

project

May 31, 2025

1 Grand Challenge: Understanding Water Pollution and Global Health Display

1.1 Introduction

Water pollution is a pressing global issue that affects not just ecosystems, but also human health, economic stability, and social equity. Contaminated water sources contribute to the spread of preventable diseases, especially in developing regions, where access to clean water and sanitation remains limited. Despite global efforts, many countries continue to face high rates of waterborne diseases and environmental degradation, while others have managed to control these risks effectively.

This project aims to explore **why some countries struggle with water pollution impacts more than others** by analyzing a comprehensive dataset of water quality indicators, health outcomes, and socioeconomic infrastructure. Using statistical modeling and visual analytics, we seek to uncover the **key relationships between pollution, disease rates, and national resilience factors** like GDP, healthcare access, and sanitation coverage.

1.2 Project Goals

- Quantify the level of water contamination across countries based on WHO/EPA thresholds.
- Correlate contamination levels with health indicators like diarrheal diseases, cholera, typhoid, and infant mortality.
- Identify whether higher GDP, better healthcare, and stronger infrastructure offer significant protection.
- Compare “high-risk” and “low-risk” countries to extract actionable insights.
- Suggest policy directions and water treatment strategies that can mitigate future risks.

By combining data science with environmental analysis, we hope to not only diagnose the problem — but also highlight what works, where, and why.

1.3 Imported Libraries – Why We Use Them

To build a complete machine learning pipeline for environmental risk assessment, we relied on the following core Python libraries:

- **pandas**: For loading, cleaning, grouping, and aggregating the dataset across countries and features.
- **numpy**: For numerical operations and efficient manipulation of arrays and vectors during calculations.

- **matplotlib.pyplot**: To create static plots such as bar charts, line graphs, and radar charts for data exploration and presentation.

1.3.1 From Scikit-learn:

- **sklearn.preprocessing**:
 - **MinMaxScaler**: To scale raw risk scores into a 0–100 range.
 - **StandardScaler**: To standardize features for clustering and modeling.
- **sklearn.cluster**:
 - **KMeans**: To identify peer country clusters based on pollution and infrastructure profiles.
- **sklearn.ensemble**:
 - **RandomForestRegressor**: To model and explain what factors most influence a country's pollution risk.
- **sklearn.model_selection**:
 - **train_test_split**: To split the dataset into training and testing sets during model evaluation.
- **scipy.spatial.distance**:
 - **cdist**: To compute distances between countries for nearest peer selection.
- **warnings**: To suppress unnecessary output and maintain clean logs during runtime.

These libraries worked seamlessly together to support everything from risk scoring and peer identification to explainable ML and actionable policy recommendations.

```
[16]: # Importing the libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

from sklearn.preprocessing import MinMaxScaler
from sklearn.preprocessing import StandardScaler
from sklearn.cluster import KMeans
from scipy.spatial.distance import cdist
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestRegressor

import warnings
warnings.filterwarnings('ignore')

# Loading the dataset
file_path = 'water_pollution_disease.csv'
df = pd.read_csv(file_path)

# Displaying the first few rows of the dataset
print("First few rows of the dataset:")
display(df.head())

# Basic shape information
print(f"\nDataset shape: {df.shape[0]} rows, {df.shape[1]} columns.")
```

```
# Summary statistics
print("\nSummary Statistics:")
display(df.describe())
```

First few rows of the dataset:

	Country	Region	Year	Water Source	Type	Contaminant Level (ppm)	\
0	Mexico	North	2015		Lake	6.06	
1	Brazil	West	2017		Well	5.24	
2	Indonesia	Central	2022		Pond	0.24	
3	Nigeria	East	2016		Well	7.91	
4	Mexico	South	2005		Well	0.12	

	pH Level	Turbidity (NTU)	Dissolved Oxygen (mg/L)	Nitrate Level (mg/L)	\
0	7.12	3.93	4.28	8.28	
1	7.84	4.79	3.86	15.74	
2	6.43	0.79	3.42	36.67	
3	6.71	1.96	3.12	36.92	
4	8.16	4.22	9.15	49.35	

