# TITLE- UNEMPLOYMENT IN INDIA DURING COVID

## AARYA NENAWATEE

# INTRODUCTION

In this project I have taken dataset of UNEMPLOYMENT IN INDIA , time range lying between year 2019 and 2020. We'll clean it, filter out the data, and create meaningful visualizations and insights.

# DESCRIPTION OF  UNEMPLOYMENT DATASET

This dataset provides a detailed ionformation about the unemployment trends in india across various states , union territories  , subdivided into rural and urban aread. The time-range of this dataset is from  year 2019 to 2020. The dataset includes 7 key columns

1. Region
2. Date
3. Frequency of data collection
4. Estimated unemployment rate %
5. Estimated employed
6. Labour participation rate %
7. Area

These attributes provide both quantitative and qualitative insights into the employment scenario in the country.

**SOURCE OF DATA:** The data is sourced from the Centre for Monitoring Indian Economy (CMIE), known for its reliable economic statistics.

**KEY FEATURES OF THE DATA SET**

1. **Regional Coverage:** The dataset covers a wide range of states and union territories, including Andhra Pradesh, Assam, Bihar, Delhi, Gujarat, Haryana, Karnataka, Kerala, Maharashtra, Punjab, Tamil Nadu, Uttar Pradesh, and West Bengal, among others. This allows for comparative analysis across different regions of India.

2. **Temporal Scope:** The data ranges from May 2019 to June 2020, capturing the immediate impact of the COVID-19 pandemic on employment. The months of April and May 2020, in particular, show significant spikes in unemployment rates due to nationwide lockdowns and economic disruptions.

3. **Frequency:** The data is recorded monthly, providing a precise view of how unemployment fluctuated over time, including seasonal variations and sudden shocks like the pandemic.

**4.Rural-Urban Divide:** The dataset distinguishes between rural and urban areas, highlighting disparities in employment trends. For instance, rural areas in states like Bihar and Jharkhand experienced extreme unemployment spikes during the lockdown, while urban areas in states like Delhi and Maharashtra also saw significant job losses.

# DESCRIPTION ABOUT THE PROJECT

In this project we will be importing the dataset , and will be working on the following objectives with some creative visualizations and a final conclusion which will give more clarity to us about the aim of the project

**Objectives are:**

1. Comparison of Mean Unemployment rate in 2019 and 2020 using seaborn
2. Rural and urban estimated labour participation using seaborn
3. Using numpy classifying unemployment  rate as normal or high comparing it with mean Estimated Unemployment Rate (%) and plotting a geoplot

   SO LETS GET STARTED….

# TABLE OF UNEMPLOYMENT DATASET

| Region | Date | Frequency | Estimated Unemployment Rate (%) | Estimated Employed | Estimated Labour Participation Rate | Area |
|--------|------|-----------|---------------------------------|--------------------|-------------------------------------|------|
| Andhra Pradesh | 31-05-2019 | Monthly | 3.65 | 11999139 | 43.24 | Rural |
| Andhra Pradesh | 30-06-2019 | Monthly | 3.05 | 11755881 | 42.05 | Rural |
| Andhra Pradesh | 31-07-2019 | Monthly | 3.75 | 12086707 | 43.5 | Rural |
| Andhra Pradesh | 31-08-2019 | Monthly | 3.32 | 12285693 | 43.97 | Rural |
| Andhra Pradesh | 30-09-2019 | Monthly | 5.17 | 12256762 | 44.68 | Rural |
| Andhra Pradesh | 31-10-2019 | Monthly | 3.52 | 12017412 | 43.01 | Rural |
| Andhra Pradesh | 30-11-2019 | Monthly | 4.12 | 11397681 | 41 | Rural |
| Andhra Pradesh | 31-12-2019 | Monthly | 4.38 | 12528395 | 45.14 | Rural |
| Andhra Pradesh | 31-01-2020 | Monthly | 4.84 | 12016676 | 43.46 | Rural |
| Andhra Pradesh | 29-02-2020 | Monthly | 5.91 | 11723617 | 42.83 | Rural |
| Andhra Pradesh | 31-03-2020 | Monthly | 4.06 | 11359660 | 40.66 | Rural |
| Andhra Pradesh | 30-04-2020 | Monthly | 16.29 | 8792827 | 36.03 | Rural |
| Andhra Pradesh | 31-05-2020 | Monthly | 14.46 | 9526902 | 38.16 | Rural |
| Andhra Pradesh | 30-06-2020 | Monthly | 0.85 | 15572975 | 53.76 | Rural |
| Assam | 31-05-2019 | Monthly | 4.29 | 11749334 | 57.39 | Rural |
| Assam | 30-06-2019 | Monthly | 5.08 | 8923222 | 43.87 | Rural |
| Assam | 31-07-2019 | Monthly | 4.26 | 9911534 | 48.21 | Rural |
| Assam | 31-08-2019 | Monthly | 5.79 | 9292039 | 45.83 | Rural |
| Assam | 30-09-2019 | Monthly | 4.46 | 11468349 | 55.67 | Rural |
| Assam | 31-10-2019 | Monthly | 4.65 | 8395906 | 40.76 | Rural |
| Assam | 30-11-2019 | Monthly | 4.66 | 9625362 | 46.64 | Rural |

