# Air Quality & Pollution

# Introduction

Did you know that air pollution causes 7 million premature deaths annually worldwide? Behind these statistics are real people whose lives could be saved through better air quality management.

Air quality has a direct impact on:

- Public health (respiratory diseases, heart problems)
- Economic productivity (healthcare costs, work absences)
- Environmental sustainability (ecosystem damage, climate effects)
- Quality of life (outdoor activities, community well-being)

# TARGET AUDIENCE

Primary Audiences
- **Environmental agencies** responsible for pollution monitoring and regulation
- **Municipal governments** making urban planning and public health decisions
- **Public health officials** developing intervention strategies for vulnerable populations
- **Industrial compliance managers** seeking to minimize environmental impact

Secondary Audiences
- **Community organizations** advocating for environmental justice
- **Research institutions** studying pollution effects on health outcomes
- **Educational institutions** teaching about environmental science and data analytics
- **General public** interested in understanding local air quality conditions

# Data Collection & Preparation

Our analysis leverages a comprehensive synthetic dataset comprising 5,000 records across 10 key environmental and socioeconomic variables:

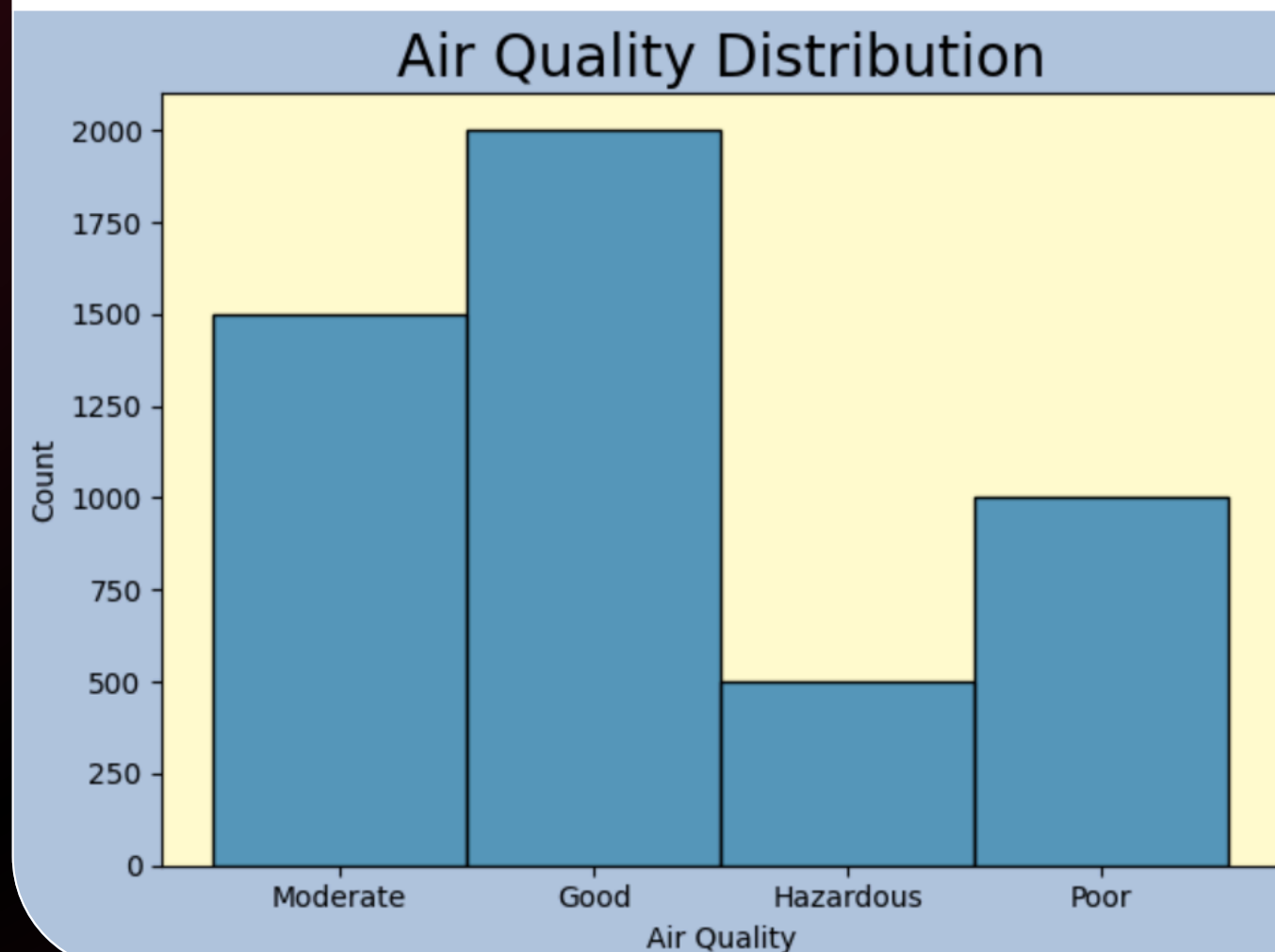| Feature Type | Variables |
|---|---|
| Environmental Factors | Temperature, Humidity, PM2.5, PM10, NO2, SO2, CO |
| Socioeconomic Factors | Proximity to Industrial Areas, Population Density |
| Target Variable | Air Quality Classification |

# Data Quality Assurance

Before analysis, we implemented a robust preprocessing pipeline:
- Standardized measurements across all sensors
- Identified and handled outliers to prevent model distortion
- Validated data integrity through cross-referencing with established monitoring stations
- Normalized features to improve model convergence and performance

