

# Global Pollution Analysis and Energy Recovery

Machine Learning Based Classification Approach

Abhisha Whaval





# Problem Statement

The accelerating pace of industrialization and energy consumption has created unprecedented environmental challenges worldwide. Rising pollution levels threaten ecosystems, public health, and sustainable development across nations.



## Global Pollution Crisis

Rapid industrialization driving increased emissions across developed and developing nations



## Classification Challenge

Difficulty in accurately categorizing pollution severity levels across diverse countries



## Need for Data-Driven Solutions

Requirement for machine learning approaches to derive actionable environmental insights

# Project Objectives

This research leverages machine learning techniques to classify global pollution patterns and provide evidence-based recommendations for environmental policy. Our approach combines data preprocessing, model comparison, and insight generation.

01

---

## Data Preprocessing

Clean and transform raw global pollution data for analysis

03

---

## Model Comparison

Evaluate multiple ML classifiers to identify optimal approach

02

---

## Severity Classification

Categorize countries into Low, Medium, and High pollution levels

04

---

## Generate Insights

Derive actionable recommendations for environmental policy

# Dataset Overview

## Data Source

## Global\_Pollution\_Analysis.csv

Comprehensive environmental dataset containing multi-year pollution metrics across countries worldwide.

## Target Variable

## Pollution Severity

Low / Medium / High

## Key Attributes

- **Geographic & Temporal:** Country, Year
- **Emissions Data:** CO<sub>2</sub> Emissions (in Megatons)
- **Pollution Indicators:** Air Quality Index, Water Contamination Levels, Soil Degradation Metrics
- **Energy Factors:** Total Energy Consumption, Renewable vs. Non-renewable Sources
- **Industrial Metrics:** Manufacturing Output, Industrial Growth Rate



# Phase 1: Data Preprocessing

Robust data preprocessing ensures model accuracy and reliability. We implemented a comprehensive pipeline to handle missing values, encode categorical variables, normalize features, and engineer the target variable.



## Missing Value Imputation

Applied mean imputation for numerical features and mode imputation for categorical attributes



## Label Encoding

Transformed categorical features into numerical representations for model compatibility



## Feature Scaling

StandardScaler normalization to ensure features contribute equally to model training



## Target Creation

Quantile-based discretization of CO<sub>2</sub> emissions to define pollution severity classes





## Phase 2: Machine Learning Models

We evaluated three distinct classification algorithms, each offering unique approaches to pattern recognition and decision-making.