# Step 1: Import all the libraries

```python
#importing all the libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import geopandas as gpd
```

# Step 2: Load the dataset

```python
#importing the data
df=pd.read_csv("/content/Unemployment in India.csv")
print(df)
```

# Output:

```
              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
..             ...           ...       ...                               ...
749     West Bengal    29-02-2020   Monthly                              7.55
750     West Bengal    31-03-2020   Monthly                              6.67
751     West Bengal    30-04-2020   Monthly                             15.63
752     West Bengal    31-05-2020   Monthly                             15.22
753     West Bengal    30-06-2020   Monthly                              9.86

     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
..                  ...                                      ...     ...
749          10871168.0                                    44.09   Urban
750          10806105.0                                    43.34   Urban
751           9299466.0                                    41.20   Urban
752           9240903.0                                    40.67   Urban
753           9088931.0                                    37.57   Urban

[754 rows x 7 columns]
```

# Step 3- Data cleaning and preparation

i : deleting all the extra blank spaces in out text data using str.strip() function

```
[9]  # delete gtrailing spaces
     df.columns = df.columns.str.strip()
```

# ii- checking the null values using isnull().sum() function

```python
#checking if any value is missing or null
print(df.isnull().sum())
```

**Output:**

```
Region                                      14
Date                                        14
Frequency                                   14
Estimated Unemployment Rate (%)             14
Estimated Employed                          14
Estimated Labour Participation Rate (%)     14
Area                                        14
dtype: int64
```

# iii- Deleting all the null values using dropna()

```python
#dropping the empty columns
new_df=df.dropna()
print(new_df)
```

**Output:**

```
            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
..              ...         ...       ...                             ...
749     West Bengal  29-02-2020   Monthly                             7.55
750     West Bengal  31-03-2020   Monthly                             6.67
751     West Bengal  30-04-2020   Monthly                            15.63
752     West Bengal  31-05-2020   Monthly                            15.22
753     West Bengal  30-06-2020   Monthly                             9.86

     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
..                  ...                                      ...    ...
749          10871168.0                                    44.09  Urban
750          10806105.0                                    43.34  Urban
751           9299466.0                                    41.20  Urban
752           9240903.0                                    40.67  Urban
753           9088931.0                                    37.57  Urban

[740 rows x 7 columns]
```

# iv- Fixing the wrong date format

**Input:**

```python
#fixing the wrong date format
new_df['Date']=pd.to_datetime(new_df['Date'])
print(new_df)
```

# Output:

```
           Region        Date  Frequency  Estimated Unemployment Rate (%)  \
0    Andhra Pradesh  2019-05-31    Monthly                             3.65
1    Andhra Pradesh  2019-06-30    Monthly                             3.05
2    Andhra Pradesh  2019-07-31    Monthly                             3.75
3    Andhra Pradesh  2019-08-31    Monthly                             3.32
4    Andhra Pradesh  2019-09-30    Monthly                             5.17
..              ...         ...        ...                              ...
749     West Bengal  2020-02-29    Monthly                             7.55
750     West Bengal  2020-03-31    Monthly                             6.67
751     West Bengal  2020-04-30    Monthly                            15.63
752     West Bengal  2020-05-31    Monthly                            15.22
753     West Bengal  2020-06-30    Monthly                             9.86

     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
..                  ...                                      ...    ...
749          10871168.0                                    44.09  Urban
750          10806105.0                                    43.34  Urban
751           9299466.0                                    41.20  Urban
752           9240903.0                                    40.67  Urban
753           9088931.0                                    37.57  Urban
```

