# Analyzing Air Quality Data: Unveiling India's Pollution Patterns

## Data Science Bootcamp - Spring 2024

**Group 4**

Jesse Noppe-Brandon

Mithuna Grandhe

Yixin Xu

Kripa Shankar S R

# Introduction

Air quality refers to **the degree to which the air is suitable or clean enough for humans or the environment**.

India is currently one of the major countries of rapid development but it is also very high prone to low air quality levels which cause a lot of diseases like lung cancer as well as even death.

# Air Quality Monitoring and Importance



AIR QUALITY INDEX

Air Quality Index (AQI) is **used for reporting daily air quality**. It tells how clean or polluted your air is, and what associated health effects might be a concern for you. AQI focuses on health effects you may experience within a few hours or days after breathing polluted air.

## 1

### Public Health Impact

Poor air quality significantly impacts public health, leading to increased respiratory issues, cardiovascular problems, and reduced life expectancy, especially for vulnerable populations.

## 2

### Environmental Regulations

Comprehensive air quality monitoring is essential for enforcing environmental regulations and ensuring compliance with air pollution standards set by government agencies.

## 3

### Pollution Source Identification

Air quality data can help identify the sources of pollution, enabling targeted interventions and policies to address the root causes of poor air quality.

## 4

### Informed Decision-Making

Real-time air quality data empowers citizens, policymakers, and urban planners to make informed decisions about transportation, industrial activities, and other factors affecting air pollution.

# Data Collection and Preprocessing

Effective air quality control requires robust data collection strategies. Networked IoT sensors capture real-time measurements of pollutants, particulates, and environmental factors across urban and industrial areas. Preprocessing this raw sensor data is crucial, involving techniques like outlier removal, data imputation, and feature scaling to ensure data integrity.



```python
#print(airq.isnull().sum())
airq_m = airq.copy()
#forward filling so2,no2, and rspm
airq_m['so2'] = airq_m['so2'].fillna(method="ffill")
airq_m['no2'] = airq_m['no2'].fillna(method="ffill")
airq_m['rspm'] = airq_m['rspm'].fillna(method="bfill")

#print(airq_m.isnull().sum())

airq_m = airq_m.drop(columns=['stn_code', 'sampling_date', 'agency', 'location_monitoring_station', 'pm2_5', 'spm'])
airq_m = airq_m.dropna()
airq_m = airq_m[~airq_m['type'].isin(['RIRUO', 'Sensitive'])]
#print(airq_m.info())
airq_m.isnull().sum()
```
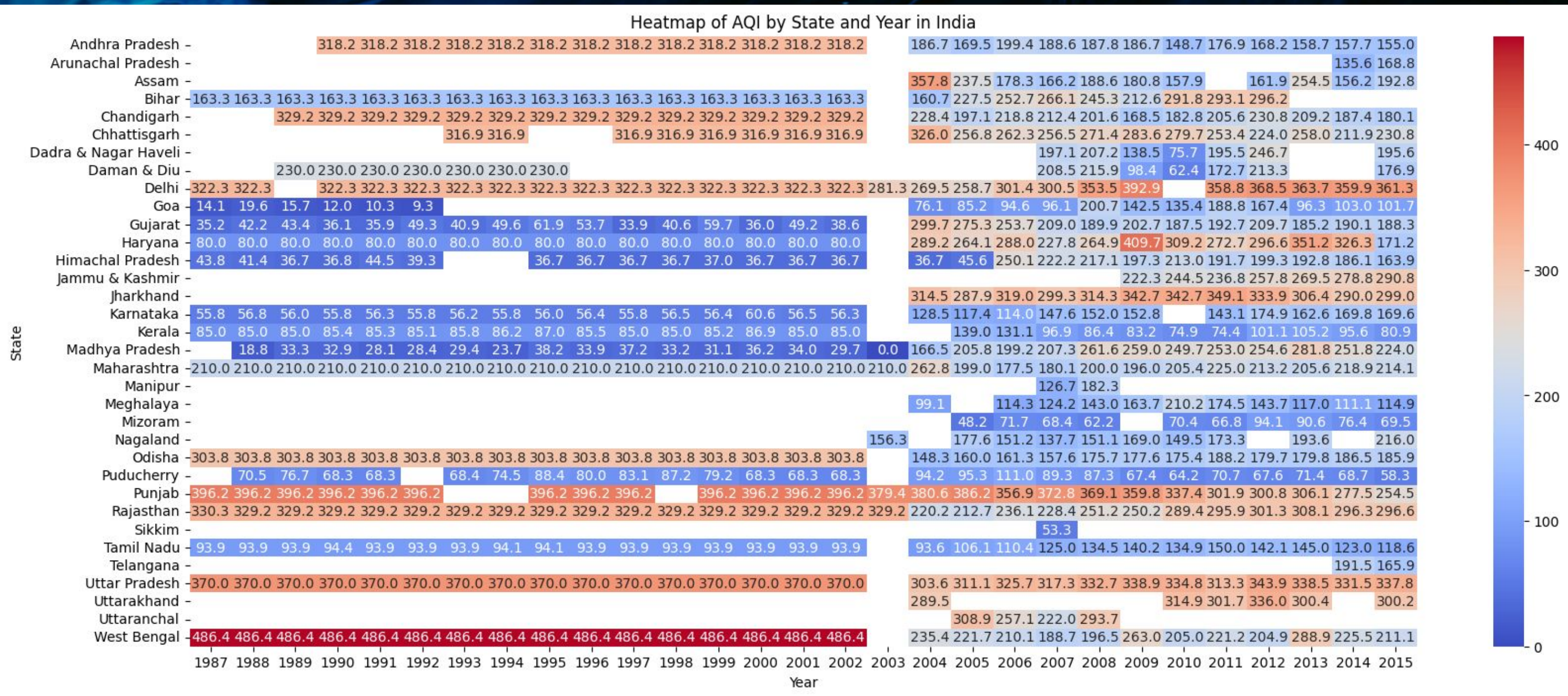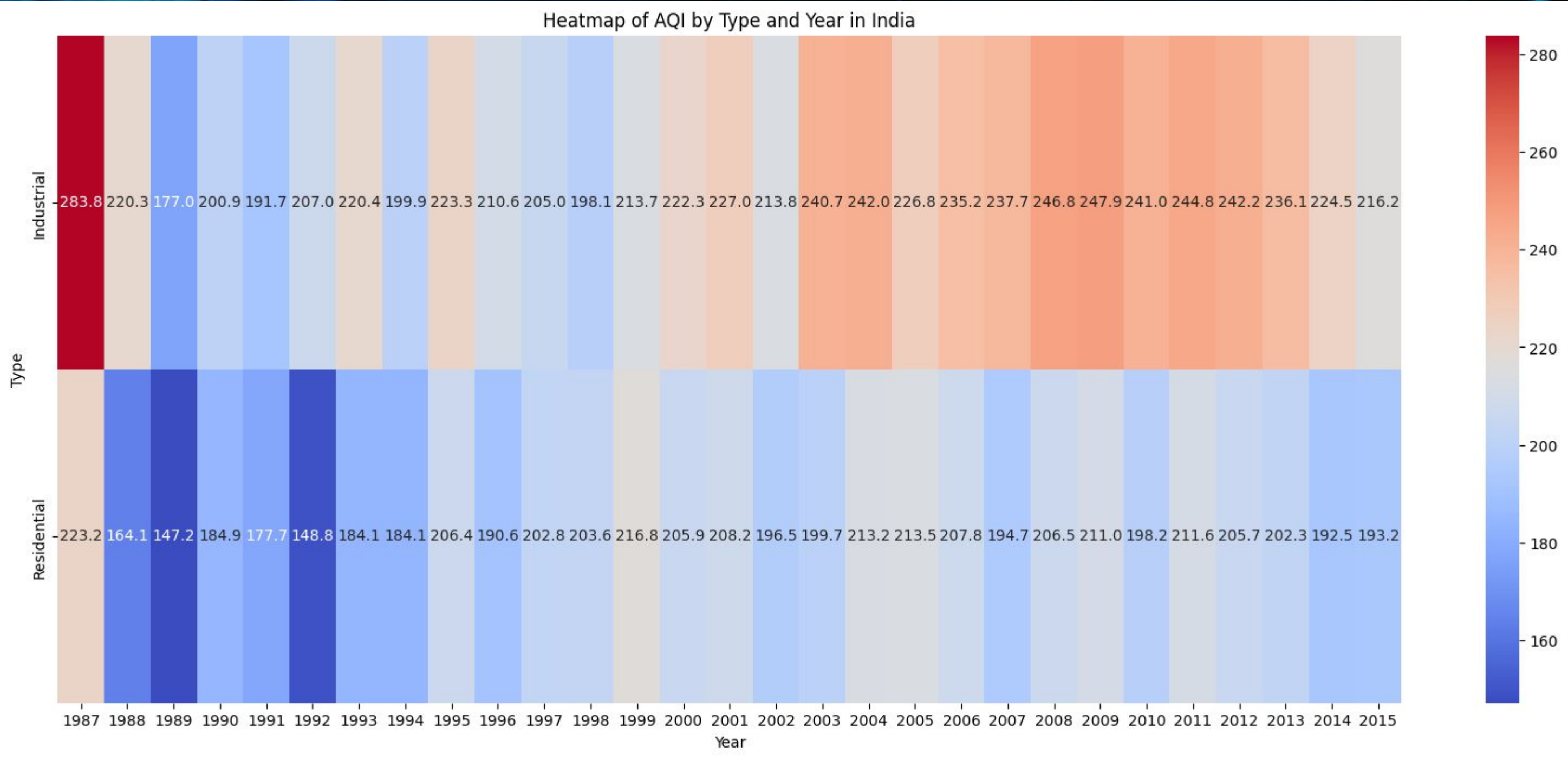
```
state        0
location     0
type         0
so2          0
no2          0
rspm         0
date         0
year         0
dtype: int64
```