# the distribution of the target variable classes

```python
sns.histplot(data=data,x='Air Quality')
plt.title('Air Quality Distribution',size=20)

plt.tight_layout()
plt.gcf().patch.set_facecolor('lightsteelblue')
plt.gca().set_facecolor('lemonchiffon')
plt.show()
```
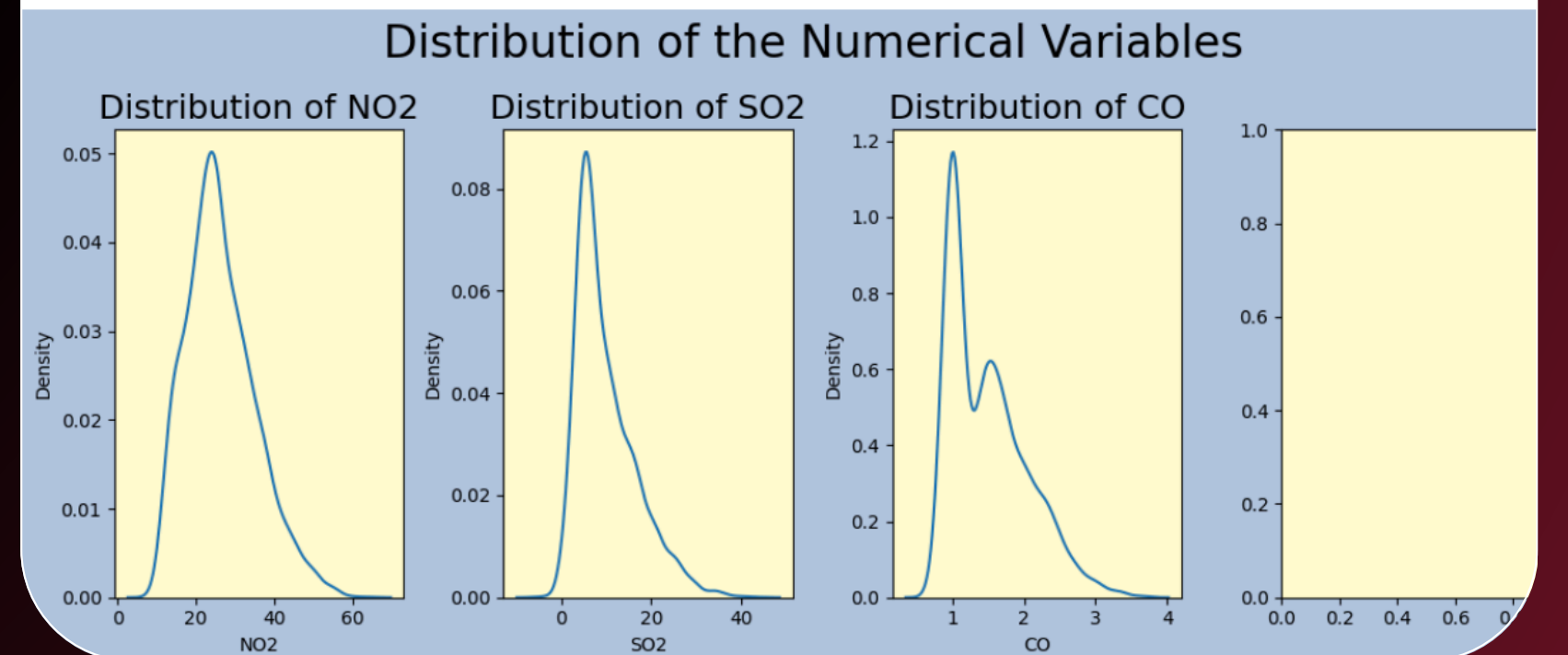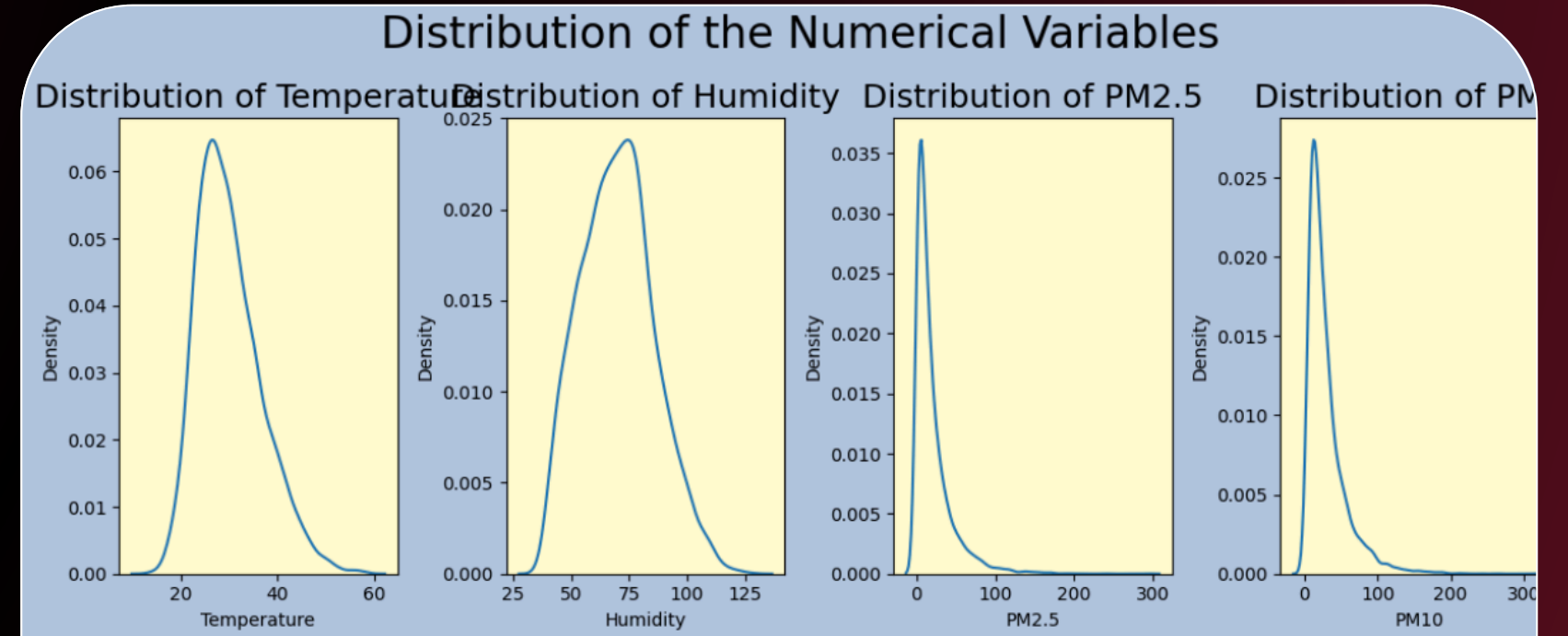
# Distribution of the Numerical Variables

```python
features = ['Temperature','Humidity','PM2.5','PM10','NO2','SO2','CO', 'Proximity_
as','Population_Density']

for i in range(2):
    fig,(ax1,ax2,ax3,ax4) = plt.subplots(ncols=4,figsize=(12,5))
    ax1 = sns.distplot(data[features[i*4]],ax=ax1,hist=False)
    ax1.set_title('Distribution of '+str(features[i*4]),fontsize=18)
    ax1.set_facecolor('lemonchiffon')
    ax2 = sns.distplot(data[features[i*4+1]],ax=ax2,hist=False)
    ax2.set_title('Distribution of '+str(features[i*4+1]),fontsize=18)
    ax2.set_facecolor('lemonchiffon')
    ax3 = sns.distplot(data[features[i*4+2]],ax=ax3,hist=False)
    ax3.set_title('Distribution of '+str(features[i*4+2]),fontsize=18)
    ax3.set_facecolor('lemonchiffon')
    if i < 1:
        ax4 = sns.distplot(data[features[i*4+3]],ax=ax4,hist=False)
        ax4.set_title('Distribution of '+str(features[i*4+3]),fontsize=18)
        ax4.set_facecolor('lemonchiffon')
    else:
        ax4.set_facecolor('lemonchiffon')

    fig.suptitle("Distribution of the Numerical Variables",fontsize=24)

    plt.tight_layout()
    fig.set_facecolor('lightsteelblue')
```