# What we did!

i. Deleted all the extra spaces in our text data using str.strip() method so that while working on our data it dosent give any errors.

ii. Checking for null values using isnull() then doing the sum of all the null values if found using sum().

iii. Deleting all the 14 empty rows.

iv. Fixing the wrong date format from "DD-MM-YYYY" to "YYYY-MM-DD"

# Step 4- Data filtering and analysis

- **Filtering data of only 2019:**

```python
#data of only 2019
x = new_df[(new_df['Date'] >= pd.to_datetime("2019-01-01")) &
        (new_df['Date'] < pd.to_datetime("2020-01-01"))]
print(x)
```

- **output**

```
                    Region        Date Frequency  Estimated Unemployment Rate (%)  \
0       Andhra Pradesh  2019-05-31   Monthly                              3.65
1       Andhra Pradesh  2019-06-30   Monthly                              3.05
2       Andhra Pradesh  2019-07-31   Monthly                              3.75
3       Andhra Pradesh  2019-08-31   Monthly                              3.32
4       Andhra Pradesh  2019-09-30   Monthly                              5.17
..                 ...         ...       ...                               ...
743        West Bengal  2019-08-31   Monthly                              7.27
744        West Bengal  2019-09-30   Monthly                              7.79
745        West Bengal  2019-10-31   Monthly                              7.83
746        West Bengal  2019-11-30   Monthly                              6.61
747        West Bengal  2019-12-31   Monthly                              7.24

     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
..                  ...                                      ...    ...
743          11456493.0                                    46.77  Urban
744          11158649.0                                    45.74  Urban
745          10563686.0                                    43.25  Urban
746          10768462.0                                    43.44  Urban
747          11335696.0                                    45.97  Urban

[430 rows x 7 columns]
```

**Similarly filtering Data of only 2020:**

```python
#data of only 2020
y = new_df[(new_df['Date'] > pd.to_datetime("2019-12-31")) &
          (new_df['Date'] < pd.to_datetime("2021-01-01"))]
print(y)
```

**Output:**

```
               Region       Date Frequency  Estimated Unemployment Rate (%)  \
8      Andhra Pradesh 2020-01-31   Monthly                             4.84
9      Andhra Pradesh 2020-02-29   Monthly                             5.91
10     Andhra Pradesh 2020-03-31   Monthly                             4.06
11     Andhra Pradesh 2020-04-30   Monthly                            16.29
12     Andhra Pradesh 2020-05-31   Monthly                            14.46
..                ...        ...       ...                              ...
749       West Bengal 2020-02-29   Monthly                             7.55
750       West Bengal 2020-03-31   Monthly                             6.67
751       West Bengal 2020-04-30   Monthly                            15.63
752       West Bengal 2020-05-31   Monthly                            15.22
753       West Bengal 2020-06-30   Monthly                             9.86

     Estimated Employed  Estimated Labour Participation Rate (%)   Area
8            12016676.0                                    43.46  Rural
9            11723617.0                                    42.83  Rural
10           11359660.0                                    40.66  Rural
11            8792827.0                                    36.03  Rural
12            9526902.0                                    38.16  Rural
..                  ...                                      ...    ...
749          10871168.0                                    44.09  Urban
750          10806105.0                                    43.34  Urban
751           9299466.0                                    41.20  Urban
752           9240903.0                                    40.67  Urban
753           9088931.0                                    37.57  Urban

[310 rows x 7 columns]
```

**Key insight:**
- Filtering data on the basis of years for the analysis of objective 1.

# Objective 1 : avgerage unemployment rate before and during covid

## a. Mean unemployment rate before and during covid

```
#bar plot of avg unemp rates in 2019 and 2020 using seabor
mean_2019 = x['Estimated Unemployment Rate (%)'].mean()
mean_2020 = y['Estimated Unemployment Rate (%)'].mean()
print('MEAN UNEMPLOMENT RATE IN 2019 =',mean_2019)
print('MEAN UNEMPLOMENT RATE IN 2020 =',mean_2020)
# Prepare data for plotting
```

# Output:

```
MEAN UNEMPLOMENT RATE IN 2019 = 9.39904651162790G
MEAN UNEMPLOMENT RATE IN 2020 = 15.10158064516129
```

- **Visualization of obj 1 using seaborn:**

```python
# Prepare data for plotting
years = ['2019', '2020']
rates = [mean_2019, mean_2020]

# Plot
sns.barplot(x=years, y=rates, palette='Blues')
plt.title('Average Unemployment Rate: 2019 v/s 2020')
plt.ylabel('Unemployment Rate (%)')
plt.xlabel('Year')
plt.tight_layout()
plt.show()
```