Heatmap of AQI by State and Year in India

**Air Quality Index (AQI) trends for India based on types- industrial and residential.**



Heatmap of AQI by Type and Year in India

**These policies aim to regulate emissions and improve air quality over the years:**

## Industrial Pollution Policies

- National Emission Standards: Set limits for pollutants from industries, requiring pollution control technologies.
- Perform, Achieve, and Trade (PAT) Scheme: Encourages energy efficiency in industries through a trading system of energy savings certificates.

## Residential Pollution Policies

1. Pradhan Mantri Ujjwala Yojana (PMUY): Provides LPG connections to households to replace polluting fuels.
2. National Biomass Cookstoves Initiative: Promotes cleaner cookstoves to reduce indoor air pollution.

## Overarching Policies Affecting Both Sectors

- National Clean Air Programme (NCAP): Introduced in 2019, the NCAP aims to reduce particulate matter (PM10 and PM2.5) levels by 20-30% by 2024 compared to 2017 levels in more than 100 non-attainment cities identified for poor air quality.
- Smart Cities Mission: This initiative indirectly contributes to reducing urban pollution by promoting sustainable and smart solutions in urban development, including waste management and sustainable urban transport

# Introduction to SOi, NOi, RSPMi

```python
def cal_RSPMI(rspm):
    rpi = 0
    if rspm <= 30:
        rpi = rspm * 50 / 30
    elif rspm > 30 and rspm <= 60:
        rpi = 50 + (rspm - 30) * 50 / 30
    elif rspm > 60 and rspm <= 90:
        rpi = 100 + (rspm - 60) * 100 / 30
    elif rspm > 90 and rspm <= 120:
        rpi = 200 + (rspm - 90) * 100 / 30
    elif rspm > 120 and rspm <= 250:
        rpi = 300 + (rspm - 120) * 100 / 130
    else:
        rpi = 400 + (rspm - 250) * 100 / 130
    return rpi
airq_m['Rpi'] = airq_m['rspm'].apply(cal_RSPMI)
data = airq_m[['rspm', 'Rpi']]
data.sample(20)
```

```python
def cal_SOi(so2):
    si=0
    if (so2<=40):
        si= so2*(50/40)
    elif (so2>40 and so2<=80):
        si= 50+(so2-40)*(50/40)
    elif (so2>80 and so2<=380):
        si= 100+(so2-80)*(100/300)
    elif (so2>380 and so2<=800):
        si= 200+(so2-380)*(100/420)
    elif (so2>800 and so2<=1600):
        si= 300+(so2-800)*(100/800)
    elif (so2>1600):
        si= 400+(so2-1600)*(100/800)
    return si
airq_m['SOi']=airq_m['so2'].apply(cal_SOi)
data= airq_m[['so2','SOi']]
data.sample(20)
```

```python
def cal_Noi(no2):
    ni=0
    if(no2<=40):
        ni= no2*50/40
    elif(no2>40 and no2<=80):
        ni= 50+(no2-40)*(50/40)
    elif(no2>80 and no2<=180):
        ni= 100+(no2-80)*(100/100)
    elif(no2>180 and no2<=280):
        ni= 200+(no2-180)*(100/100)
    elif(no2>280 and no2<=400):
        ni= 300+(no2-280)*(100/120)
    else:
        ni= 400+(no2-400)*(100/120)
    return ni
airq_m['Noi']=airq_m['no2'].apply(cal_Noi)
data= airq_m[['no2','Noi']]
data.sample(20)
```