1

### Naive Bayes

**Approach:** Probabilistic classifier based on Bayes' theorem

**Strength:** Fast training, effective with independent features

**Application:** Baseline model for probabilistic classification

2

### K-Nearest Neighbors

**Approach:** Instance-based learning using distance metrics

**Strength:** Non-parametric, captures local patterns

**Application:** Distance-based similarity classification

3

### Decision Tree

**Approach:** Hierarchical rule-based decision structure

**Strength:** Interpretable, handles non-linear relationships

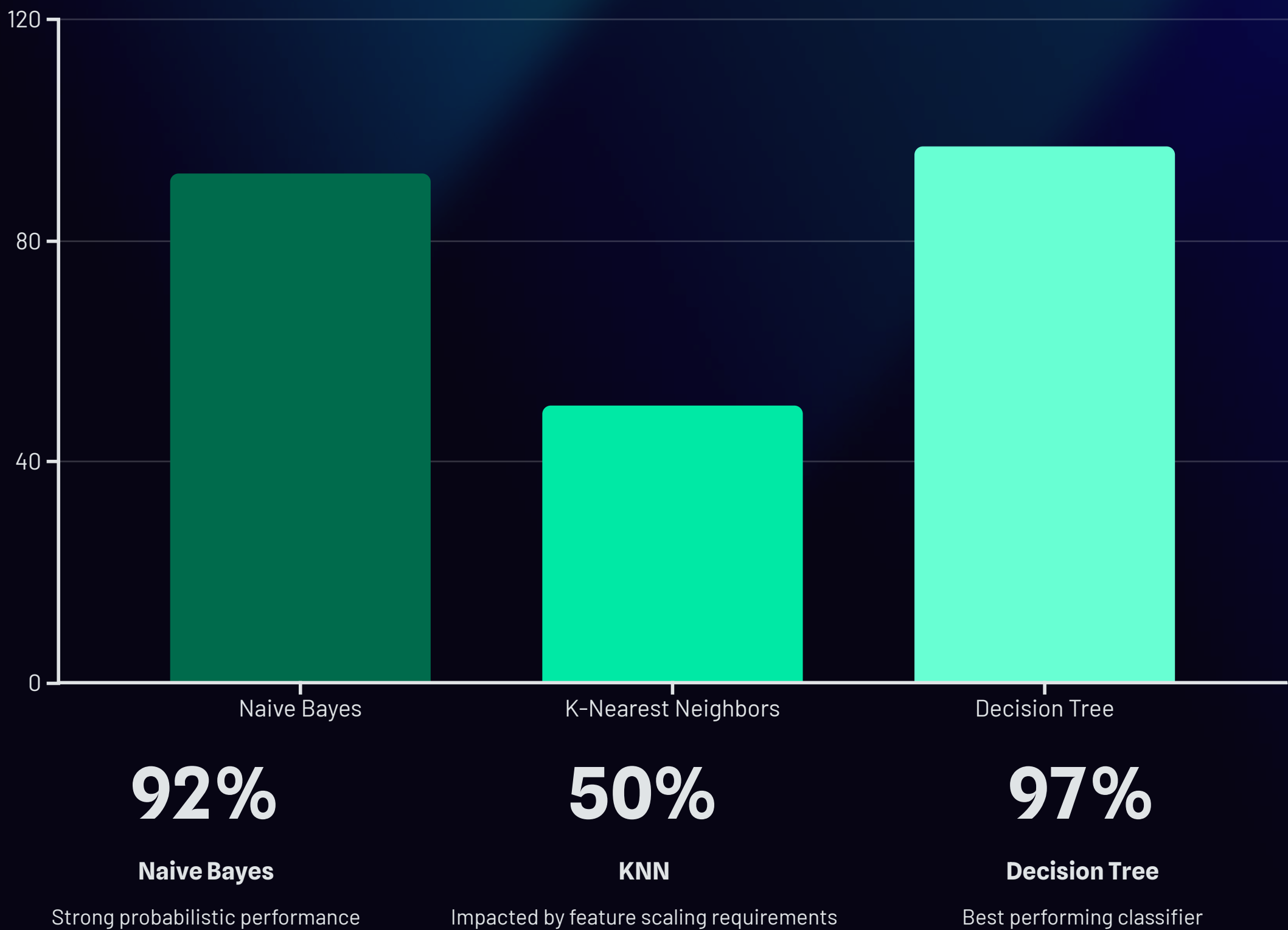
**Application:** Feature importance and transparent decisions

## Evaluation Framework

All models assessed using: **Accuracy**, **Confusion Matrix**, **Precision**, **Recall**, and **F1-Score**

# Model Performance Results

Comparative analysis reveals significant performance differences across classification approaches. The Decision Tree model demonstrated superior accuracy in predicting pollution severity levels.



# Confusion Matrix Analysis

Confusion matrices provide detailed insight into model prediction patterns, revealing where each classifier succeeds and struggles with pollution severity categorization.