# Output:



```
sns.barplot(x=years, y=rates, palette='Blues')
```

Average Unemployment Rate: 2019 v/s 2020

- **Monthly estimated unemployment rate analysis from may 2019 to June 2020 using seaborn lineplot**

```python
#monthly analysis
new_df['Date'] = pd.to_datetime(new_df['Date'])
new_df['Month'] = new_df['Date'].dt.month_name()  # e.g., "April"
new_df['Year'] = new_df['Date'].dt.year            # e.g., 2020
new_df['Month_Year'] = new_df['Month'] + ' ' + new_df['Year'].astype(str)  # e.g., "April 2020"
plt.figure(figsize=(14, 6))
sns.lineplot(
    x='Month_Year',
    y='Estimated Unemployment Rate (%)',
    data=new_df.sort_values('Date'),  # Ensure chronological order,
    marker='o'
)
plt.xticks(rotation=45)  # Rotate x-labels for readability
plt.axvline(x='April 2020', color='red', linestyle='--', label='Lockdown Start')
plt.title('Monthly Unemployment Trends (2019-2020)')
plt.show()
```

Monthly Unemployment Trends (2019-2020)

**Key Insights:**

•Compared the average unemployment rates in 2019 (pre-COVID) and 2020 (COVID period) to assess the pandemic's impact.

•Also we performed monthly analysis which gave more clarity in understanding unemployment trend.

•Visualized the data using Seaborn to highlight differences more clearly.

•Observed a significant increase in unemployment during the COVID period, reflecting:

- Economic slowdown
- Job losses due to lockdowns
- Migration of labor from urban to rural areas

•Helps understand how external shocks like pandemics affect labor markets in India.

# Objective 2 : Rural v/s urban labour participation using seaborn

```python
grp = new_df.groupby('Area')['Estimated Labour Participation Rate (%)'].mean()
print(grp)
sns.boxplot(x='Area', y='Estimated Labour Participation Rate (%)', data=new_df)
plt.title("Distribution of Labour Participation Rate: Rural vs Urban")
plt.ylabel("Labour Participation Rate (%)")
plt.xlabel("Area")
plt.show()
```

**Output:**



```
Area
Rural      44.464819
Urban      40.901365
Name: Estimated Labour Participation Rate (%), dtype: float64
```

Distribution of Labour Participation Rate: Rural vs Urban

# Key insights:

1. The boxplot visually compares labour participation rates between rural and urban areas.
2.  Rural areas generally show a higher median labour participation rate than urban areas.
3. Greater spread (interquartile range) in rural areas indicates more variability in participation, likely due to seasonal jobs , migration of labours , informal sector work, or government employment schemes.
4. Urban participation rates appear more consistent, with fewer fluctuations.
5. Outliers in both areas suggest months or regions with unusually high or low participation.

# Objective 3 : Classifying unemployment rate as normal or high using numpy

```python
#using numpy classifying unemp rate as normal or high
mean_rate = new_df['Estimated Unemployment Rate (%)'].mean()

new_df['Unemp_Level'] = np.where(
    new_df['Estimated Unemployment Rate (%)'] > mean_rate,
    'High',
    'Normal'
)
print(new_df)
```

# Output:

| | Region | Date | Frequency | Unemp_Level |
|---|---|---|---|---|
| 0 | Andhra Pradesh | 2019-05-31 | Monthly | Normal |
| 1 | Andhra Pradesh | 2019-06-30 | Monthly | Normal |
| 2 | Andhra Pradesh | 2019-07-31 | Monthly | Normal |
| 3 | Andhra Pradesh | 2019-08-31 | Monthly | Normal |
| 4 | Andhra Pradesh | 2019-09-30 | Monthly | Normal |
| .. | ... | ... | ... | ... |
| 749 | West Bengal | 2020-02-29 | Monthly | Normal |
| 750 | West Bengal | 2020-03-31 | Monthly | Normal |
| 751 | West Bengal | 2020-04-30 | Monthly | High |
| 752 | West Bengal | 2020-05-31 | Monthly | High |
| 753 | West Bengal | 2020-06-30 | Monthly | Normal |