# Calculation of AQI - Air Quality Index

```python
def cal_aqi(si, ni, rpi):
    aqi = 0
    if (si >= ni and si >= rpi):
        aqi = si
    elif (ni >= si and ni >= rpi):
        aqi = ni
    elif (rpi >= si and rpi >= ni):
        aqi = rpi
    return aqi

airq_m['AQI'] = airq_m.apply(lambda x: cal_aqi(x['SOi'], x['Noi'], x['Rpi']), axis=1)
data = airq_m[['state', 'SOi', 'Noi', 'Rpi', 'AQI']]
data.head()
```
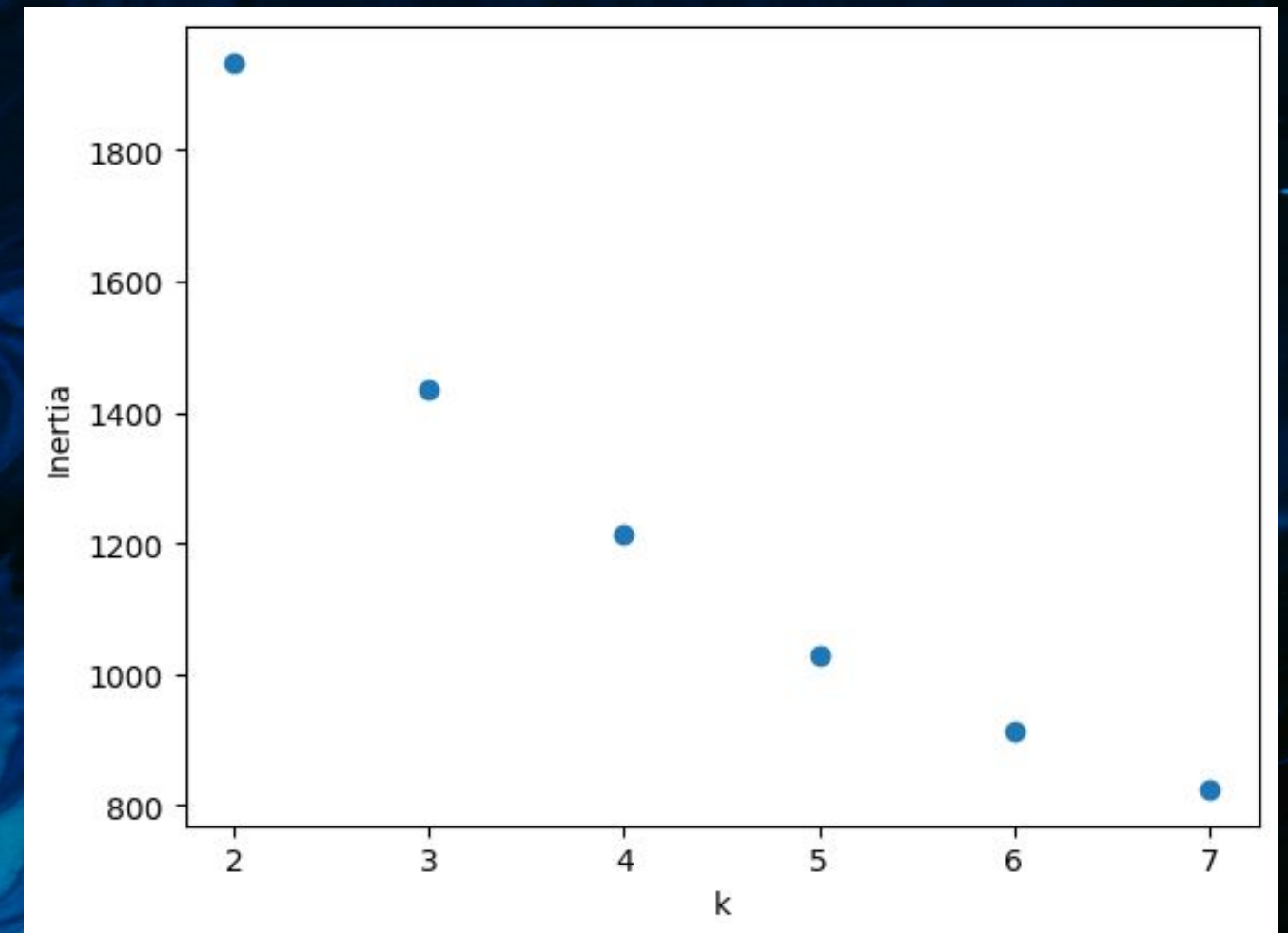
```python
def AQI_Range(x):
    if x<=50:
        return "Good"
    elif x>50 and x<=100:
        return "Moderate"
    elif x>100 and x<=200:
        return "Poor"
    elif x>200 and x<=300:
        return "Unhealthy"
    elif x>300 and x<=400:
        return "Very unhealthy"
    elif x>400:
        return "Hazardous"

airq_m['AQI_Range'] = airq_m['AQI'] .apply(AQI_Range)
airq_m.tail(20)
```
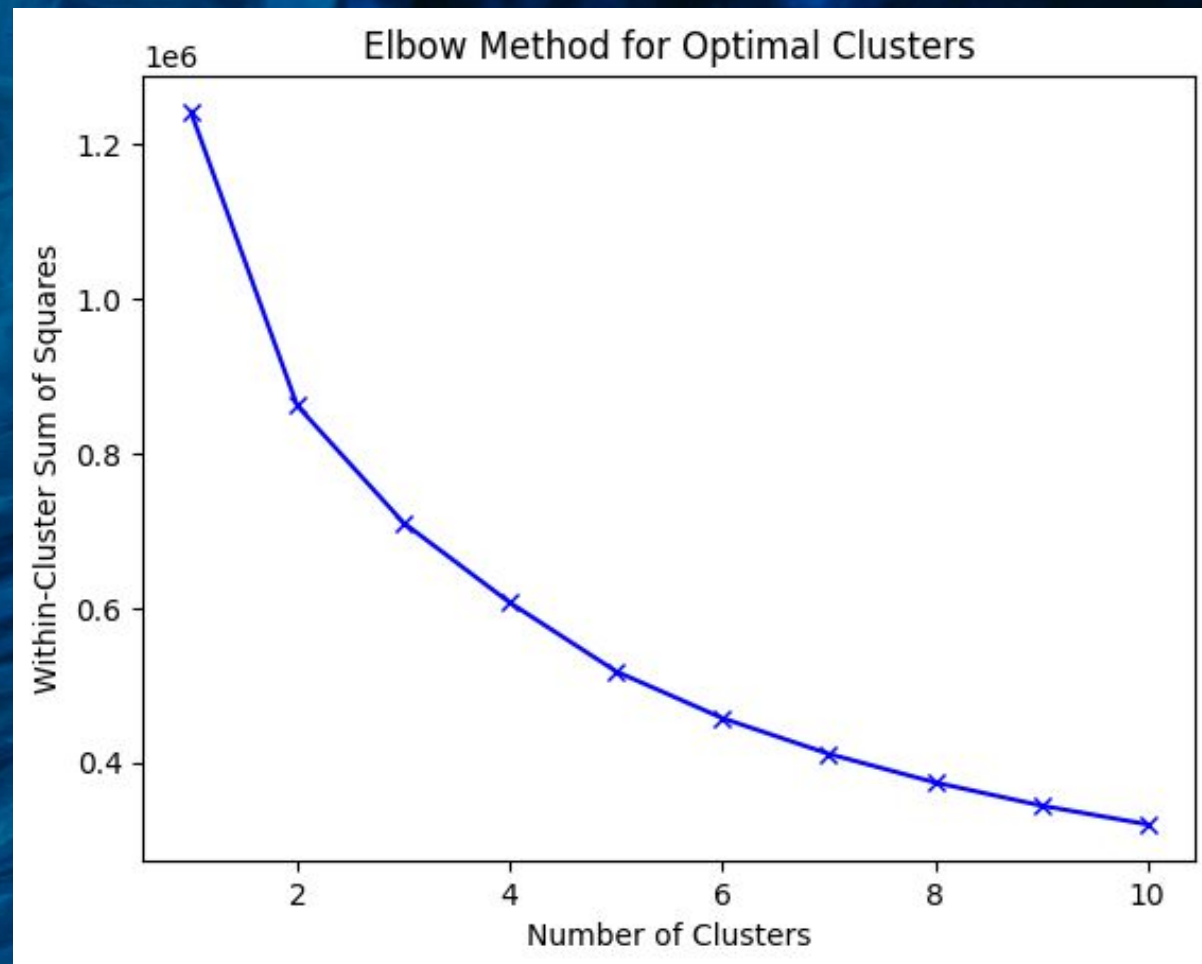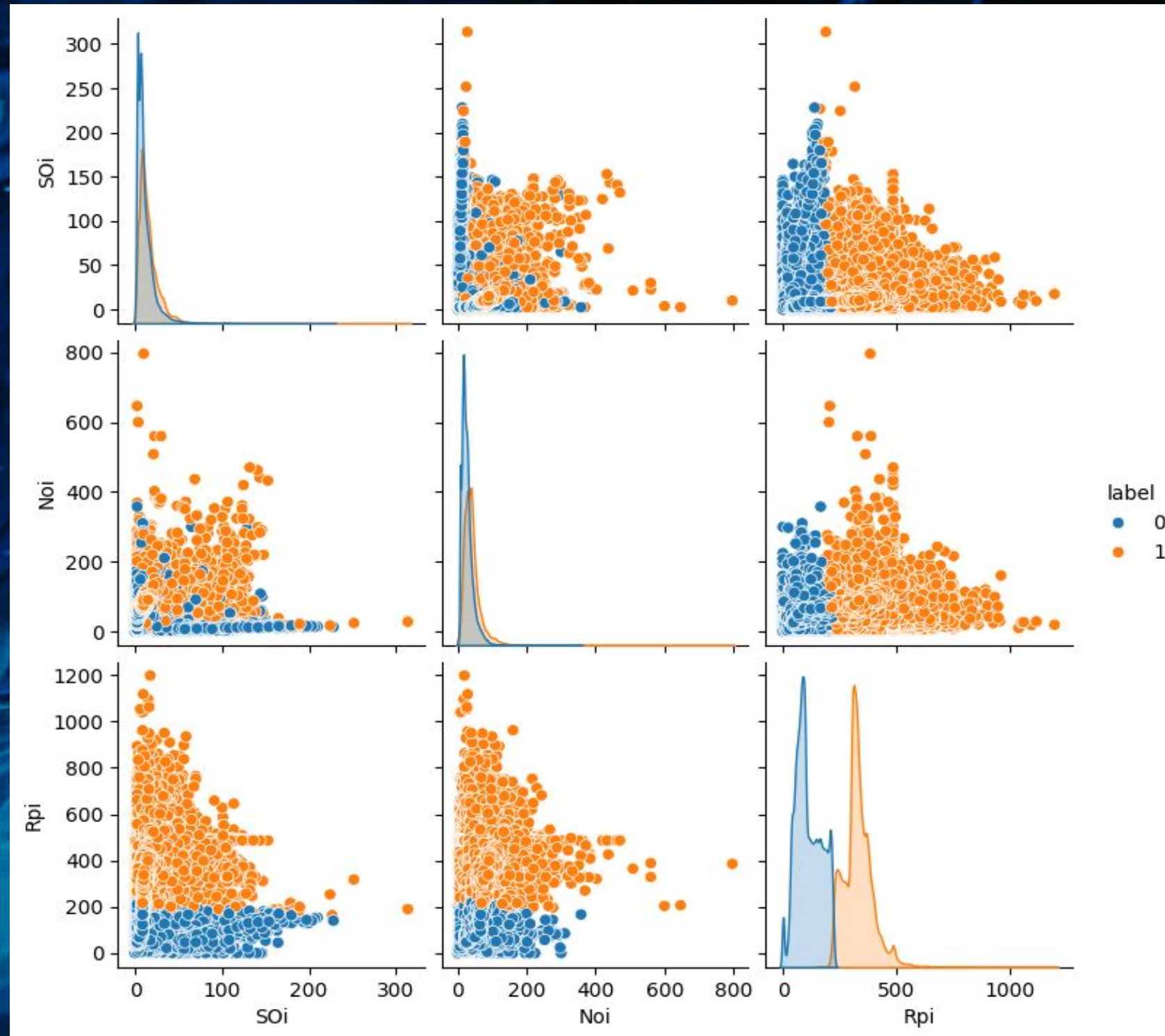
# K-means Clustering:

```
k = 2
kmeans = cluster.KMeans(n_clusters=k)
kmeans.fit(X_scaled)

labels = kmeans.labels_
centroids = kmeans.cluster_centers_
inertia = kmeans.inertia_

df['label'] = labels
```
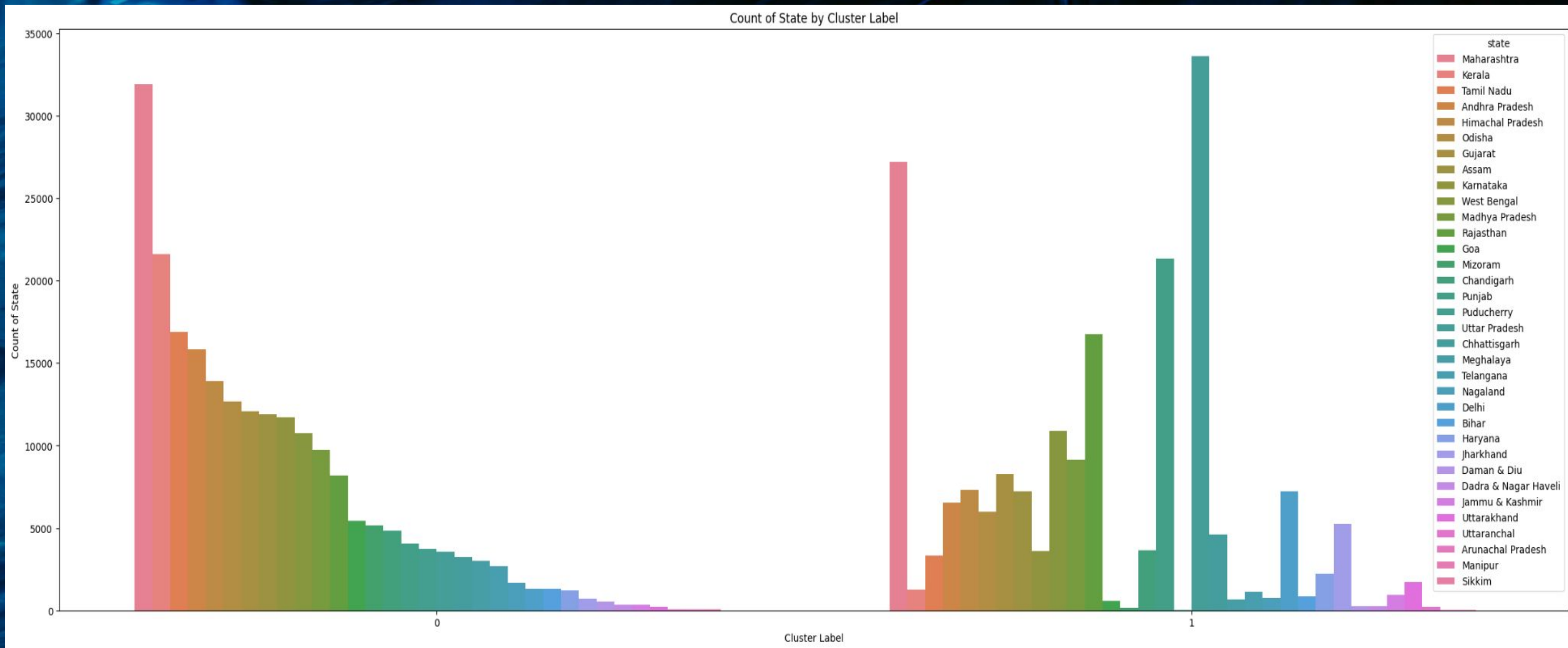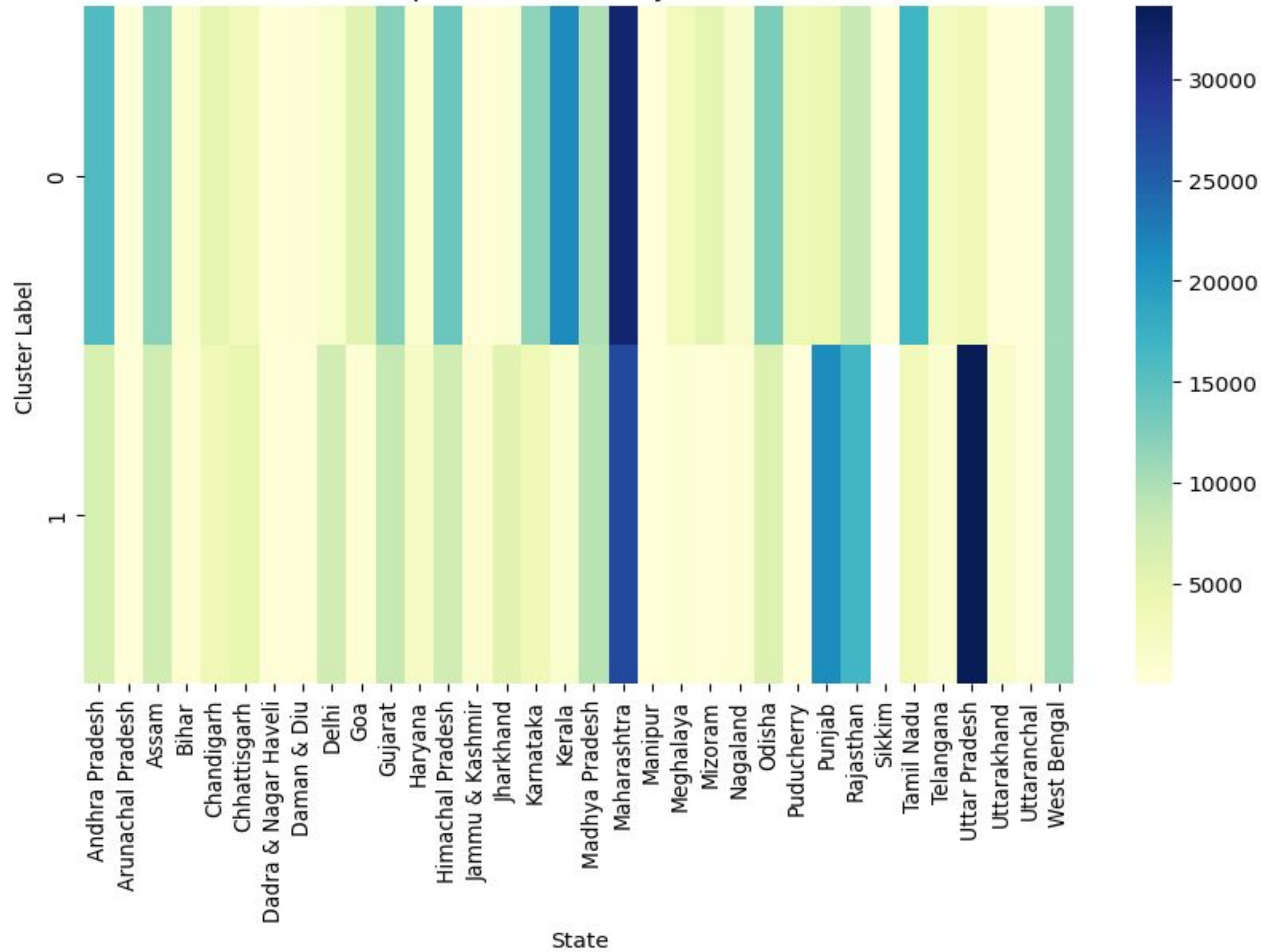
K-means clustering on SOI, NOI, RPI using two labels:

# Count of state by Cluster Label Graph:

Heatmap of State Count by Cluster Label

# ML models used to Train

```
[ ]  RF=RandomForestRegressor().fit(X_train,Y_train)
```

```
▼ LinearRegression
LinearRegression()
```

```
#predicting train
train_preds1=RF.predict(X_tra
#predicting on test
test_preds1=RF.predict(X_test
```

```
▼ DecisionTreeRegressor
DecisionTreeRegressor()
```

```
#predicting train
train_pred=model.predict(
#predicting on test
test_pred=model.predict(X
```

```
#predicting train
train_preds=DT.predict(X_train)
#predicting on test
test_preds=DT.predict(X_test)
```

```
[ ] RMSE_train=(np.sqrt(metri
    RMSE_test=(np.sqrt(metrics
    print("RMSE TrainingData
    print("RMSE TestData = ",
    print('-'*50)
    print('RSquared value on
    print('RSquared value on
```

```
[ ] RMSE_train=(np.sqrt(metrics.me
    RMSE_test=(np.sqrt(metrics.mea
    print("RMSE TrainingData = ",
    print("RMSE TestData = ",str(
    print('-'*50)
    print('RSquared value on trai
    print('RSquared value on test
```

```
RMSE TrainingData =  0.4237155
RMSE TestData =  1.09287093397
----------------------------
RSquared value on train: 0.999
RSquared value on test: 0.9999
```

```
[ ] RMSE_train=(np.sqrt(metrics.mean_squared_error(Y_train,train_preds)))
    RMSE_test=(np.sqrt(metrics.mean_squared_error(Y_test,test_preds)))
    print("RMSE TrainingData = ",str(RMSE_train))
    print("RMSE TestData = ",str(RMSE_test))
    print('-'*50)
    print('RSquared value on train:',DT.score(X_train, Y_train))
    print('RSquared value on test:',DT.score(X_test, Y_test))
```

```
RMSE TrainingData =  5.26¢
RMSE TestData =  5.103765276394435
------------------------------
RSquared value on train: 0.9981916296381512
RSquared value on test: 0.9983033594660915
```

```
RMSE TrainingData =  2.447578149985567e-12
RMSE TestData =  1.2083953346584535
--------------------------------------------------
RSquared value on train: 1.0
RSquared value on test: 0.9999048899231755
```