## Naïve Bayes



**92% Accuracy**

Strong overall performance with occasional Medium-High confusion

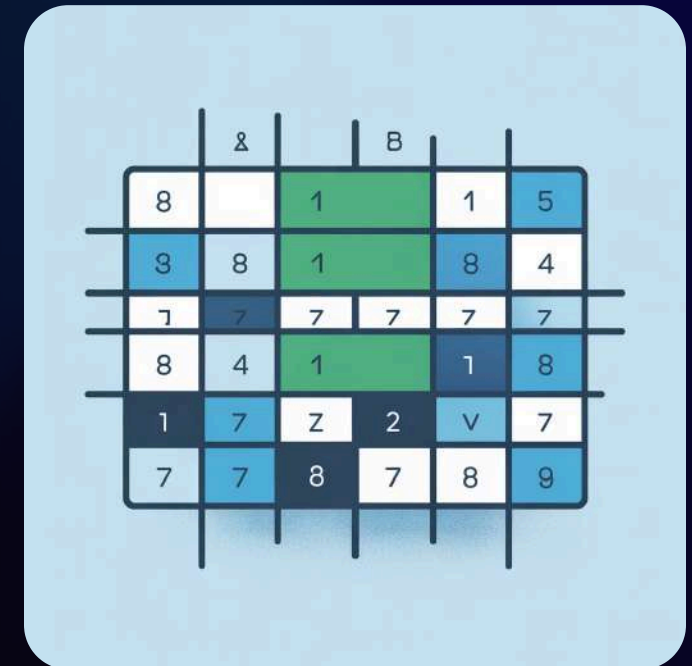
## K-Nearest Neighbors



**50% Accuracy**

Significant misclassification across all severity levels

## Decision Tree



**97% Accuracy**

Minimal errors, excellent class separation

**Key Observation:** Decision Tree demonstrates superior ability to correctly classify pollution severity with minimal false positives and false negatives across all categories.



# Comprehensive Model Comparison

Model	Accuracy	Strengths	Limitations
Naive Bayes	92%	Fast, probabilistic reasoning	Assumes feature independence
K-Nearest Neighbors	50%	Simple, non-parametric	Sensitive to scaling, computationally expensive
Decision Tree	97%	Highest accuracy, interpretable rules	Potential overfitting without pruning

## Critical Finding

The Decision Tree classifier outperforms alternatives by **5% over Naive Bayes** and **47% over KNN**, demonstrating its effectiveness in capturing complex pollution patterns and non-linear relationships between environmental variables.

# Key Research Findings

Our analysis uncovered critical relationships between pollution severity and environmental factors, providing evidence-based insights for policy development.

CO <sub>2</sub> Emissions Driver	Energy Impact	Scaling Importance	Model Interpretability
Strong positive correlation between high CO <sub>2</sub> emissions and severe pollution classification	Total energy consumption emerged as a primary predictor of pollution severity levels	Feature standardization significantly improved KNN model performance	Decision Trees provide transparent, explainable classification rules for stakeholders

# Actionable Insights & Recommendations

Based on our machine learning analysis, we propose evidence-based strategies for reducing global pollution severity and promoting sustainable development.



## Clean Energy Transition

Accelerate adoption of renewable energy sources (solar, wind, hydroelectric) to reduce carbon emissions

*Target: 50% renewable energy by 2035*



## Emission Controls

Implement strict regulatory frameworks with mandatory emission limits for industrial sectors

*Focus: Manufacturing, transportation, energy production*



## Industrial Efficiency

Promote energy-efficient technologies and sustainable manufacturing practices across industries

*Incentivize green innovation through tax benefits*



## Data-Driven Monitoring

Deploy ML-based systems for real-time pollution tracking and early warning systems

Enable proactive intervention and policy adjustment



# Conclusion



## Research Impact

This project demonstrates the powerful application of machine learning in environmental science, successfully classifying global pollution severity with **97% accuracy** using Decision Tree algorithms.

## Key Achievements

- Developed robust preprocessing pipeline for environmental data
- Compared three distinct ML classification approaches
- Identified Decision Tree as optimal classifier for pollution analysis
- Generated actionable insights supporting environmental policy

## Future Implications

This methodology supports **data-driven environmental policymaking**, enabling governments and organizations to make informed decisions for sustainable development and pollution mitigation.



# Thank You

## Questions & Discussion

Abhisha | [abhishawhaval@gmail.com](mailto:abhishawhaval@gmail.com)

Global Pollution Analysis and Energy Recovery using Machine Learning

