows x 7 columns]>

In [7]:
```python
#Summary statistics in the dataset "umemployment rate till 2020"
print("Shape of the data set ",df1.shape)
print("Size of the data set",df1.size)
print("\n")
print("Info of the dataset \n",df1.info)
print("\n")
print("Descriptive statistics of the dataset \n",df1.describe)
```

Shape of the data set  (267 9)

Shape of the data set (267, 9)
Size of the data set 2403


Info of the dataset
 <bound method DataFrame.info of                    Region        Date  Frequency    Est
imated Unemployment Rate (%)  \
0     Andhra Pradesh   31-01-2020         M                    5.48
1     Andhra Pradesh   29-02-2020         M                    5.83
2     Andhra Pradesh   31-03-2020         M                    5.79
3     Andhra Pradesh   30-04-2020         M                   20.51
4     Andhra Pradesh   31-05-2020         M                   17.43
..               ...          ...       ...                     ...
262      West Bengal   30-06-2020         M                    7.29
263      West Bengal   31-07-2020         M                    6.83
264      West Bengal   31-08-2020         M                   14.87
265      West Bengal   30-09-2020         M                    9.35
266      West Bengal   31-10-2020         M                    9.98

        Estimated Employed   Estimated Labour Participation Rate (%)  Region.1  \
0                 16635535                                     41.02     South
1                 16545652                                     40.90     South
2                 15881197                                     39.18     South
3                 11336911                                     33.10     South
4                 12988845                                     36.46     South
..                     ...                                       ...       ...
262               30726310                                     40.39      East
263               35372506                                     46.17      East
264               33298644                                     47.48      East
265               35707239                                     47.73      East
266               33962549                                     45.63      East

      longitude   latitude
0       15.9129     79.740
1       15.9129     79.740
2       15.9129     79.740
3       15.9129     79.740
4       15.9129     79.740
..          ...        ...
262     22.9868     87.855
263     22.9868     87.855
264     22.9868     87.855
265     22.9868     87.855
266     22.9868     87.855

[267 rows x 9 columns]>


Descriptive statistics of the dataset
 <bound method NDFrame.describe of                    Region        Date  Frequency    E
stimated Unemployment Rate (%)  \
0     Andhra Pradesh   31-01-2020         M                    5.48
1     Andhra Pradesh   29-02-2020         M                    5.83
2     Andhra Pradesh   31-03-2020         M                    5.79
3     Andhra Pradesh   30-04-2020         M                   20.51
4     Andhra Pradesh   31-05-2020         M                   17.43
..               ...          ...       ...                     ...
262      West Bengal   30-06-2020         M                    7.29
263      West Bengal   31-07-2020         M                    6.83
264      West Bengal   31-08-2020         M                   14.87
265      West Bengal   30-09-2020         M                    9.35

```
266      West Bengal   31-10-2020              M                              9.98

         Estimated Employed   Estimated Labour Participation Rate (%) Region.1  \
0                  16635535                                     41.02    South
1                  16545652                                     40.90    South
2                  15881197                                     39.18    South
3                  11336911                                     33.10    South
4                  12988845                                     36.46    South
..                      ...                                       ...      ...
262                30726310                                     40.39     East
263                35372506                                     46.17     East
264                33298644                                     47.48     East
265                35707239                                     47.73     East
266                33962549                                     45.63     East

     longitude   latitude
0      15.9129    79.740
1      15.9129    79.740
2      15.9129    79.740
3      15.9129    79.740
4      15.9129    79.740
..         ...       ...
262    22.9868    87.855
263    22.9868    87.855
264    22.9868    87.855
265    22.9868    87.855
266    22.9868    87.855

[267 rows x 9 columns]>
```

In [8]:
```python
#Names of columns in both datasets
print("Column names in the dataset umemployment in India: \n",df.columns)
print("\n \n")
print("Column names in the dataset umemployment rate till 2020: \n",df1.columns)
```

```
Column names in the dataset umemployment in India:
 Index(['Region', ' Date', ' Frequency', ' Estimated Unemployment Rate (%)',
       ' Estimated Employed', ' Estimated Labour Participation Rate (%)',
       'Area'],
      dtype='object')



Column names in the dataset umemployment rate till 2020:
 Index(['Region', ' Date', ' Frequency', ' Estimated Unemployment Rate (%)',
       ' Estimated Employed', ' Estimated Labour Participation Rate (%)',
       'Region.1', 'longitude', 'latitude'],
      dtype='object')
```