# Training the ML model

```python
features = ['SOi', 'Noi', 'Rpi', 'state', 'year','AQI']
target = 'AQI'
pred = airq_m[features].copy()

pred = pd.get_dummies(pred, columns=['state'], prefix='state')  # Prefix to identify the state columns

# Standardize the numerical features
scaler = StandardScaler()
numerical_features = ['SOi', 'Noi', 'Rpi', 'year']  # The features to standardize
pred[numerical_features] = scaler.fit_transform(pred[numerical_features])


# Selecting relevant columns for modeling
# The original 'state' column is now replaced with one-hot encoded columns.
# We need to redefine features to include these new columns.
# You can get the list of one-hot encoded 'state' columns.
state_columns = [col for col in pred.columns if col.startswith('state_')]

# Redefine features to include the correct column names
features = ['SOi', 'Noi', 'Rpi', 'year'] + state_columns

# Now 'target' should refer to the original 'AQI' column
target = 'AQI'  # The target variable

# Split the data into training and test sets
X_train, X_test, y_train, y_test = train_test_split(
    pred[features], pred[target], test_size=0.2, random_state=42
)
```

```python
31]  # Select and train a model (e.g., Random Forest)
     model = RandomForestRegressor(n_estimators=100, random_state=42)
     model.fit(X_train, y_train)
```
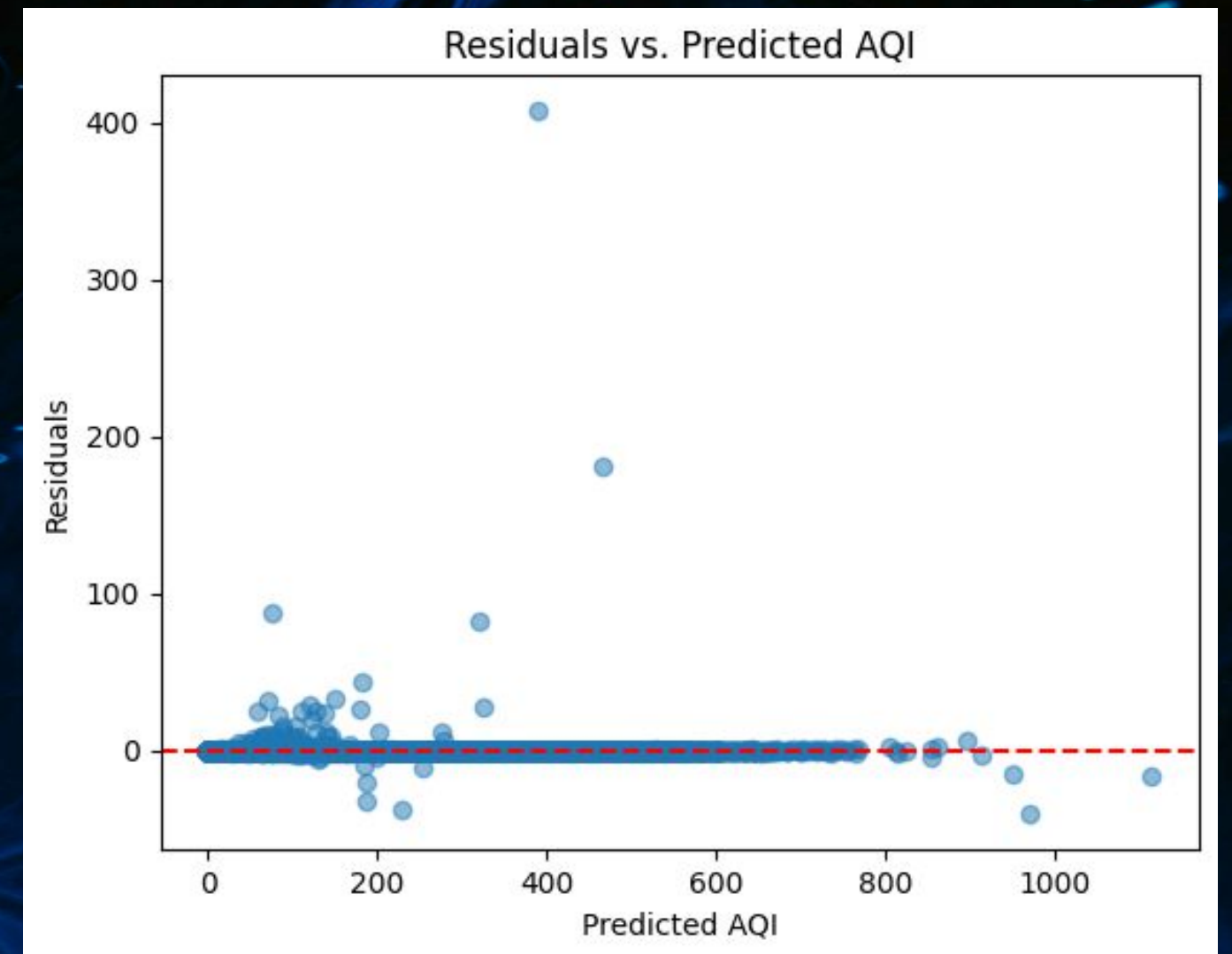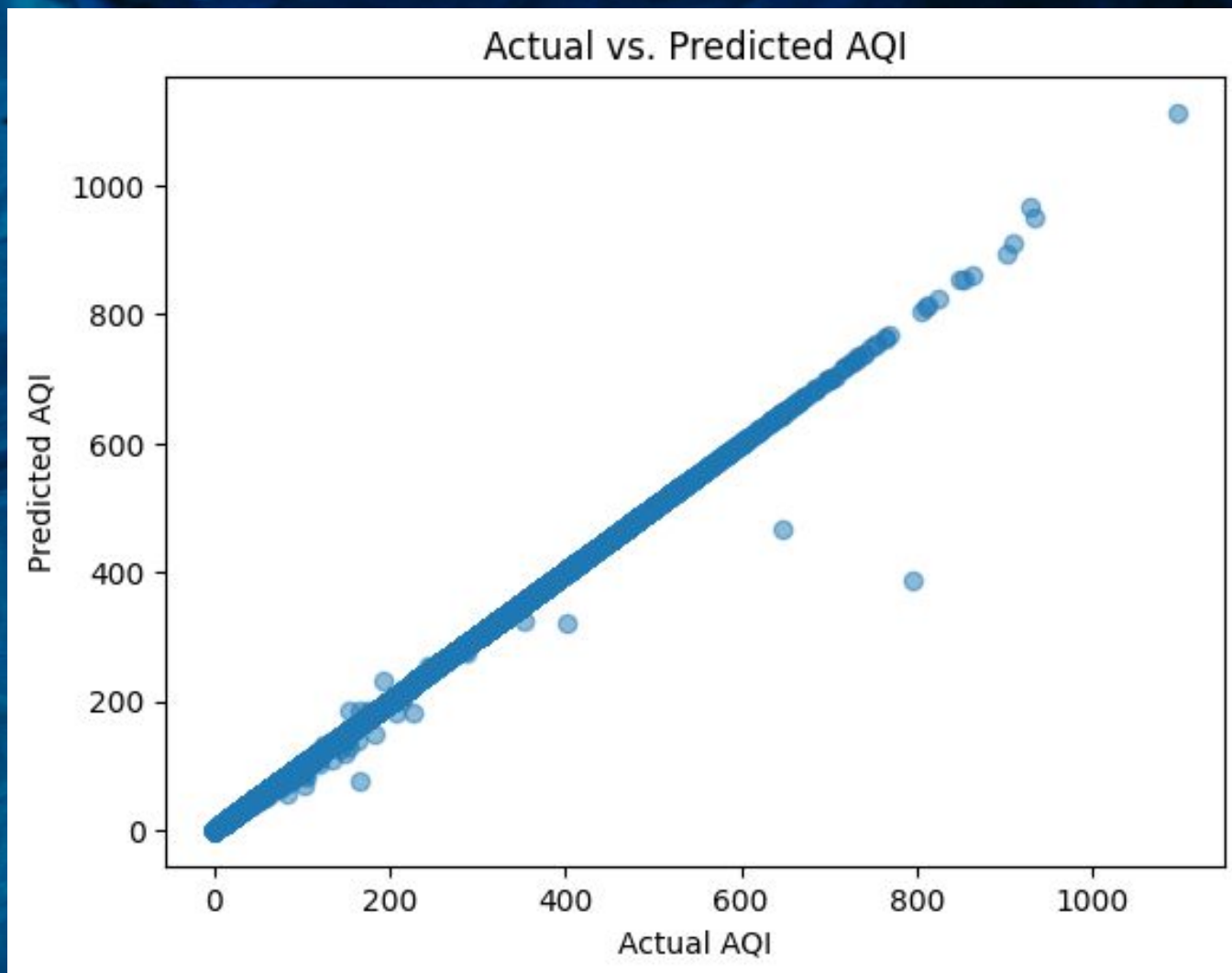
```
▼              RandomForestRegressor
RandomForestRegressor(random_state=42)
```

```python
# Predict on the test set
y_pred = model.predict(X_test)

# Calculate performance metrics
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

print("MSE:", mse)
print("R-squared:", r2)
```

```
MSE: 2.7960007282483312
R-squared: 0.9998179752769996
```

Actual vs. Predicted AQI — Residuals vs. Predicted AQI

```python
from sklearn.model_selection import cross_val_score

# Perform 5-fold cross-validation on the model
cross_val_scores = cross_val_score(model, X_train, y_train, cv=5, scoring='neg_mean_squared_error')

# Convert to positive values and calculate the mean
mean_mse = -cross_val_scores.mean()

print("Mean MSE from Cross-Validation:", mean_mse)
```