# Outlier Detection & Handling

Outliers can significantly impact model performance, particularly for air quality monitoring where extreme events are critical to detect:

- **Methodology:** We employed the Interquartile Range (IQR) method to identify statistical anomalies
- **Verification:** Each outlier was cross-verified against known pollution events
- **Preservation strategy:** Extreme but valid readings were retained to ensure model robustness to pollution spikes
- **Transformation impact:** Removing statistical anomalies improved model generalization while maintaining sensitivity to unusual events
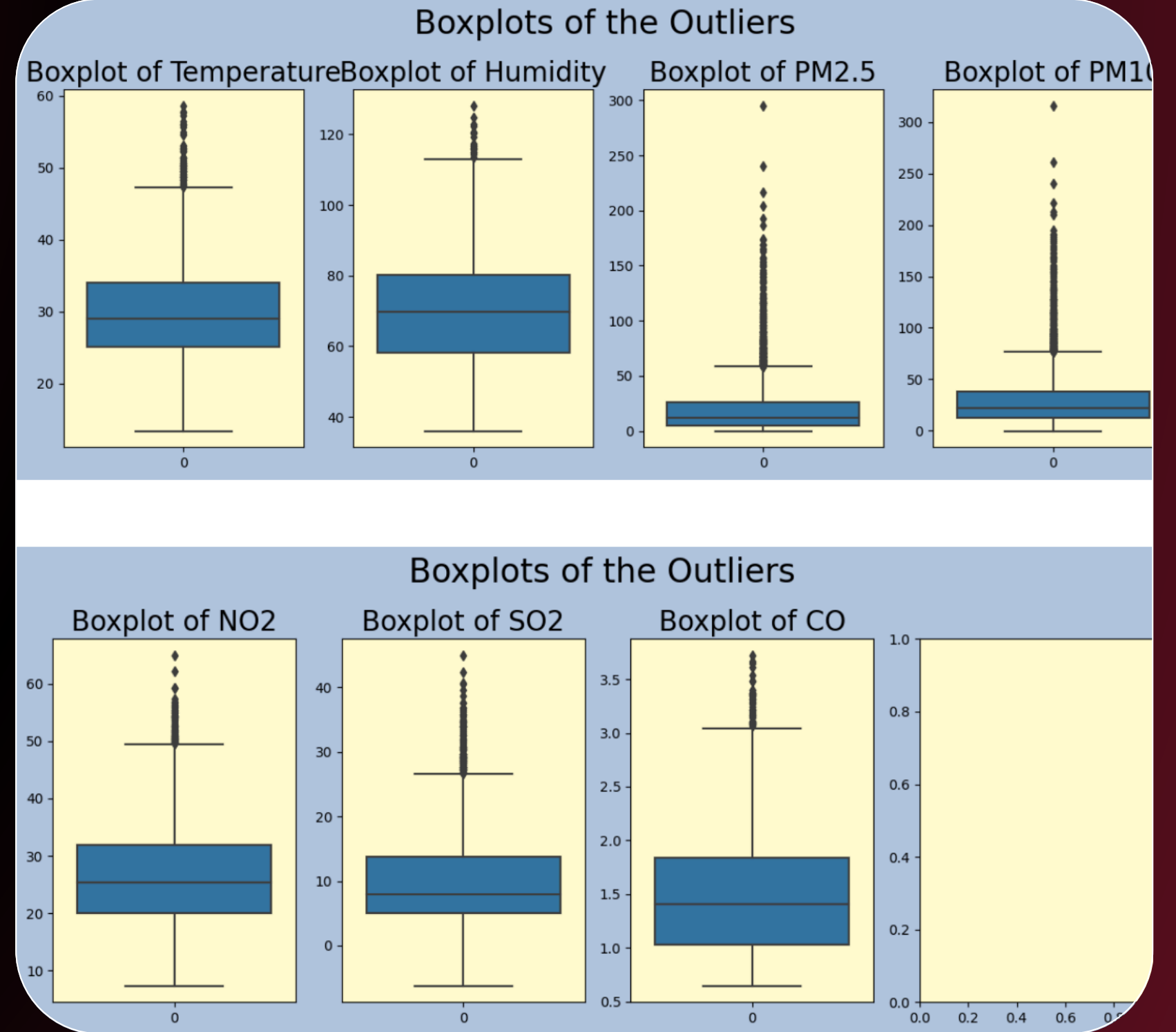
```python
for i in range(2):
    fig,(ax1,ax2,ax3,ax4) = plt.subplots(ncols=4,figsize=(12,5))
    ax1 = sns.boxplot(data[features[i*4]],ax=ax1)
    ax1.set_title('Boxplot of '+str(features[i*4]),fontsize=20)
    ax1.set_facecolor('lemonchiffon')
    ax2 = sns.boxplot(data[features[i*4+1]],ax=ax2)
    ax2.set_title('Boxplot of '+str(features[i*4+1]),fontsize=20)
    ax2.set_facecolor('lemonchiffon')
    ax3 = sns.boxplot(data[features[i*4+2]],ax=ax3)
    ax3.set_title('Boxplot of '+str(features[i*4+2]),fontsize=20)
    ax3.set_facecolor('lemonchiffon')
    if i < 1:
        ax4 = sns.boxplot(data[features[i*4+3]],ax=ax4)
        ax4.set_title('Boxplot of '+str(features[i*4+3]),fontsize=20)
        ax4.set_facecolor('lemonchiffon')
    else:
        ax4.set_facecolor('lemonchiffon')

    fig.suptitle("Boxplots of the Outliers",fontsize=24)

    plt.tight_layout()
    fig.set_facecolor('lightsteelblue')
```
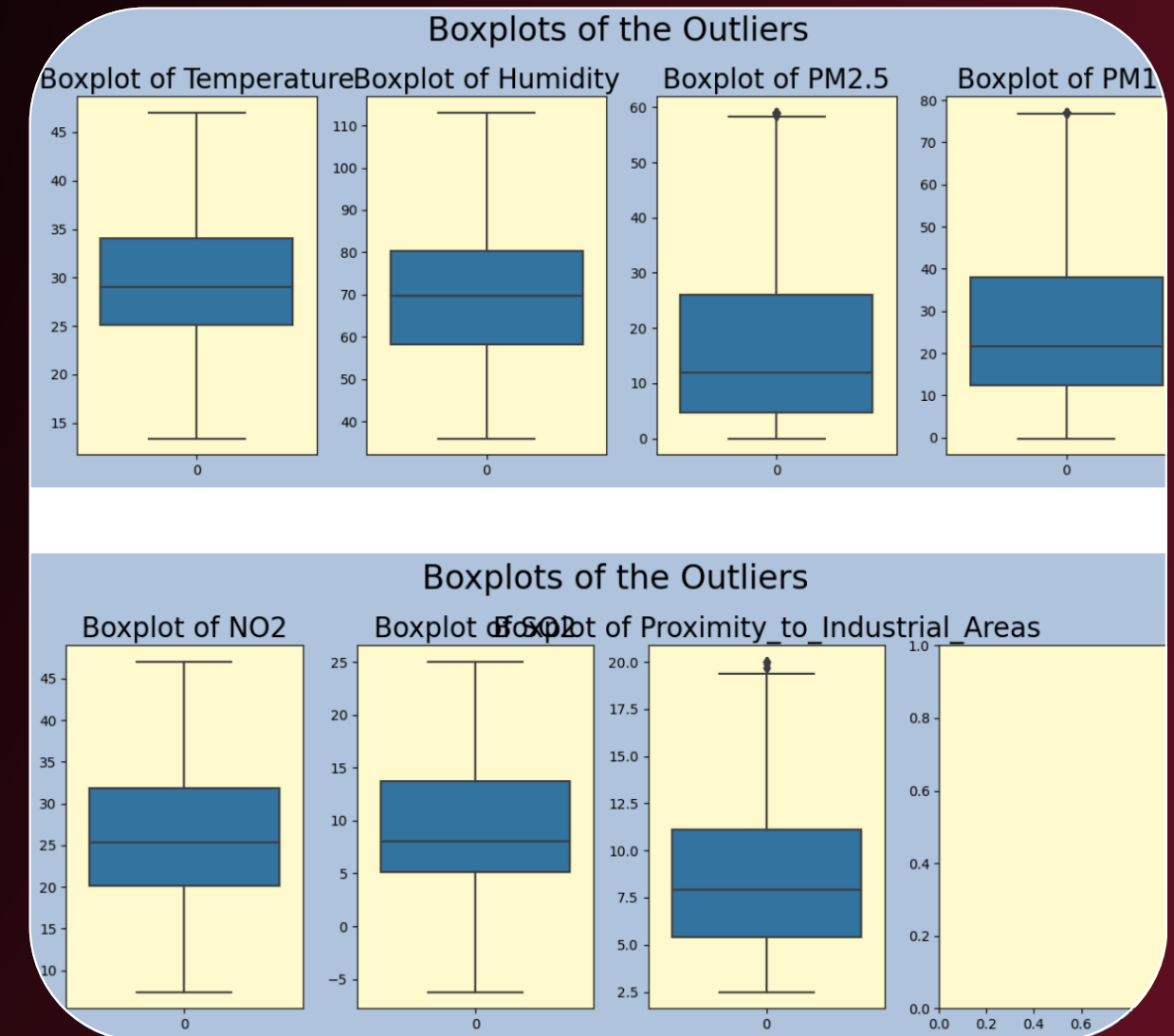
outliers