	Lead Concentration (µg/L)	...	Cholera Cases per 100,000 people	\
0	7.89	...	33	
1	14.68	...	27	
2	9.96	...	39	
3	6.77	...	33	
4	12.51	...	31	

	Typhoid Cases per 100,000 people	\
0	44	
1	8	
2	50	
3	13	
4	68	

	Infant Mortality Rate (per 1,000 live births)	GDP per Capita (USD)	\
0	76.16	57057	
1	77.30	17220	
2	48.45	86022	
3	95.66	31166	
4	58.78	25661	

	Healthcare Access Index (0-100)	Urbanization Rate (%)	\
0	96.92	84.61	
1	84.73	73.37	
2	58.37	72.86	
3	39.07	71.07	
4	23.03	55.55	

	Sanitation Coverage (% of Population)	Rainfall (mm per year)	\
0	63.23	2800	
1	29.12	1572	
2	93.56	2074	
3	94.25	937	
4	69.23	2295	

	Temperature (°C)	Population Density (people per km ²)
0	4.94	593
1	16.93	234
2	21.73	57
3	3.79	555
4	31.44	414

[5 rows x 24 columns]

Dataset shape: 3000 rows, 24 columns.

Summary Statistics:

	Year	Contaminant Level (ppm)	pH Level	Turbidity (NTU)	\
count	3000.000000	3000.000000	3000.000000	3000.000000	
mean	2012.012667	4.954390	7.255847	2.480023	
std	7.229287	2.860072	0.720464	1.419984	
min	2000.000000	0.000000	6.000000	0.000000	
25%	2006.000000	2.560000	6.630000	1.257500	
50%	2012.000000	4.950000	7.280000	2.460000	
75%	2018.000000	7.400000	7.870000	3.660000	
max	2024.000000	10.000000	8.500000	4.990000	

	Dissolved Oxygen (mg/L)	Nitrate Level (mg/L)	\
count	3000.000000	3000.000000	
mean	6.492850	25.08025	
std	2.027966	14.50517	
min	3.000000	0.05000	
25%	4.710000	12.52500	
50%	6.490000	24.79000	
75%	8.252500	37.91000	
max	10.000000	49.99000	

	Lead Concentration (µg/L)	Bacteria Count (CFU/mL)	\
count	3000.000000	3000.000000	
mean	10.047913	2488.477333	
std	5.798238	1431.421553	
min	0.000000	0.000000	
25%	5.120000	1268.000000	
50%	10.065000	2469.000000	

75%	15.032500	3736.250000
max	20.000000	4998.000000

Access to Clean Water (% of Population) \		
count	3000.000000	
mean	64.612333	
std	20.308463	
min	30.010000	
25%	47.027500	
50%	64.780000	
75%	82.302500	
max	99.990000	

Diarrheal Cases per 100,000 people Cholera Cases per 100,000 people \		
count	3000.000000	3000.000000
mean	249.776667	24.25100
std	144.111543	14.33259
min	0.000000	0.00000
25%	124.000000	12.00000
50%	248.000000	24.00000
75%	378.000000	37.00000
max	499.000000	49.00000

Typhoid Cases per 100,000 people \	
count	3000.000000
mean	49.270000
std	28.984165
min	0.000000
25%	24.000000
50%	49.000000
75%	75.000000
max	99.000000

Infant Mortality Rate (per 1,000 live births) GDP per Capita (USD) \		
count	3000.000000	3000.000000
mean	50.811900	50036.196667
std	28.465323	28598.750508
min	2.060000	521.000000
25%	26.467500	25010.250000
50%	50.230000	49621.500000
75%	76.260000	74778.250000
max	99.990000	99948.000000