Mean MSE from Cross-Validation: 1.35019912265623

Developing the future data for the model to work on:

```python
# Define the states and year range
states = ["Maharashtra", "Uttar Pradesh", "Andhra Pradesh", "Punjab", "Rajasthan", "Kerala",
          "Himachal Pradesh", "West Bengal", "Gujarat", "Tamil Nadu", "Madhya Pradesh", "Assam",
          "Odisha", "Karnataka", "Delhi", "Chandigarh", "Chhattisgarh", "Goa", "Jharkhand",
          "Mizoram", "Telangana", "Meghalaya", "Puducherry", "Haryana", "Nagaland", "Bihar",
          "Uttarakhand", "Jammu & Kashmir", "Daman & Diu", "Dadra & Nagar Haveli", "Uttaranchal",
          "Arunachal Pradesh", "Manipur", "Sikkim"]

years = list(range(2016, 2025))

# Create a DataFrame with all state-year combinations
future_data = pd.DataFrame(list(itertools.product(states, years)), columns=["state", "year"])

# Generate random values within derived ranges
np.random.seed(42)  # For reproducibility
future_data["SOi"] = np.random.uniform(soi_min, soi_max, len(future_data))
future_data["Noi"] = np.random.uniform(noi_min, noi_max, len(future_data))
future_data["Rpi"] = np.random.uniform(rpi_min, rpi_max, len(future_data))
```
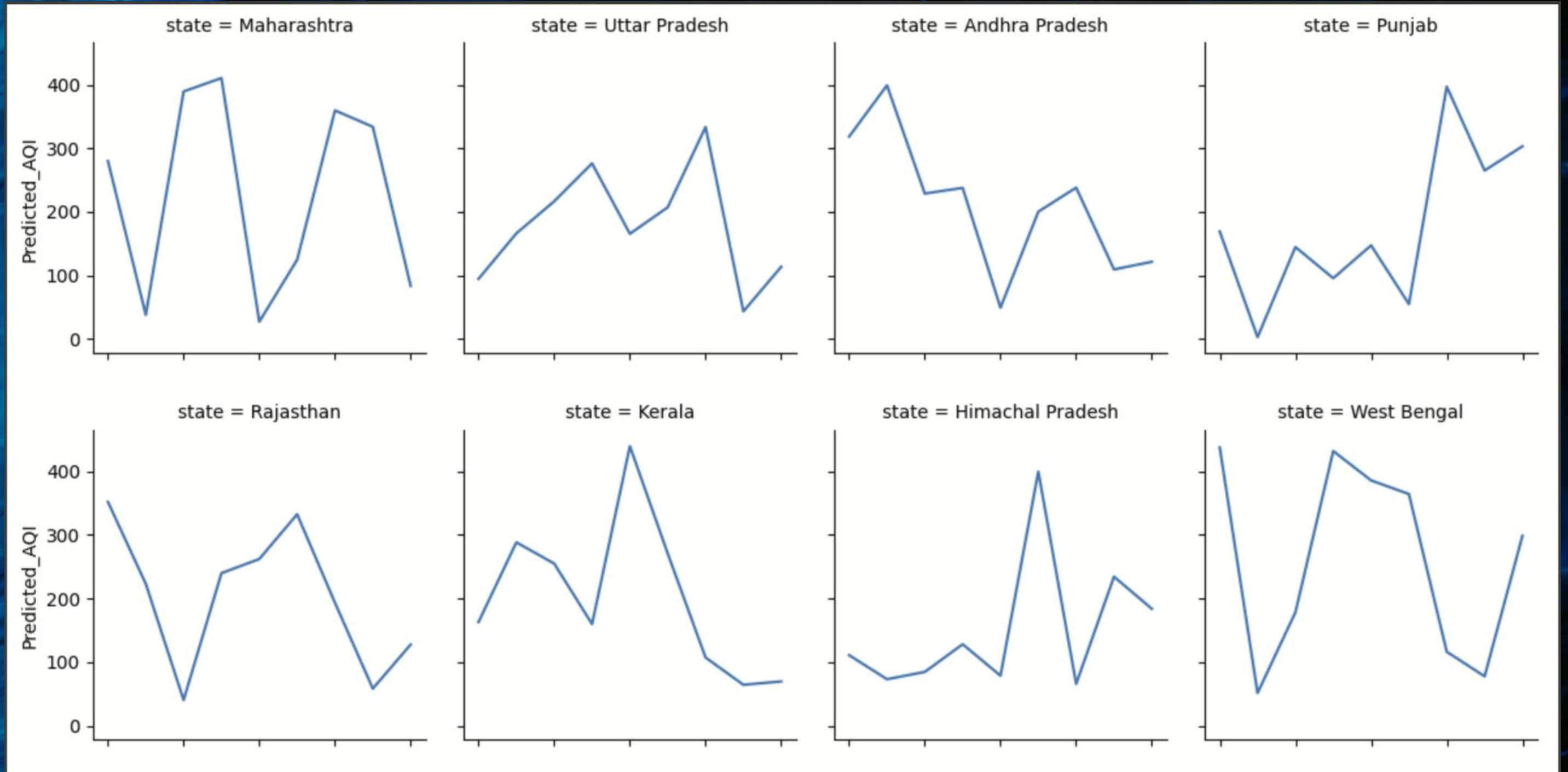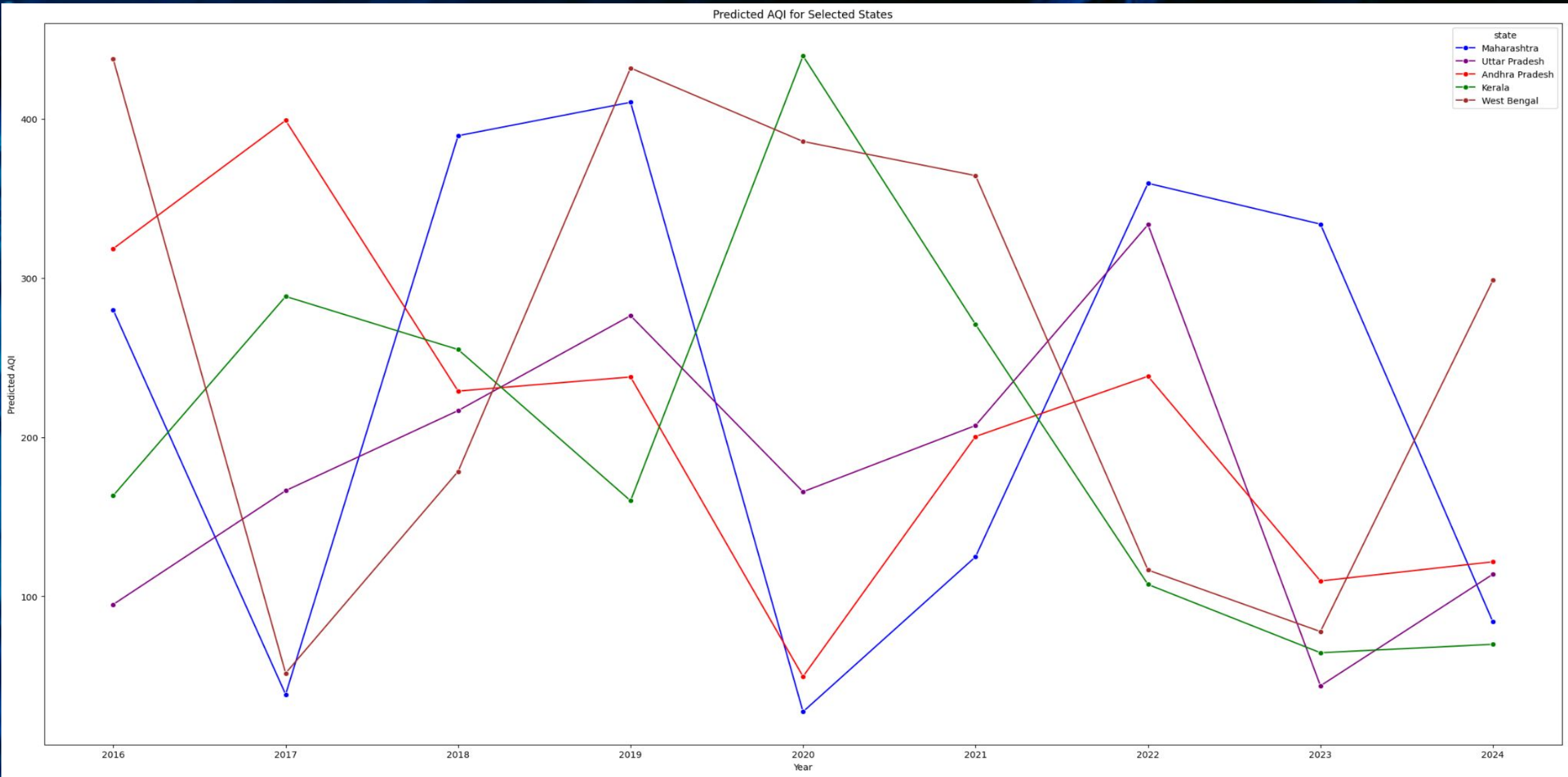
```python
[39] # Predict future AQI levels
     future_aqi_predictions = model.predict(future_data_encoded)

     # Add the predictions to the future data
     future_data["Predicted_AQI"] = future_aqi_predictions
```

Predicted Values for Individual States from 2016-2024

Predicted AQI for Selected States

state
Maharashtra
Uttar Pradesh
Andhra Pradesh
Kerala
West Bengal

Real-time Andhra Pradesh
Most polluted city ranking          Orange

| # | CITY | US AQI |
|---|------|--------|
| 1 | Anantapur | 113 |
| 2 | Chittoor | 99 |

Real-time Uttar Pradesh
Most polluted city ranking          Purple

| # | CITY | US AQI |
|---|------|--------|
| 1 | Baghpat | 118 |
| 2 | Firozabad | 92 |