# After we remove outliers

```python
def outlier_imputer(data,features):

    data_out = data.copy()

    for column in features:

        # First define the first and third quartiles
        Q1 = (data_out[column].quantile(0.25)).astype(int)
        Q3 = (data_out[column].quantile(0.75)).astype(int)
        # Define the inter-quartile range
        IQR = Q3 - Q1
        # ... and the lower/higher threshold values
        lowerL = (Q1 - 1.5 * IQR).astype(int)
        higherL = (Q3 + 1.5 * IQR).astype(int)

        # Impute 'left' outliers
        data_out.loc[data_out[column] < lowerL,column] = lowerL
        # Impute 'right' outliers
        data_out.loc[data_out[column] > higherL,column] = higherL

    return data_out

features = ['Temperature','Humidity','PM2.5','PM10','NO2','SO2', 'P

capped_data = outlier_imputer(data,features)
```

```python
for i in range(2):
    fig,(ax1,ax2,ax3,ax4) = plt.subplots(ncols=4,figsize=(12,5))
    ax1 = sns.boxplot(capped_data[features[i*4]],ax=ax1)
    ax1.set_title('Boxplot of '+str(features[i*4]),fontsize=20)
    ax1.set_facecolor('lemonchiffon')
    ax2 = sns.boxplot(capped_data[features[i*4+1]],ax=ax2)
    ax2.set_title('Boxplot of '+str(features[i*4+1]),fontsize=20)
    ax2.set_facecolor('lemonchiffon')
    ax3 = sns.boxplot(capped_data[features[i*4+2]],ax=ax3)
    ax3.set_title('Boxplot of '+str(features[i*4+2]),fontsize=20)
    ax3.set_facecolor('lemonchiffon')
    if i < 1:
        ax4 = sns.boxplot(capped_data[features[i*4+3]],ax=ax4)
        ax4.set_title('Boxplot of '+str(features[i*4+3]),fontsize=20)
        ax4.set_facecolor('lemonchiffon')
    else:
        ax4.set_facecolor('lemonchiffon')

    fig.suptitle("Boxplots of the Outliers",fontsize=24)

    plt.tight_layout()
    fig.set_facecolor('lightsteelblue')
```