Healthcare Access Index (0-100) Urbanization Rate (%) \		
count	3000.000000	3000.000000
mean	50.029193	50.062480
std	28.896676	22.779125
min	0.190000	10.030000

25%	24.982500	30.557500
50%	50.390000	49.795000
75%	74.817500	69.727500
max	99.980000	89.980000

	Sanitation Coverage (% of Population)	Rainfall (mm per year) \
count	3000.000000	3000.000000
mean	60.371007	1591.849000
std	23.159678	817.502434
min	20.010000	200.000000
25%	40.440000	865.750000
50%	60.580000	1572.000000
75%	80.420000	2308.250000
max	99.990000	2999.000000

	Temperature (°C)	Population Density (people per km ²)
count	3000.000000	3000.000000
mean	20.130917	505.390333
std	11.689244	283.275224
min	0.060000	10.000000
25%	9.840000	254.750000
50%	20.175000	513.000000
75%	30.672500	745.000000
max	39.990000	999.000000

1.4 Dataset Overview

This dataset captures the multifaceted nature of water pollution and its impacts on public health and infrastructure across various countries. It spans 3000 data points with 24 columns, including:

- **Water quality indicators** (e.g. contaminant level, pH, turbidity, dissolved oxygen, etc.)
- **Health statistics** (e.g. diarrheal, cholera, typhoid cases, infant mortality)
- **Socioeconomic factors** (e.g. GDP per capita, healthcare access, urbanization)
- **Infrastructure access** (e.g. access to clean water, sanitation coverage, water treatment method)

The data spans from the year 2000 to 2024, providing a temporal dimension for analysis.

A preview of the first few rows helps us understand the data structure, formatting, and potential cleaning needs.

```
[18]: # Cleaning the dataset
df.isnull().sum()[df.isnull().sum() > 0]
```

```
[18]: Series([], dtype: int64)
```

1.5 Handling Missing Values

We identified missing values in the **Water Treatment Method** column. Since this is a categorical feature and dropping rows could lead to loss of important data, we fill the missing entries with

"Unknown" to preserve dataset integrity.

```
[19]: # Fill missing values with 'Unknown'
df = df.fillna('Unknown')

# Verify no missing values remain
print("\nMissing values after filling:")
print(df.isnull().sum()[df.isnull().sum() > 0])

# Save the processed dataset
df.to_csv('water_pollution_disease_processed.csv', index=False)
print("\nProcessed dataset saved as 'water_pollution_disease_processed.csv'")
```

Missing values after filling:
Series([], dtype: int64)

Processed dataset saved as 'water_pollution_disease_processed.csv'

1.6 Defining Risk Based on WHO Thresholds

To quantify pollution risk, we base our system on safety thresholds established by the World Health Organization (WHO) and EPA. Each pollutant is evaluated relative to its defined safe range.

Rather than applying a binary penalty (0 for safe, 1 for unsafe), we use a **distance-based penalty** that measures **how far** a value deviates from its safe boundary:

Pollutant Factor	Safe Range
pH Level	6.5 – 8.5
Contaminant Level (ppm)	3
Turbidity (NTU)	1
Dissolved Oxygen (mg/L)	5
Nitrate Level (mg/L)	50
Lead Concentration (µg/L)	10
Bacteria Count (CFU/mL)	= 0

1.7 Pollution Risk Scoring System

1. Deviation calculation

For each pollutant, compute the absolute difference from its safe bound (e.g. $\max(\text{val} - \text{upper}, 0)$).

2. Raw risk score

Sum the seven deviations into a single continuous value.

3. Normalize

Apply Min–Max scaling to map the raw score into a **0–100** range.

Higher scores correspond to more severe contamination.