LIVE AQI CITY RANKING

Real-time West Bengal
Most polluted city ranking          Brown
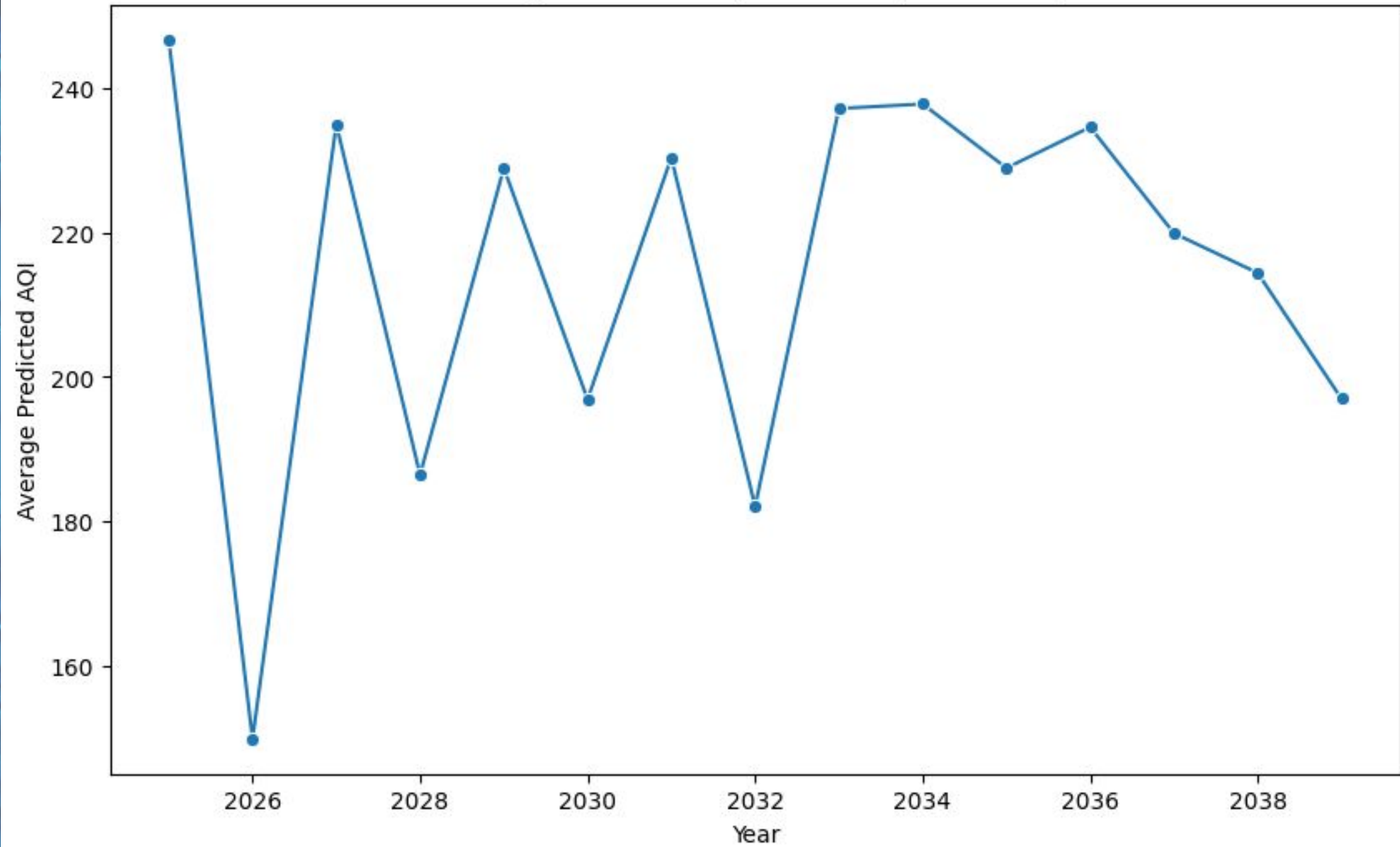
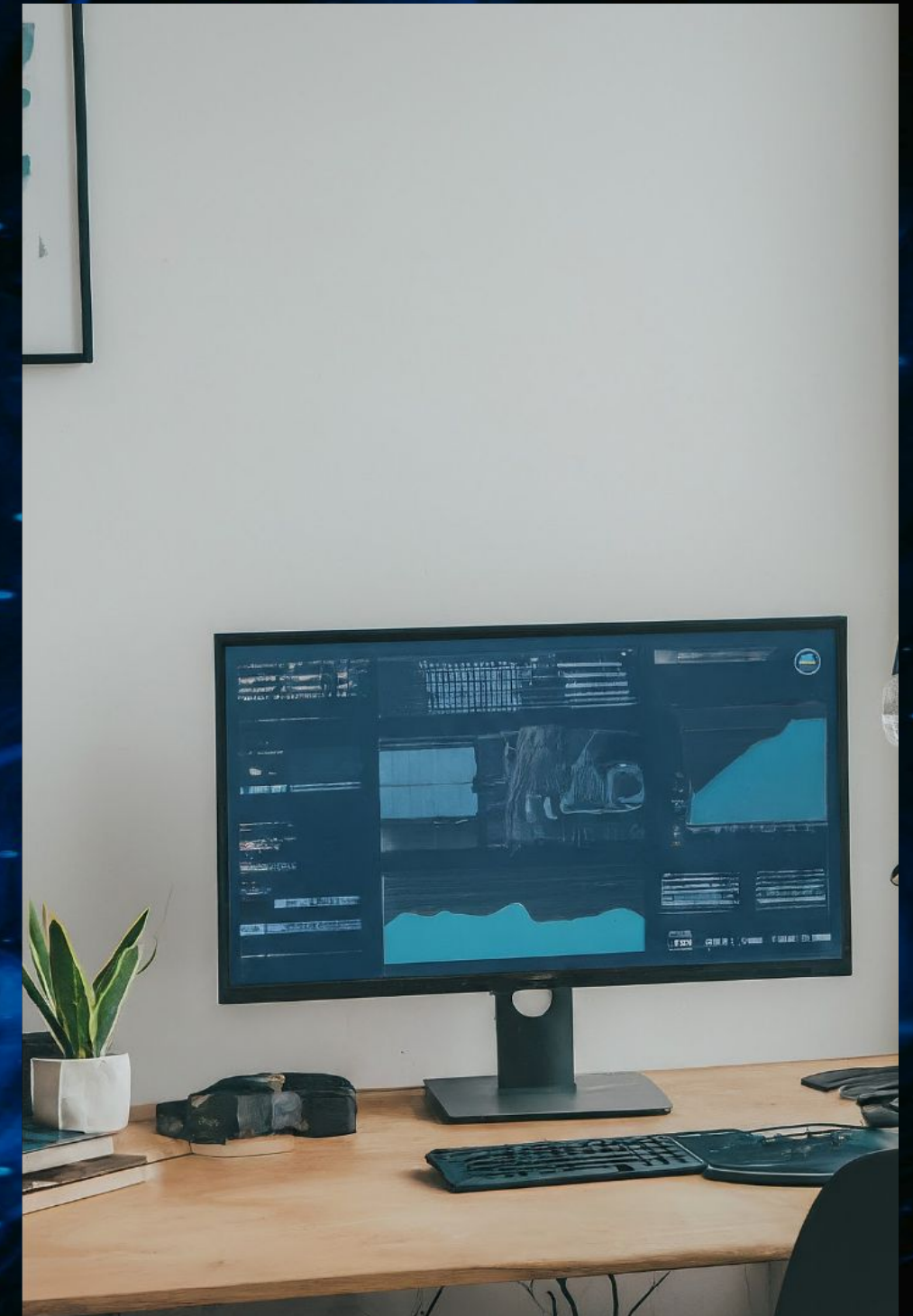| # | CITY | US AQI |
|---|------|--------|
| 1 | Barakpur | 154 |
| 2 | Kolkata | 154 |

Average AQI for India (1987-2024)

Average Predicted AQI for India (2025-2039)

# Conclusion and Future Directions

In conclusion, the application of data science and machine learning algorithms has proven to be a powerful approach for effective air quality control. By leveraging predictive modeling, anomaly detection, and optimization techniques, we can better monitor, forecast, and optimize air quality. **Looking ahead**, integrating air quality data with IoT sensor networks will enable real-time, hyper-local insights to drive even more precision in air quality management.

# Thank You