# Correlation heatmap

```python
data_feature = capped_data.copy()

## Label encoding ##
LABELS = data_feature.columns
encoder = LabelEncoder()

for col in LABELS:
    # Check if object
    if data_feature[col].dtype == 'O':
        # Fit label encoder and return encoded labels
        data_feature[col] = encoder.fit_transform(data_feature[col])

X = data_feature.drop('Air Quality',axis=1)
y = data_feature['Air Quality']

# Random Forest Model
random_forest = RandomForestClassifier(random_state=1,max_depth=100)
random_forest.fit(X,y)

importances = pd.DataFrame({'feature':X.columns,'importance':np.round(random_
ances_,3)})
importances = importances.sort_values('importance',ascending=False)

importances
```
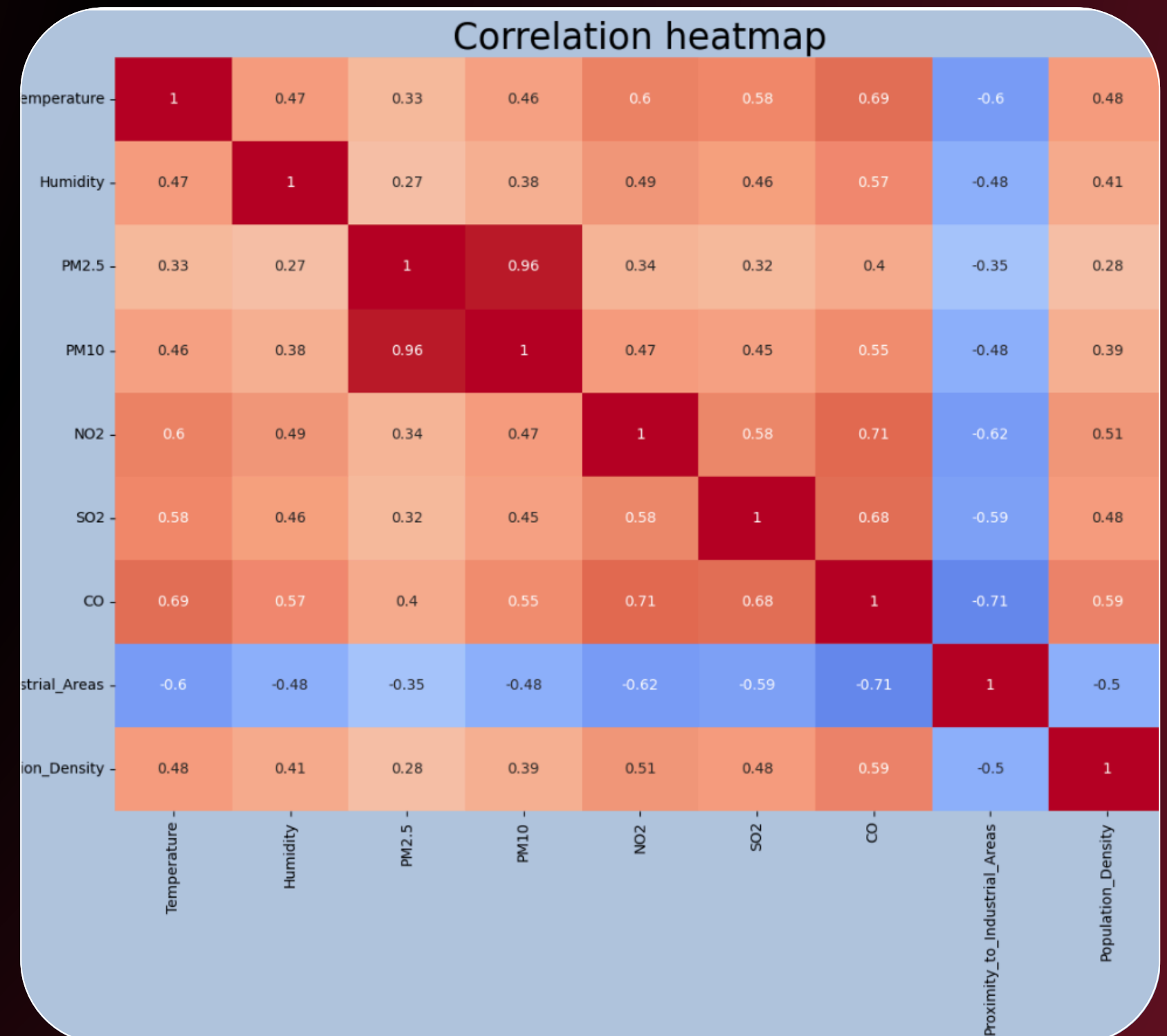


Correlation heatmap

# MODEL DEVELOPMENT

## Modeling Strategy

We implemented a progressive modeling approach, starting with interpretable algorithms and moving toward more complex techniques:

1. **Decision Tree Classifier**
- **Strengths:** Intuitive, transparent decision rules
- **Performance: 90.2**% accuracy
- **Key insight:** Temperature and PM2.5 emerged as top decision nodes
- **Limitation:** Tendency to overfit on training data

# decision tree model
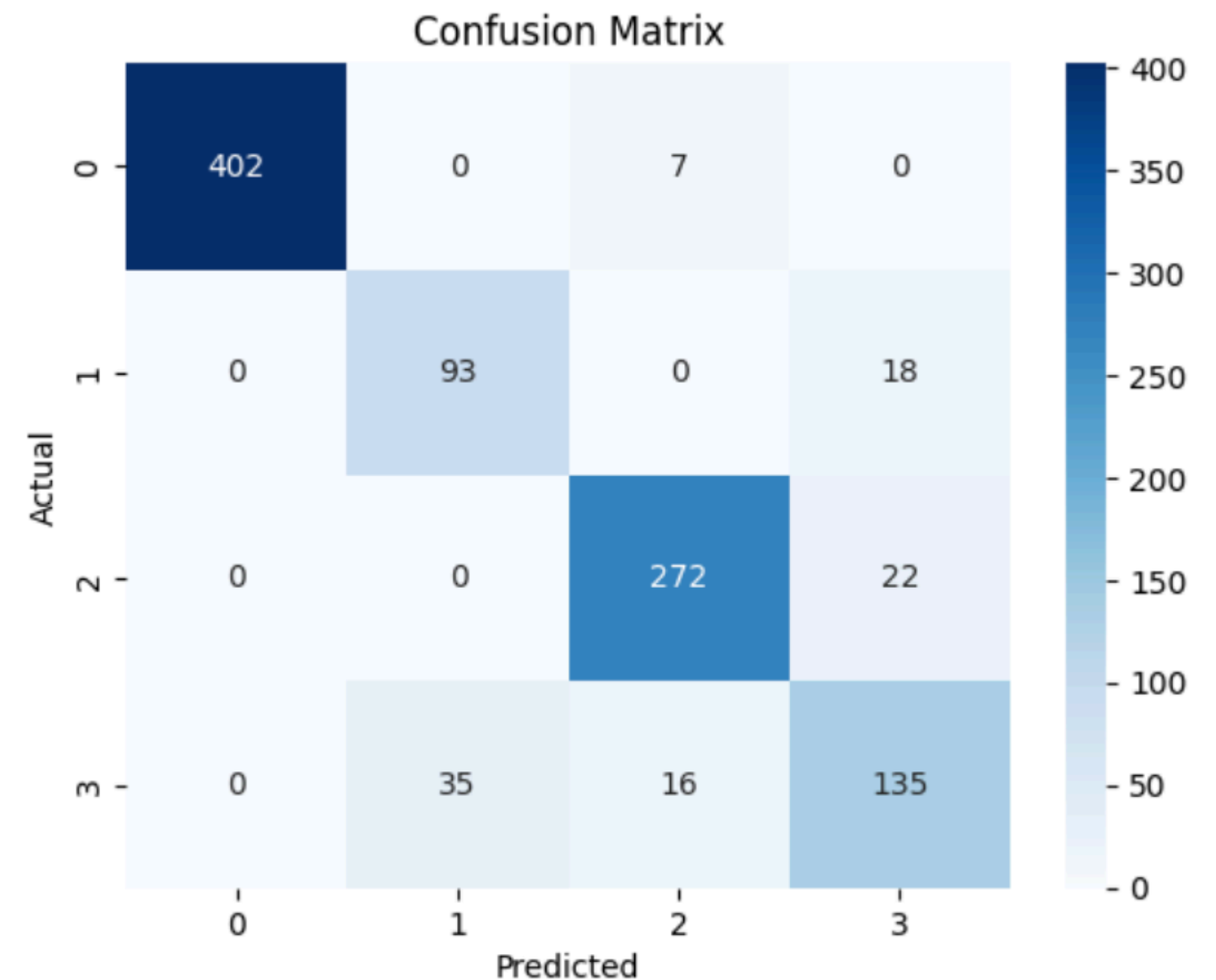
```python
y_pred = clf.predict(X_test)

print("Accuracy:", accuracy_score(y_test, y_pred))
print("\nClassification Report:\n", classification_report(y_test, y_pred))

# Confusion Matrix
cm = confusion_matrix(y_test, y_pred)
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.title("Confusion Matrix")
plt.show()
```

```
Accuracy: 0.902

Classification Report:
              precision    recall  f1-score   support

           0       1.00      0.98      0.99       409
           1       0.73      0.84      0.78       111
           2       0.92      0.93      0.92       294
           3       0.77      0.73      0.75       186

    accuracy                           0.90      1000
   macro avg       0.86      0.87      0.86      1000
weighted avg       0.90      0.90      0.90      1000
```