1.8 Multi-Threshold High-Risk Frequency Analysis

To understand both **how often** and **how severely** water quality fails, we classify samples as “high-risk” at three escalating thresholds on the 0–100 scale:

- **Threshold 60**: general violation frequency
- **Threshold 70**: severe violation frequency
- **Threshold 80**: critical violation frequency

We then compute, for each country, the **percentage of samples** exceeding each threshold. In our bar-chart subplots:

- Every bar is annotated **inside** with its exact percent (“X.X%”) in white.
 - Countries that remain high across all thresholds exhibit **chronic, systemic failures**.
 - Steep drop-offs indicate **borderline or situational** risks.
 - A mid-range bump at the highest threshold flags **localized extreme hotspots**.
-

1.9 Why This Matters

- **Frequency + severity** gives a fuller picture than a single average.
- Policy makers can see whether a nation needs **system-wide overhaul**, **targeted hotspot intervention**, or **routine maintenance**.
- The annotated multi-threshold chart is both **precise** and **intuitive** for decision-grade reporting.

```
[20]: # Distance-based penalty functions
def penalty_ph(val):
    lower, upper = 6.5, 8.5
    if val < lower:
        return lower - val
    elif val > upper:
        return val - upper
    return 0
def penalty_contaminant(val): return max(val - 3, 0)
def penalty_turbidity(val): return max(val - 1, 0)
def penalty_do(val): return max(5 - val, 0)
def penalty_nitrate(val): return max(val - 50, 0)
def penalty_lead(val): return max(val - 10, 0)
def penalty_bacteria(val): return val / 1000

# Applying the penalty functions to the dataset
```



```

df['penalty_ph'] = df['pH Level'].apply(penalty_ph)
df['penalty_contaminant'] = df['Contaminant Level (ppm)'].
    ↳apply(penalty_contaminant)
df['penalty_turbidity'] = df['Turbidity (NTU)'].apply(penalty_turbidity)
df['penalty_do'] = df['Dissolved Oxygen (mg/L)'].apply(penalty_do)
df['penalty_nitrate'] = df['Nitrate Level (mg/L)'].apply(penalty_nitrate)
df['penalty_lead'] = df['Lead Concentration (µg/L)'].
    ↳apply(penalty_lead)
df['penalty_bacteria'] = df['Bacteria Count (CFU/mL)'].
    ↳apply(penalty_bacteria)

# Summing the penalty columns to get the raw risk score
penalty_cols = ['penalty_ph', 'penalty_contaminant', 'penalty_turbidity',
                'penalty_do', 'penalty_nitrate', 'penalty_lead',
                ↳'penalty_bacteria']
df['raw_risk_score'] = df[penalty_cols].sum(axis=1)

# Scale the raw risk score to a 0-100 range
scaler = MinMaxScaler(feature_range=(0, 100))
df['risk_score'] = scaler.fit_transform(df[['raw_risk_score']])

# Define thresholds and compute high-risk percentages
thresholds = [60, 70, 80]
high_risk_pcts = {}
for thr in thresholds:
    col = f'is_high_risk_{thr}'
    df[col] = df['risk_score'] >= thr
    high_risk_pcts[thr] = (df.groupby('Country')[col].mean() * 100).
        ↳sort_values(ascending=False)

# Plot with annotations
fig, axs = plt.subplots(1, 3, figsize=(18, 6), sharey=True)

for ax, thr in zip(axs, thresholds):
    pct = high_risk_pcts[thr]
    bars = ax.barh(pct.index, pct.values, color='steelblue')
    ax.invert_yaxis()
    ax.set_title(f'High-Risk % (Threshold {thr})')
    ax.set_xlabel('% of Samples')
    ax.grid(True, axis='x', linestyle='--', alpha=0.5)

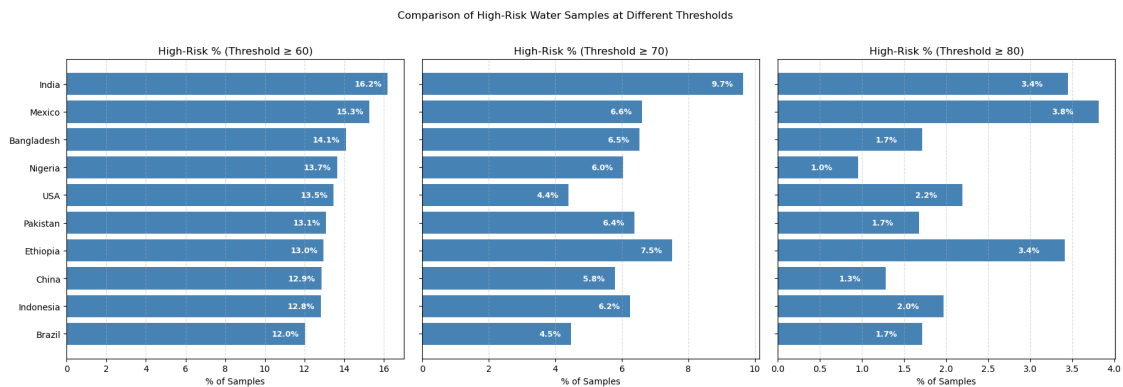
# Annotate each bar inside in white
for bar in bars:
    w = bar.get_width()
    y = bar.get_y() + bar.get_height() / 2
    ax.text(
        w - 0.3,      # shift left by 0.3 units