In [9]:
```python
#To check if both the datasets are null or not
print(df.isnull())
print("\n")
print(df1.isnull())
```

```
     Region   Date  Frequency  Estimated Unemployment Rate (%)  \
0     False  False      False                            False
1     False  False      False                            False
2     False  False      False                            False
3     False  False      False                            False
```

```
4        False  False          False                                        False
..        ...    ...            ...                                          ...
763      True   True           True                                         True
764      True   True           True                                         True
765      True   True           True                                         True
766      True   True           True                                         True
767      True   True           True                                         True

     Estimated Employed  Estimated Labour Participation Rate (%)   Area
0               False                                      False  False
1               False                                      False  False
2               False                                      False  False
3               False                                      False  False
4               False                                      False  False
..                ...                                        ...    ...
763              True                                       True   True
764              True                                       True   True
765              True                                       True   True
766              True                                       True   True
767              True                                       True   True

[768 rows x 7 columns]

     Region   Date   Frequency   Estimated Unemployment Rate (%)  \
0    False  False       False                             False
1    False  False       False                             False
2    False  False       False                             False
3    False  False       False                             False
4    False  False       False                             False
..     ...    ...         ...                               ...
262  False  False       False                             False
263  False  False       False                             False
264  False  False       False                             False
265  False  False       False                             False
266  False  False       False                             False

     Estimated Employed  Estimated Labour Participation Rate (%)  Region.1  \
0               False                                      False     False
1               False                                      False     False
2               False                                      False     False
3               False                                      False     False
4               False                                      False     False
..                ...                                        ...       ...
262             False                                      False     False
263             False                                      False     False
264             False                                      False     False
265             False                                      False     False
266             False                                      False     False

     longitude  latitude
0        False     False
1        False     False
2        False     False
3        False     False
4        False     False
..         ...       ...
262      False     False
263      False     False
264      False     False
265      False     False
```
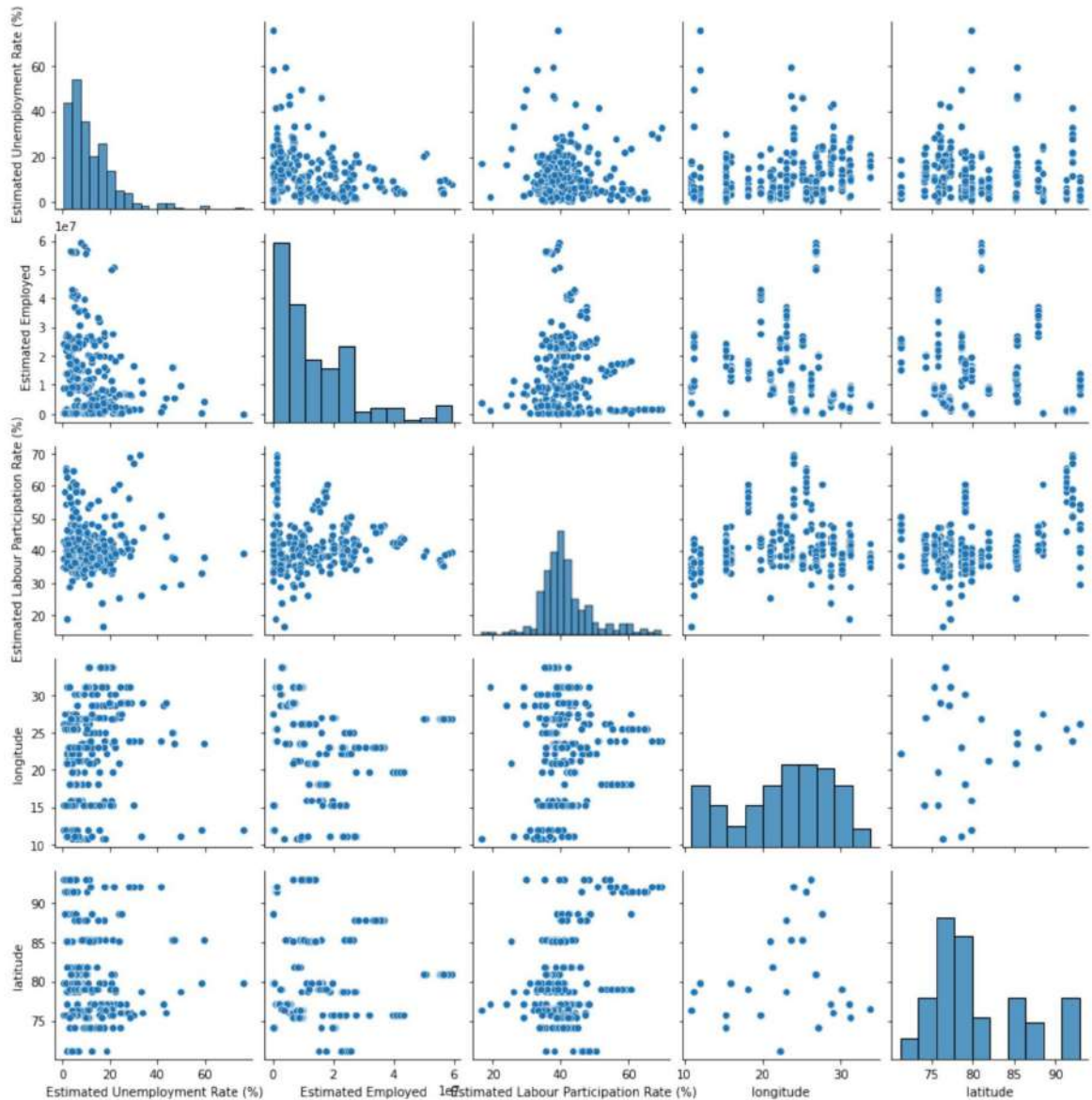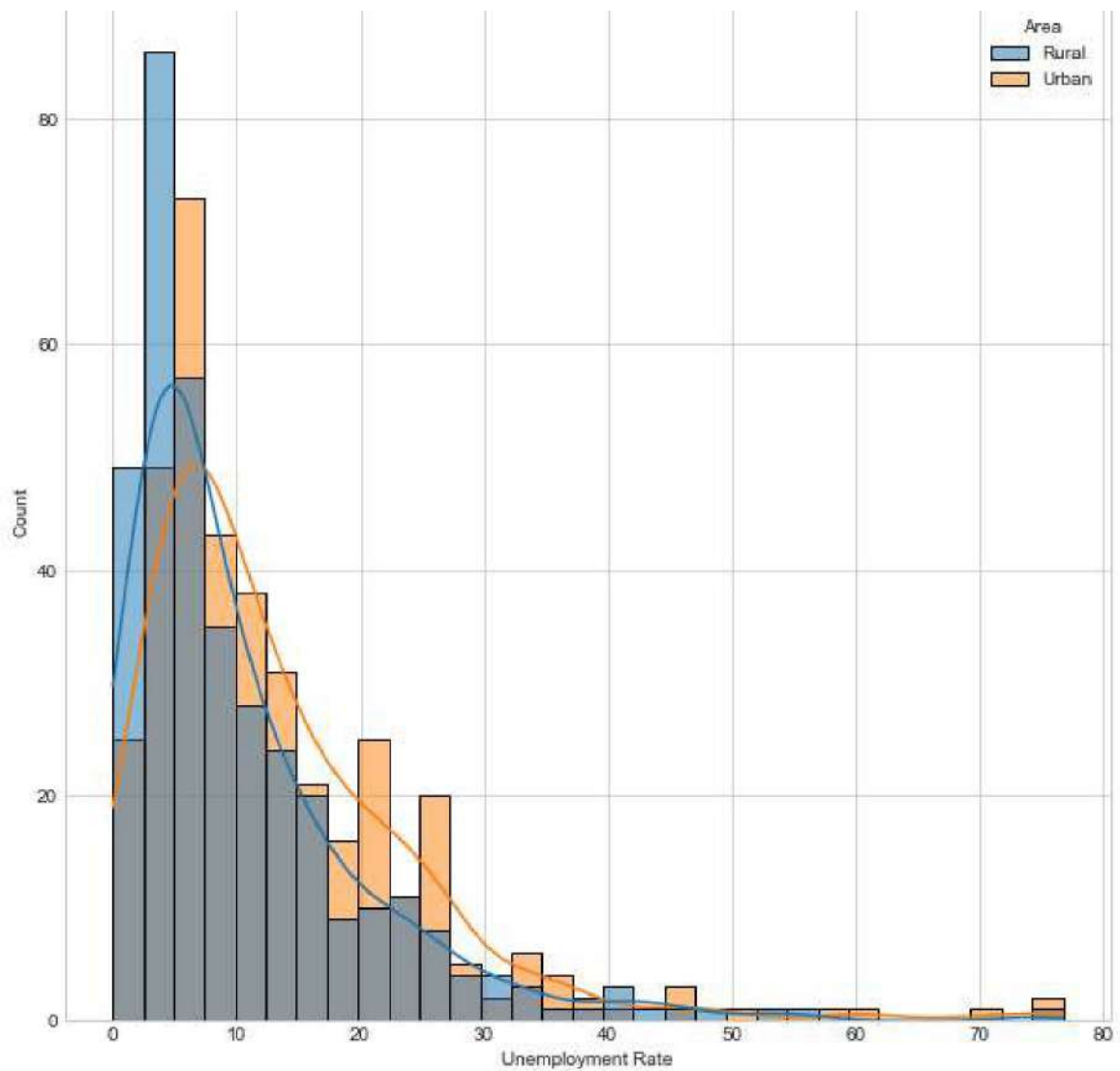
```
265        False        False
266        False        False

[267 rows x 9 columns]
```