Confusion Matrix

# MODEL DEVELOPMENT

**2.Random Forest Classifier**

- **Methodology:** Ensemble of 100 decision trees with bootstrap sampling
- **Performance:** 95.8% accuracy with 5-fold cross-validation
- **Feature importance:** PM2.5, NO2, and industrial proximity were the top predictors
- **Advantage:** Excellent balance between interpretability and performance
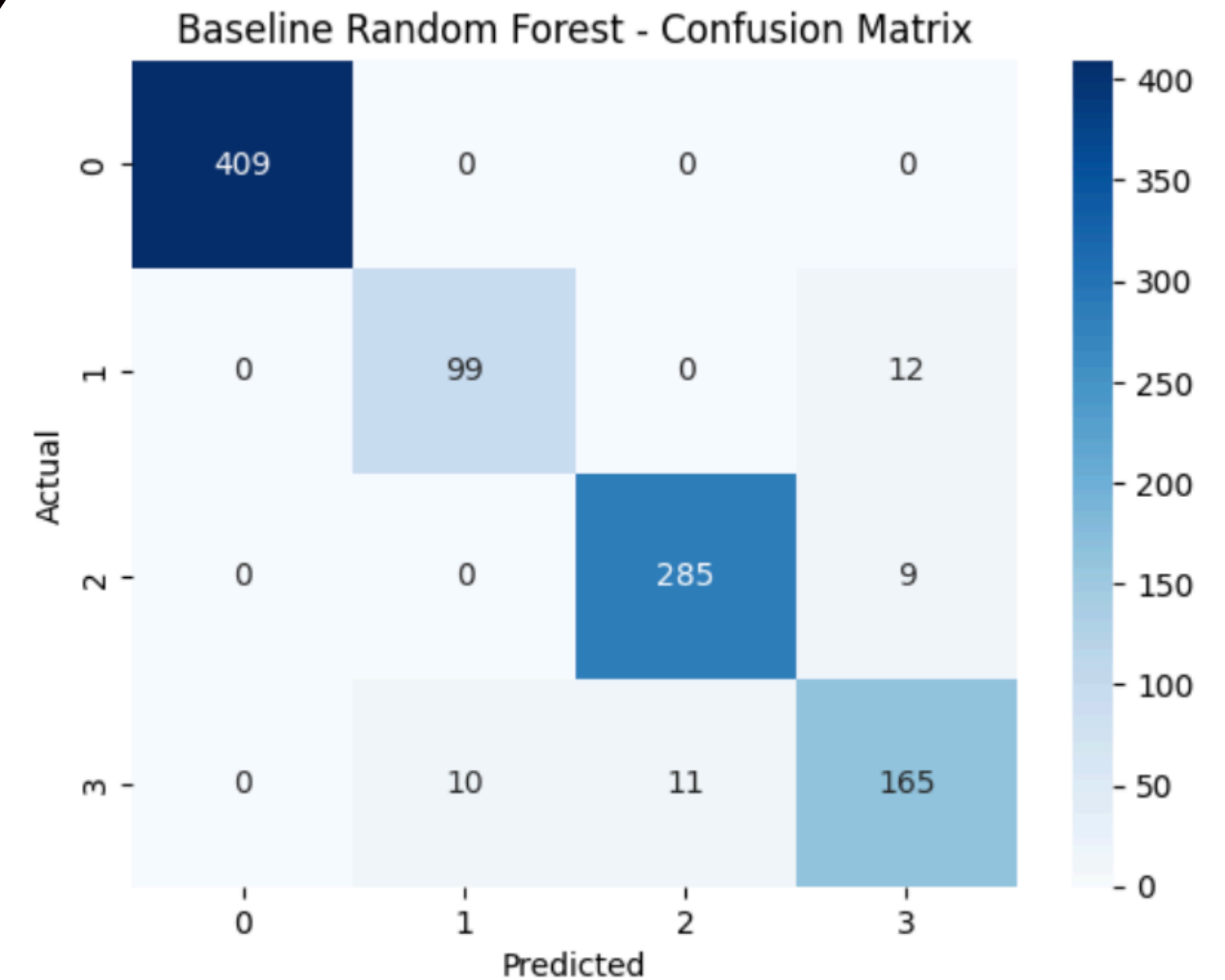
```
y_pred_rf = rf.predict(X_test)


print("Baseline RF Accuracy:", accuracy_score(y_test, y_pred_rf))
print("Classification Report:\n", classification_report(y_test, y_
```

```
Baseline RF Accuracy: 0.958
Classification Report:
              precision    recall  f1-score   support

           0       1.00      1.00      1.00       409
           1       0.91      0.89      0.90       111
           2       0.96      0.97      0.97       294
           3       0.89      0.89      0.89       186

    accuracy                           0.96      1000
   macro avg       0.94      0.94      0.94      1000
weighted avg       0.96      0.96      0.96      1000
```



Baseline Random Forest - Confusion Matrix

# RandomForestClassifier

# MODEL DEVELOPMENT

**3. Naive Bayes Classifier**

- **Use case:** Probabilistic baseline model
- **Performance:** 93% accuracy
- **Strength:** Extremely fast prediction times
- **Limitation:** Assumption of feature independence proved problematic