```

```

        y,
        f'{w:.1f}%', # format label
        va='center',
        ha='right',
        color='white',
        fontsize=9,
        fontweight='bold'
    )
plt.suptitle('Comparison of High-Risk Water Samples at Different Thresholds',
            ↪y=1.02)
plt.tight_layout()
plt.show()

```



1.10 Interpreting Country Water-Quality Performance

Using our multi-threshold analysis of high-risk samples, we classify each country along two dimensions:

1. **Frequency of general violations** (percentage of samples with `risk_score ≥ 60`)
2. **Persistence of critical failures** (percentage of samples with `risk_score ≥ 80`)

Country	60%	80%	Classification	Meaning
	Unsafe	Critical		
India	16.2 %	3.4 %	Doing Bad	Chronic, system-wide contamination requiring major infrastructure overhaul.
Mexico	15.3 %	3.8 %	Doing Bad	Widespread and severe pollution—urgent, large-scale intervention needed.
Ethiopia	13.0 %	3.4 %	Hotspot Risk	Generally acceptable water, but isolated regions suffer catastrophic failures.
Bangladesh	14.1 %	1.7 %	Borderline Risk	Frequent mild violations; very few extreme cases—focus on preventative maintenance.
Nigeria	13.7 %	1.0 %	Borderline Risk	Many samples just over safe limits; critical failures are rare.

Country	60% Unsafe	80% Critical	Classification	Meaning
Pakistan	13.1 %	1.7 %	Borderline Risk	Steady minor issues; occasional serious spikes warrant targeted fixes.
Indonesia	12.8 %	2.0 %	Borderline Risk	Moderate frequency of unsafe samples with some critical outliers.
USA	13.5 %	2.2 %	Borderline Risk	Common mild issues plus sporadic extreme spikes—enhance monitoring.
China	12.9 %	1.3 %	Doing Fair	Most water is within or near safe limits; maintain current protocols.
Brazil	12.0 %	1.7 %	Doing Fair	Infrequent violations and few critical cases—continue existing practices.

1.10.1 Key Takeaways

- **“Doing Bad” (India, Mexico):** High rates of both frequent and critical failures. Require systemic infrastructure investment and regulatory overhaul.
- **“Hotspot Risk” (Ethiopia):** Water quality is generally good, but specific areas exhibit extreme contamination. Targeted remediation in hotspots is essential.
- **“Borderline Risk” (Bangladesh, Nigeria, Pakistan, Indonesia, USA):** Widespread minor violations with few critical incidents. Focus on preventative maintenance, routine monitoring, and incremental improvements.
- **“Doing Fair” (China, Brazil):** Relatively safe water systems with low violation rates. Emphasis should be on sustaining and fine-tuning current water-quality management strategies.