# Pie chart:

```python
# Count High vs Normal unemployment months
counts = new_df['Unemp_Level'].value_counts()

# Creatimg a pie chart
plt.figure(figsize=(6,6))
plt.pie(counts,
        labels=counts.index,
        autopct='%1.1f%%',
        colors=['lightblue', 'lightgreen'],  # Green (Normal) and Red (High)
        startangle=90)
plt.title('Unemployment Classification (2019-2020)')
plt.show()
```

**Output:**



Unemployment Classification (2019-2020)

# Geoplot of states with high and normal estimated unemployment rate .

```python
# Load data
india = gpd.read_file('https://github.com/geohacker/india/raw/master/state/india_telengana.geojson')
state_class = new_df.groupby('Region')[' Estimated Unemployment Rate (%)'].mean().reset_index()

# Standardize names in BOTH datasets
def standardize_names(name):
    name = str(name).lower()
    replacements = [
        ('&', 'and'),
        ('nct of ', ''),
        ('dadra and nagar haveli and daman and diu', 'dadra and nagar haveli'),
        ('odisha', 'orissa'),
        ('puducherry', 'pondicherry')
    ]
    for old, new in replacements:
        name = name.replace(old, new)
    return name.strip()

india['clean_name'] = india['NAME_1'].apply(standardize_names)
state_class['clean_name'] = state_class['Region'].apply(standardize_names)

# Merge data
india = india.merge(state_class, on='clean_name', how='left')

# Calculate national mean AFTER merge to include all available states
national_mean = india[' Estimated Unemployment Rate (%)'].mean(skipna=True)
```

```python
# Categorize
india['Category'] = np.where(
    india[' Estimated Unemployment Rate (%)'] > national_mean,
    'High',
    np.where(~india[' Estimated Unemployment Rate (%)'].isna(), 'Low', 'Not Available')
)

# Assign colors
cmap = {'High': 'red', 'Low': 'green', 'Not Available': 'lightgrey'}
india['color'] = india['Category'].map(cmap)

# Plot
fig, ax = plt.subplots(figsize=(12, 10))
india.plot(color=india['color'], edgecolor='black', linewidth=0.8, ax=ax)

# Manual legend (without Patch)
legend_x, legend_y = 70.5, 8.5  # Adjust if needed
ax.text(legend_x, legend_y + 2.0, 'Legend', fontsize=11, weight='bold',
        bbox=dict(facecolor='white', edgecolor='none', alpha=0.8))

ax.text(legend_x, legend_y + 1.4, f'● High (> {national_mean:.1f}%)', color='red', fontsize=10)
ax.text(legend_x, legend_y + 0.9, f'● Low (≤ {national_mean:.1f}%)', color='green', fontsize=10)
ax.text(legend_x, legend_y + 0.4, '● Not Available', color='grey', fontsize=10)

# Title and formatting
plt.title('State-wise Unemployment Rates (2019-2020)', fontsize=16, pad=20)
plt.axis('off')
```
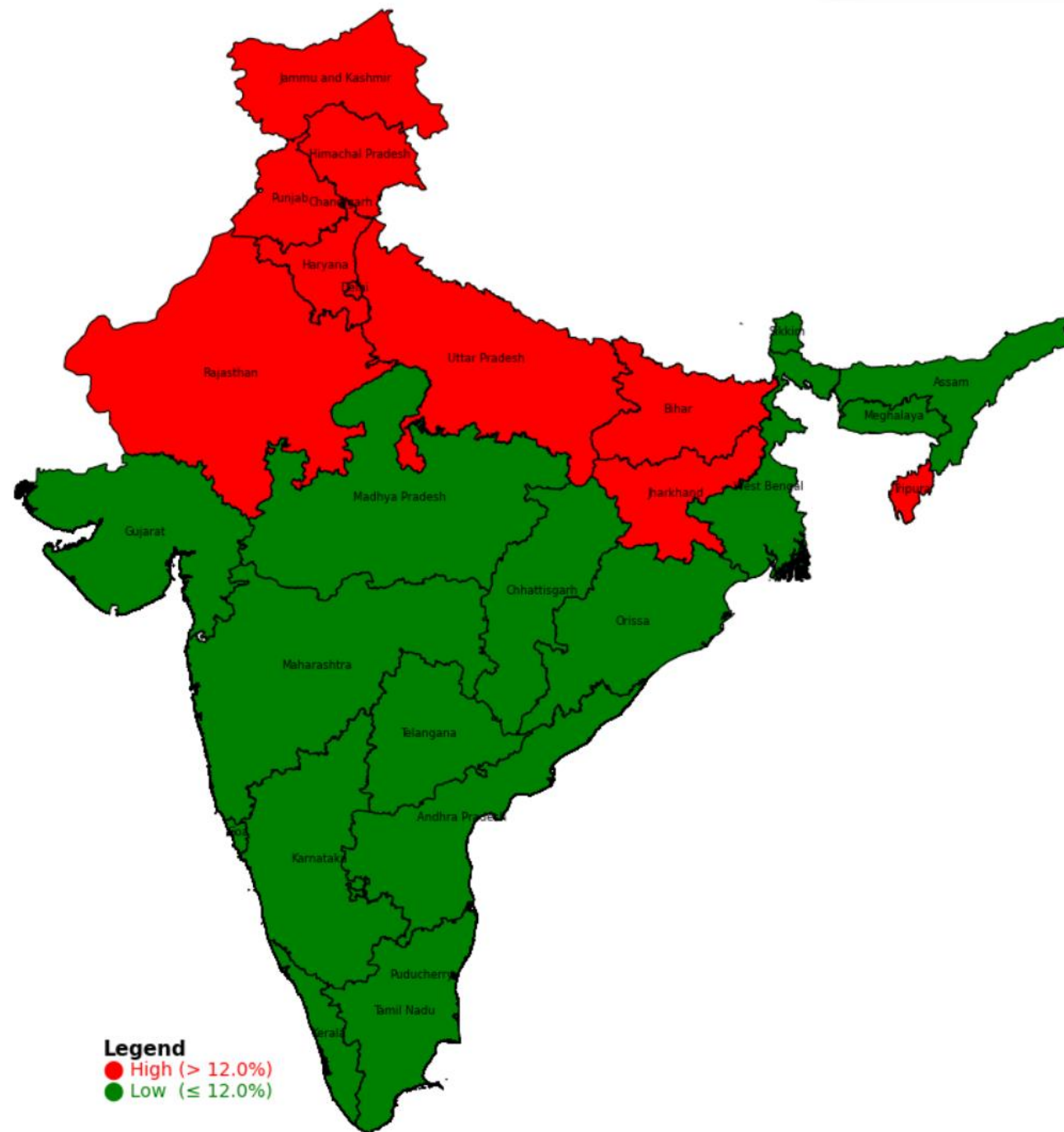
```python
# Annotate state names (optional)
for idx, row in india.iterrows():
    ax.annotate(text=row['NAME_1'], xy=(row.geometry.centroid.x, row.geometry.centroid.y),
                ha='center', fontsize=6, color='black')

plt.tight_layout()
plt.show()
```