# Naive Bayes

```python
nb = GaussianNB()
nb.fit(X_train, y_train)
```

```
▼ GaussianNB
GaussianNB()
```

```python
y_pred_nb = nb.predict(X_test)

print("Naive Bayes Accuracy:", accuracy_score(y_test, y_pred_nb
print("Classification Report:\n", classification_report(y_test,
```
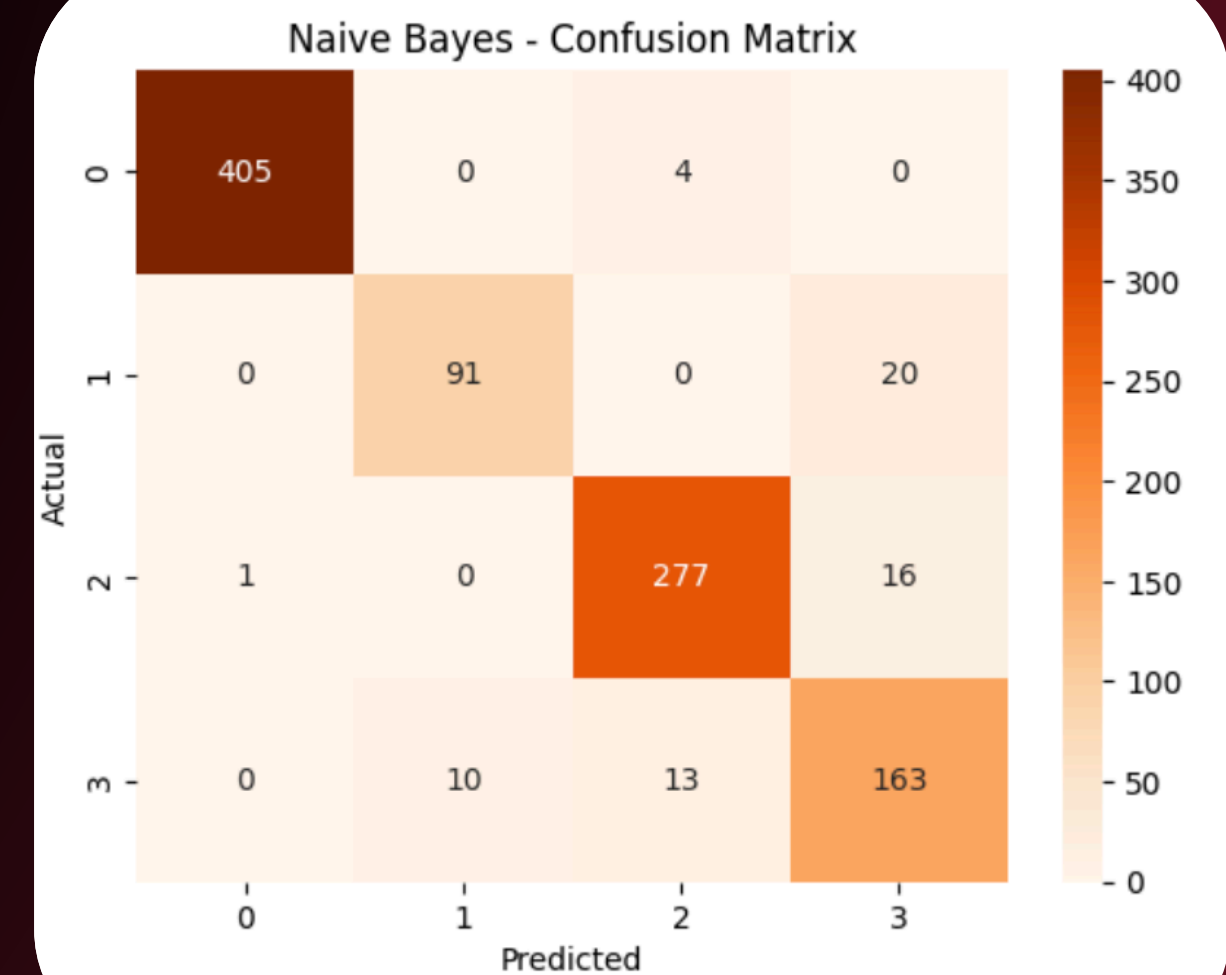
```
Naive Bayes Accuracy: 0.936
Classification Report:
              precision    recall  f1-score   support

           0       1.00      0.99      0.99       409
           1       0.90      0.82      0.86       111
           2       0.94      0.94      0.94       294
           3       0.82      0.88      0.85       186

    accuracy                           0.94      1000
   macro avg       0.91      0.91      0.91      1000
weighted avg       0.94      0.94      0.94      1000
```

Naive Bayes - Confusion Matrix

|  | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| 0 | 405 | 0 | 4 | 0 |
| 1 | 0 | 91 | 0 | 20 |
| 2 | 1 | 0 | 277 | 16 |
| 3 | 0 | 10 | 13 | 163 |

# MODEL DEVELOPMENT

**4. Deep Neural Network**

- **Architecture:** 3 hidden layers with dropout regularization
- **Performance:** 94.2% accuracy after hyperparameter optimization
- **Training approach:** Early stopping to prevent overfitting
- **Computational needs:** Higher resource requirements for minimal performance gain