In [10]:
```python
#Value counts for dataset
print("For the dataset-Unemployment in India: ")
print(df.isnull().value_counts())
print("\n")
print("For the dataset-Unemployment rate till 2020: ")
print(df1.isnull().value_counts())
```

```
For the dataset-Unemployment in India:
Region   Date   Frequency   Estimated Unemployment Rate (%)   Estimated Employed
Estimated Labour Participation Rate (%)   Area
False   False  False        False                                              False
False                                             False     740
True    True   True         True                                               True
True                                              True      28
dtype: int64


For the dataset-Unemployment rate till 2020:
Region   Date   Frequency   Estimated Unemployment Rate (%)   Estimated Employed
Estimated Labour Participation Rate (%)   Region.1   longitude   latitude
False   False  False        False                                              False
False                                             False     False       False      267
dtype: int64
```

In [11]:
```python
sns.pairplot(df)
```

Out[11]:   <seaborn.axisgrid.PairGrid at 0x1b3d44dae80>

```
In [12]:    sns.pairplot(df1)
```

Out[12]:    <seaborn.axisgrid.PairGrid at 0x1b3d979c460>



```
In [69]:    fig = plt.figure(figsize = (10, 10))
            sns.histplot(x=' Estimated Unemployment Rate (%)', data=df, kde=True, hue='Area'
            plt.title('Unemployment according to Area')
            plt.xlabel('Unemployment Rate')
            plt.show()
```

Unemployment according to Area

```
In [70]:  fig = plt.figure(figsize = (5, 5))
          sns.lineplot(y=' Estimated Unemployment Rate (%)', x=' Date', data=df)
          plt.title('Unemployment according to Date')
          plt.xlabel('Date')
          plt.xticks(rotation=90)
          plt.ylabel('Unemployment Rate')
          plt.show()
```



Unemployment according to Date

```
31-05-2019 30-06-2019 31-07-2019 31-08-2019 30-09-2019 31-10-2019 30-11-2019 31-12-2019 31-01-2020 29-02-2020 31-03-2020 30-04-2020 31-05-2020 30-06-2020
Date
```

