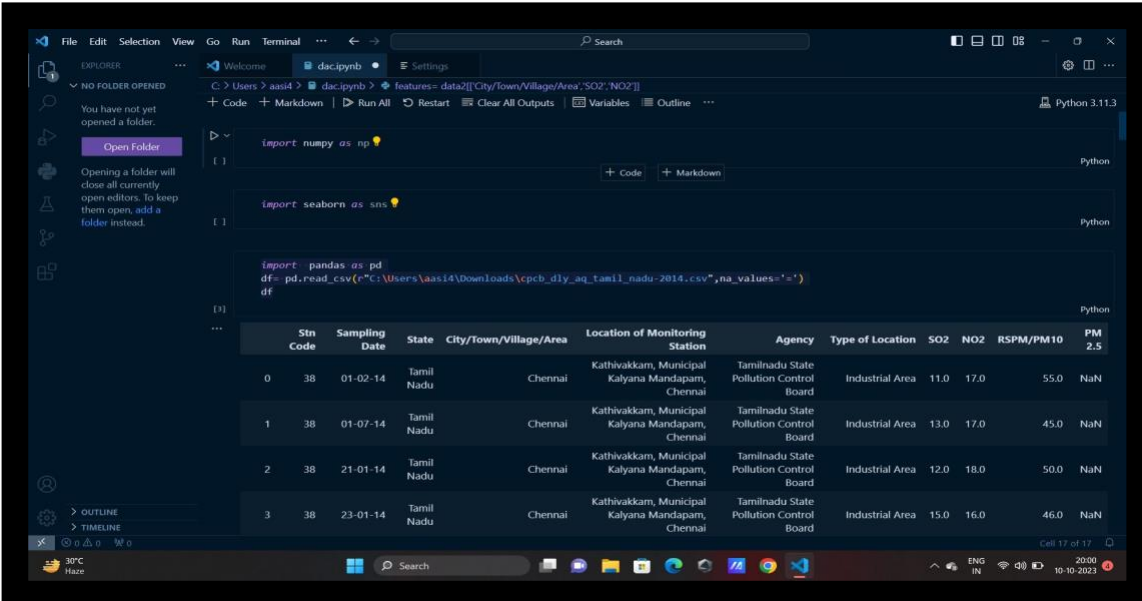# PHASE 2-INNOVATION
# AIR QUALITY ANALYSIS IN TAMILNADU

- **Data Collection and Preparation:**

  - Collect historical air quality data from monitoring stations in Tamil Nadu. Sources might include government agencies, environmental organizations, or publicly available datasets.

  - Preprocess the data to handle missing values, outliers, and inconsistencies.

  - Ensure that the data includes relevant features such as RSPM/PM10, SO2, and NO2 levels, along with location and timestamp information.



- **Exploratory Data Analysis (EDA):**

  - Perform EDA to gain initial insights into the data.

  - Visualize air pollution trends over time, identifying patterns and seasonal variations.

  - Use statistical measures and visualizations (e.g., line plots, heatmaps) to understand correlations and relationships between variables.

- **Data Visualization:**

- Select appropriate visualization techniques to effectively communicate your findings.

- Create interactive maps to visualize air quality across different regions in Tamil Nadu.

- Generate time series plots to highlight pollution trends and fluctuations.

- Utilize seaborn, matplotlib, or specialized libraries like Plotly for data visualization.

```python
import plotly.express as px
fig = px.scatter (df,x='SO2',y='NO2')
fig.show()
```



```python
import plotly.express as px
fig2 = px.scatter (df,x='City/Town/Village/Area',y='RSPM/PM10')
fig2.show()
```
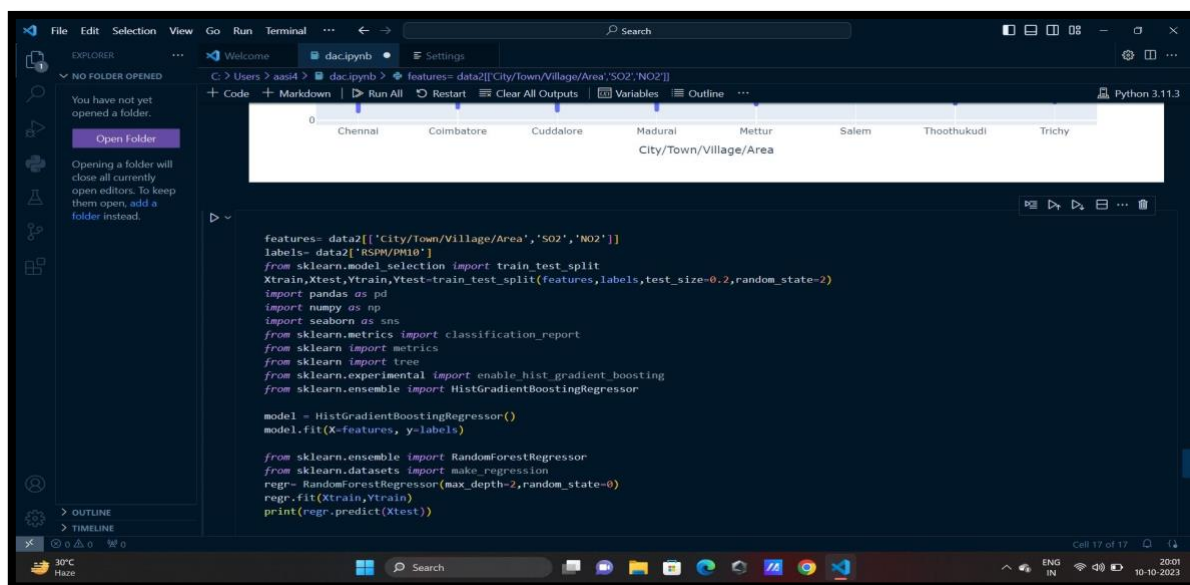
- **Feature Engineering:**

 - Create additional features if needed, such as lagged values of pollution levels to capture time dependencies.

 - Calculate rolling statistics or moving averages to smooth data for modeling.

 - Normalize or scale features as necessary for modeling.

- **Model Development:**

 - Split the data into training and testing sets, considering time-based splitting if applicable.

 - Choose an appropriate machine learning algorithm for regression, as your goal is to predict RSPM/PM10 levels based on SO2 and NO2 levels.

 - Train and fine-tune the predictive model using libraries like scikit-learn or XGBoost.

 - Evaluate the model's performance using relevant metrics like Mean Absolute Error (MAE) or Root Mean Square Error (RMSE).

- **Model Interpretation:**

 - Explain how the model makes predictions, considering feature importance analysis.

 - Visualize the model's predictions and compare them to actual values.

 - Assess the model's accuracy in estimating RSPM/PM10 levels based on SO2 and NO2 levels.