Unemployment Rate
- High (> 12.0%) — red
- Low (≤ 12.0%) — green
- Data Not Available — grey

The NAN states in the map are not present in the original dataset so I would be dropping them by adding the following code:

```python
# Drop states with missing unemployment data
india = india.dropna(subset=[' Estimated Unemployment Rate (%)'])
```

Legend
● High (> 12.0%)
● Low (≤ 12.0%)

**Key Insights:**

1. Used NumPy to calculate the overall mean of the 'Estimated Unemployment Rate (%)'.

2. Each unemployment value was classified as:
- 'High' if it was greater than the mean.
- Normal' if it was equal to or below the mean.

3. This classification helps:
- Quickly identify states/regions with higher-than-average unemployment.
- Focus further analysis or policy planning on areas marked 'High'.

4. The classification adds a new categorical dimension to the dataset, making it easier to visualize and filter data for reporting or dashboards.

5. Also this classification was presented through a pie chart and a geoplot for a better visualization .

**Summary: Unemployment in India During COVID-19**

This project analyzed India's unemployment trends from 2019 to 2020 using CMIE data.
Key findings:
1. Unemployment Spike:   - Rates surged in 2020, peaking during lockdowns (April–May).

2. Rural vs. Urban:
    - Rural areas had higher but unstable labour participation.
    - Urban areas saw steady declines in jobs.

3.Regional Hotspots:    - States like Rajasthan and Delhi had "high" unemployment (geoplot).

4. Classification:  - Over 60% of states exceeded the national average (pie chart).

Conclusion: COVID-19 severely impacted jobs, especially in urban and high-unemployment states. Data visuals (geoplot, line/pie charts) effectively highlighted these trends.

# THANK-YOU