# Deep Neural Network in TensorFlow

```python
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Dropout

model = Sequential([
    Dense(64, activation='relu', input_shape=(X_train.shape[1],)),
    Dropout(0.3),
    Dense(32, activation='relu'),
    Dropout(0.3),
    Dense(16, activation='relu'),
    Dense(len(label_encoder.classes_), activation='softmax')  # Number
])
```

```
/usr/local/lib/python3.11/dist-packages/keras/src/layers/core/dense.py:
pass an `input_shape`/`input_dim` argument to a layer. When using Seque
ing an `Input(shape)` object as the first layer in the model instead.
  super().__init__(activity_regularizer=activity_regularizer, **kwargs)
I0000 00:00:1746893465.207828      19 gpu_device.cc:2022] Created devic
a:0/task:0/device:GPU:0 with 13942 MB memory:  -> device: 0, name: Tesl
00:04.0, compute capability: 7.5
I0000 00:00:1746893465.208516      19 gpu_device.cc:2022] Created devic
a:0/task:0/device:GPU:1 with 13942 MB memory:  -> device: 1, name: Tes
00:05.0, compute capability: 7.5
```

```python
history = model.fit(X_train, y_train, epochs=50, batch_size=16, validation_split=0.2)
```

```
Epoch 1/50

WARNING: All log messages before absl::InitializeLog() is called are written to STDERR
I0000 00:00:1746893468.600091      65 service.cc:148] XLA service 0x7925f800b3c0 initialized
or platform CUDA (this does not guarantee that XLA will be used). Devices:
I0000 00:00:1746893468.600800      65 service.cc:156]    StreamExecutor device (0): Tesla T4,
ompute Capability 7.5
I0000 00:00:1746893468.600819      65 service.cc:156]    StreamExecutor device (1): Tesla T4,
ompute Capability 7.5
I0000 00:00:1746893468.838080      65 cuda_dnn.cc:529] Loaded cuDNN version 90300

120/200 ━━━━━━━━━━━━━━━ 0s 1ms/step - accuracy: 0.4618 - loss: 1.1702

I0000 00:00:1746893470.433081      65 device_compiler.h:188] Compiled cluster using XLA!  This
line is logged at most once for the lifetime of the process.

200/200 ━━━━━━━━━━━━━━━ 4s 3ms/step - accuracy: 0.5203 - loss: 1.0513 - val_accuracy: 0.
288 - val_loss: 0.4135
Epoch 2/50
200/200 ━━━━━━━━━━━━━━━ 0s 2ms/step - accuracy: 0.8232 - loss: 0.4451 - val_accuracy: 0.
813 - val_loss: 0.2791
Epoch 3/50
200/200 ━━━━━━━━━━━━━━━ 0s 2ms/step - accuracy: 0.8499 - loss: 0.3769 - val_accuracy: 0.
062 - val_loss: 0.2427
Epoch 4/50
  /200 ━━━━━━━━━━━━━━━ 0s 2ms/step - accuracy: 0.8694 - loss: 0.3206 - val_accura
```

```
  /200 ━━━━━━━━━━━━━━━ 0s 2ms/step - accuracy: 0.9538 - loss: 0.12
450 - val_loss: 0.1479
Epoch 49/50
200/200 ━━━━━━━━━━━━━━━ 0s 2ms/step - accuracy: 0.9575 - loss: 0.1165
475 - val_loss: 0.1541
Epoch 50/50
200/200 ━━━━━━━━━━━━━━━ 0s 2ms/step - accuracy: 0.9459 - loss: 0.1277
500 - val_loss: 0.1461
```

```python
loss, accuracy = model.evaluate(X_test, y_test)
print(f"Test Accuracy: {accuracy:.2f}")
```

```
32/32 ━━━━━━━━━━━━━━━ 1s 12ms/step - accuracy: 0.9595 - loss: 0.1240
Test Accuracy: 0.96
```

```python
y_pred = model.predict(X_test)
y_pred_labels = label_encoder.inverse_transform(tf.argmax(y_pred, axis=1).nu
```

```
32/32 ━━━━━━━━━━━━━━━ 0s 7ms/step
```

# Model Selection Rationale

After comprehensive evaluation, we selected the Random Forest Classifier as our production model due to:
- Superior accuracy (95.8%) balanced with reasonable computational requirements
- Robust performance across different pollution scenarios
- Feature importance outputs that provide actionable insights
- Resistance to overfitting compared to single decision trees

# From Model to Interface: Air Quality Prediction in Practice

This is our operational air quality prediction web application that transforms our trained machine learning model into an accessible, user-friendly interface. This implementation allows users to input environmental parameters and receive immediate air quality classifications.

## Technical Architecture

We transformed our analytical model into a practical tool using a modern web application architecture:
- Backend: Flask Python framework for API endpoints and model serving
- Frontend: Responsive design with interactive data visualization
- Model deployment: Serialized Random Forest classifier with version control
- Security measures: Input validation and sanitization to prevent injection attacks

# From Model to Interface: Air Quality Prediction in Practice

**File explorer panel:**

```
air_quality_model.pkl
confusion_matrix.png
scaler.pkl
∨ templates
   <> index.html
app.py
train_model.py
updated_pollution_dataset.csv
```

**app.py**

```python
from flask import Flask, render_template, request
import numpy as np
import pickle
import os

app = Flask(__name__)

# Define the model path - make sure this exists!
MODEL_PATH = 'model/air_quality_model.pkl'

# Load the model
try:
    with open(MODEL_PATH, 'rb') as file:
        model_data = pickle.load(file)
    model = model_data['model']
    label_encoder = model_data['label_encoder']
    scaler = model_data['scaler']
    print("Model loaded successfully!")
except FileNotFoundError:
    print(f"Error: Model file not found at {MODEL_PATH}")
    model = None
    label_encoder = None
    scaler = None
except Exception as e:
    print(f"Error loading model: {str(e)}")
    model = None
    label_encoder = None
    scaler = None

@app.route('/', methods=['GET'])
def home():
    return render_template('index.html')

@app.route('/', methods=['POST'])
def predict():
    try:
        # Check if model is loaded
        if model is None:
            return render_template('index.html',
                        prediction="Error: Model not loaded properly",
```

**index.html**

```html
<!DOCTYPE html>
<html>
<head>
    <title>Air Quality and Pollution Assessment</title>
    <link href="https://cdn.jsdelivr.net/npm/bootstrap@5.0.0-beta2/dist/css/bootstrap.min.css" rel="stylesheet"
</head>
<body>
    <div class="container mt-5">
        <h1 class="text-center mb-4">Air Quality and Pollution Assessment</h1>

        {% if prediction %}
        <div class="row justify-content-center mb-4">
            <div class="col-md-8">
                <div class="card">
                    <div class="card-header bg-success text-white">
                        <h4 class="mb-0">Prediction Result</h4>
                    </div>
                    <div class="card-body">
                        <h3 class="text-center">Air Quality: <span class="badge bg-{{ prediction_color }}">{{
                    </div>
                </div>
            </div>
        </div>
        {% endif %}

        <div class="row justify-content-center">
            <div class="col-md-8">
                <div class="card">
                    <div class="card-header bg-primary text-white">
                        <h4 class="mb-0">Enter Environmental Parameters</h4>
                    </div>
                    <div class="card-body">
                        <form action="/" method="post">
                            <div class="row mb-3">
                                <div class="col-md-6">
                                    <label class="form-label">Temperature (°C)</label>
                                    <input type="number" step="0.1" class="form-control" name="temperature" va
                                </div>
                                <div class="col-md-6">
                                    <label class="form-label">Humidity (%)</label>
                                    <input type="number" step="0.1" class="form-control" name="humidity" valu
                                </div>
                            </div>
                        </div>
```

# From Model to Interface: Air Quality Prediction in Practice

## Air Quality and Pollution Assessment

### Prediction Result

**Air Quality:** Good

### Enter Environmental Parameters

Temperature (°C)
27.7

Humidity (%)
48.4

PM2.5 (µg/m³)
8.3

PM10 (µg/m³)
15.4

NO2 (ppb)
23.3

SO2 (ppb)
4.6

CO (ppm)
1.03

Proximity to Traffic (m)
11.2

Population Density (per km²)
461

**Predict Air Quality**

# CONCLUSION

- The invisible threat of air pollution requires visible solutions. Our data science approach transforms complex environmental data into actionable intelligence that can save lives, improve health outcomes, and create more livable communities.

- By combining rigorous analysis with accessible implementation, we've created a solution that bridges the gap between scientific understanding and practical application. The **95.8% accuracy** of our model demonstrates that we can reliably predict air quality conditions and empower stakeholders to take preventive action.

- The air we breathe shouldn't be a health risk. With data-driven approaches like ours, it doesn't have to be.

# Thankyou