In [22]:
```python
df1.columns
```

Out[22]:
```
Index(['Region', ' Date', ' Frequency', ' Estimated Unemployment Rate (%)',
       ' Estimated Employed', ' Estimated Labour Participation Rate (%)',
       'Region.1', 'longitude', 'latitude'],
      dtype='object')
```

In [32]:
```python
fig = plt.figure(figsize = (30, 15))
plt.scatter(df1[' Date'], df1[' Estimated Employed'])

plt.title('Unemployment according to Region')
plt.xlabel('Region')
plt.ylabel('Unemployment Rate')
plt.show()
```



In [34]:
```python
fig = plt.figure(figsize = (30, 15))
sns.histplot(x=' Estimated Labour Participation Rate (%)', data=df, kde=True, hu
plt.title('Labour Participation according to Area')
plt.xlabel('Labour Participation Rate')
plt.show()
```

In [36]:
```python
fig = plt.figure(figsize = (7, 7))
sns.lineplot(y=' Estimated Labour Participation Rate (%)', x=' Date', data=df)
plt.title('Labour Participation according to Date')
plt.xlabel('Date')
plt.xticks(rotation=90)
plt.ylabel('Labour Participation Rate')
plt.show()
```



In [71]:
```python
y=df[' Estimated Unemployment Rate (%)']
x=df['Region']
plt_1 = plt.figure(figsize=(10, 10))
plt.title('Umemployment Rate', fontweight='bold' ,fontsize=20)
plt.xlabel("States",fontweight='bold',fontsize=20)
plt.ylabel("Estimated Unemployment rate",fontweight='bold',fontsize=20)
plt.xticks(rotation='vertical',fontsize=12)
```

```
sns.histplot(x, color='lavender')
```

`<AxesSubplot:title={'center':'Umemployment Rate'}, xlabel='States', ylabel='Esti mated Unemployment rate'>`



**Umemployment Rate**

```
fig = plt.figure(figsize = (9, 9))
plt.scatter(df1['Region.1'], df1[' Estimated Labour Participation Rate (%)'])
plt.title('Labour Participation according to Region')
plt.xlabel('Region')
plt.ylabel('Labour Participation Rate')
plt.show()
```



Labour Participation according to Region

```python
sns.histplot(x=' Estimated Employed', data=df, kde=True, hue='Area')
plt.title('Employment according to Area')
plt.xlabel('Employment Rate')
plt.show()
```



Employment according to Area

```python
fig = plt.figure(figsize = (9, 9))
sns.lineplot(y=' Estimated Employed', x=' Date', data=df)
plt.title('Employment according to Date')
plt.xlabel('Date')
plt.xticks(rotation=90)
```

```
plt.ylabel('Employment Rate')
plt.show()
```



Employment according to Date

In [54]:
```
fig = plt.figure(figsize = (9, 9))
plt.scatter(df1['Region.1'], df1[' Estimated Employed'])
plt.title('Employment according to Region')
plt.xlabel('Region')
plt.ylabel('Employment Rate')
plt.show()
```
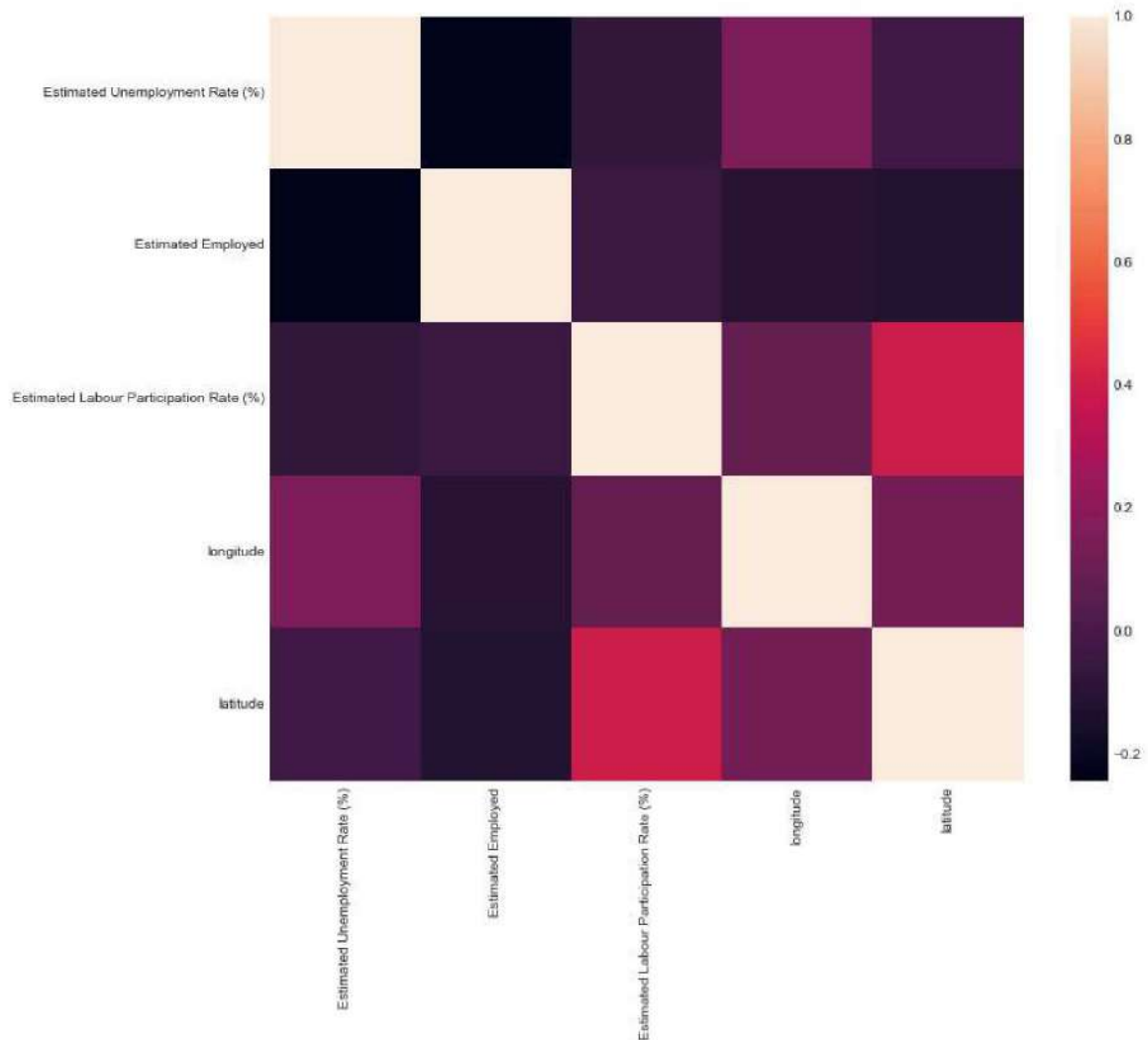


Employment according to Region

```
#Now let's have a look at the correlation between the features of this dataset:
plt.style.use('seaborn-whitegrid')
plt.figure(figsize=(12, 10))
sns.heatmap(df.corr())
plt.show()
```
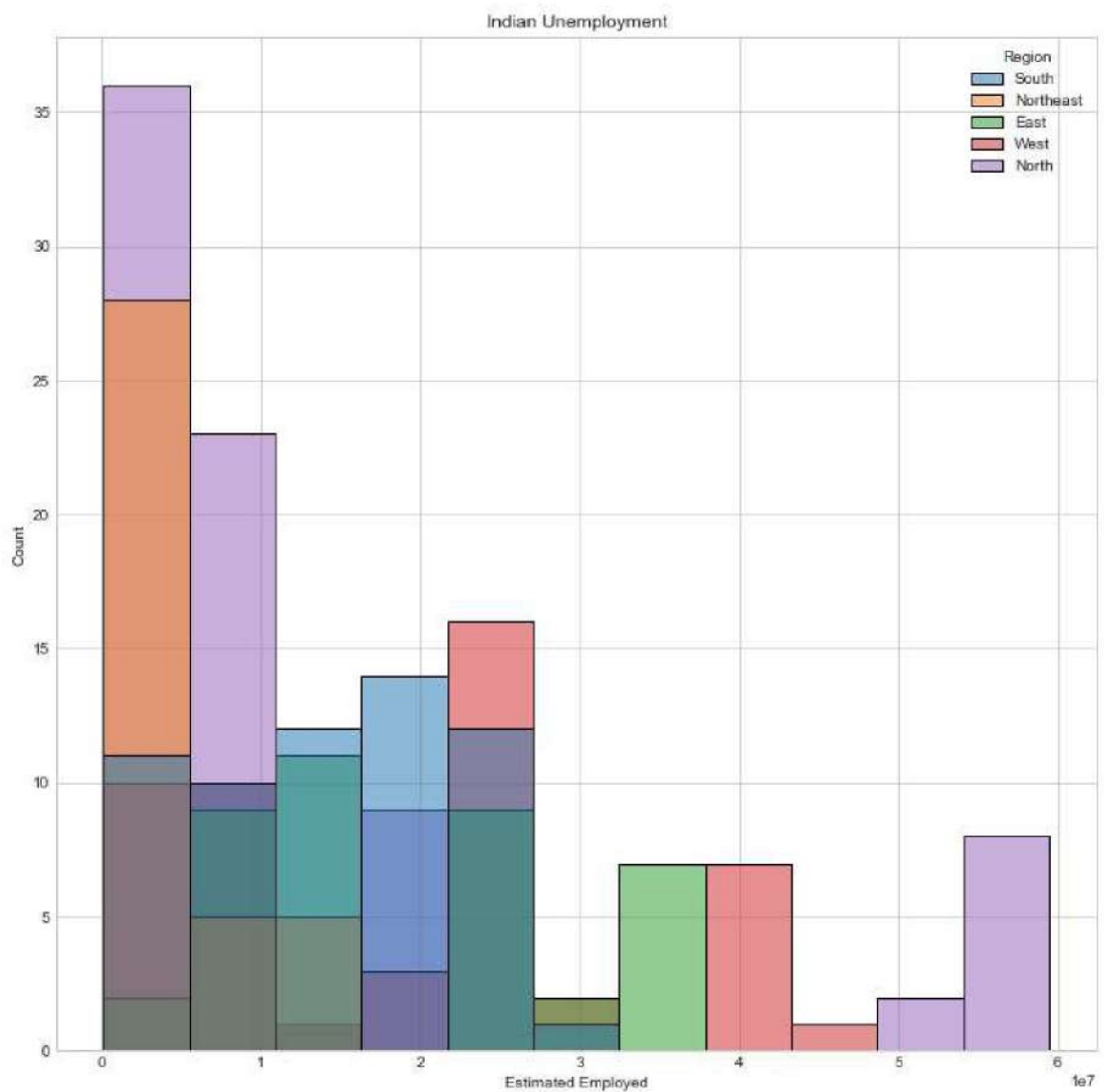
In [56]:
```python
#Now let's have a look at the correlation between the features of this dataset:
plt.style.use('seaborn-whitegrid')
plt.figure(figsize=(12, 10))
sns.heatmap(df1.corr())
plt.show()
```



In [61]:
```python
df.columns
```

Out[61]:
```
Index(['Region', ' Date', ' Frequency', ' Estimated Unemployment Rate (%)',
       ' Estimated Employed', ' Estimated Labour Participation Rate (%)',
       'Area'],
      dtype='object')
```

In [65]:
```python
plt.figure(figsize=(12, 12))
df1.columns= ["States","Date","Frequency",
              "Estimated Unemployment Rate","Estimated Employed",
              "Estimated Labour Participation Rate","Region",
              "longitude","latitude"]
plt.title("Indian Unemployment")
sns.histplot(x="Estimated Employed", hue="Region", data=df1)
```

```
sns.histplot(x="Estimated Employed", hue="Region", data=df1)
plt.show()
```



Indian Unemployment

In [66]:
```
#Now let's see the unemployment rate according to different regions of India:
plt.figure(figsize=(12, 11))
plt.title("Indian Unemployment")
sns.histplot(x="Estimated Unemployment Rate", hue="Region", data=df1)
plt.show()
```



Indian Unemployment