This classification framework combines frequency and severity metrics to prioritize tailored policy and engineering responses for each country.

1.11 Radar-Chart Pairwise Comparisons

To dissect *which specific pollutants* drive the gap between our worst- and best-performing countries, we draw four radar-chart overlays. Each subplot shows:

1. A **green dotted heptagon** marking the normalized WHO/EPA “safe” threshold for each contaminant.
2. A **blue/orange filled web** for two countries, scaled 0–1 relative to the global maximum of each pollutant.

The four comparisons are:

1. **India vs. China**
 - India (blue) is the worst performer; China (orange) is the best.
 - Where India’s web extends beyond the green outline, it exceeds safe limits most severely; where China stays inside, it maintains compliance.
2. **India vs. Brazil**

- Contrasts India against Brazil, our second-best country.
 - Highlights whether Brazil’s contamination profile tracks China’s or shows different strengths (e.g. lower turbidity vs. lead).
3. **Mexico vs. China**
- Places Mexico (second-worst) alongside China.
 - Pinpoints any pollutant gaps Mexico shares with India, and where China’s best practices offer a model.
4. **Mexico vs. Brazil**
- Compares Mexico to Brazil.
 - Reveals whether Brazil’s mitigation strategies differ enough from China’s to provide alternative benchmarks.

1.11.1 How to Read the Charts

- **Radial axis (0–1):** Each spoke is the country’s *mean pollutant level* divided by that pollutant’s global maximum.
- **Green dotted heptagon:** The normalized safe-threshold for each pollutant, so any colored area outside this line indicates a breach of WHO/EPA guidelines.
- **Blue/Orange webs:** Country profiles—larger area = worse relative contamination.

By overlaying each “bad” country against a “fair” counterpart, and layering in the safe boundary, these radar charts make it easy to see **which pollutants** require the most urgent attention and **which best practices** can be adapted from peer nations.

```
[22]: # Defining the pairs of countries and features to compare
pairs = [
    ('India', 'China'),
    ('India', 'Brazil'),
    ('Mexico', 'China'),
    ('Mexico', 'Brazil')
]
features = [
    'pH Level', 'Contaminant Level (ppm)', 'Turbidity (NTU)',
    'Dissolved Oxygen (mg/L)', 'Nitrate Level (mg/L)',
    'Lead Concentration (µg/L)', 'Bacteria Count (CFU/mL)'
]
short_labels = ['pH', 'Contam.', 'Turb.', 'DO', 'NO ', 'Pb', 'Bact.']
safe_vals = {
    'pH Level': (6.5+8.5)/2, # mid-point of safe pH
    'Contaminant Level (ppm)': 3,
    'Turbidity (NTU)': 1,
    'Dissolved Oxygen (mg/L)': 5,
    'Nitrate Level (mg/L)': 50,
```

```

    'Lead Concentration (µg/L)': 10,
    'Bacteria Count (CFU/mL)': 0
}

# Compute global maxima and per-country means
max_vals = df[features].max()
means = df.groupby('Country')[features].mean()
normed = means.divide(max_vals, axis=1)

# Normalize the safe values
safe_norm = [ safe_vals[f] / max_vals[f] for f in features ]
safe_norm += safe_norm[:1]

# Preparing 7 axes
N = len(features)
angles = np.linspace(0, 2 * np.pi, N, endpoint=False).tolist()
angles += angles[:1]

# Plotting a 2x2 grid of radar charts
fig, axs = plt.subplots(2, 2, figsize=(12, 12), subplot_kw = {'polar': True})
axs = axs.flatten()

for ax, (c1, c2) in zip(axs, pairs):
    v1 = normed.loc[c1].tolist(); v1 += v1[:1]
    v2 = normed.loc[c2].tolist(); v2 += v2[:1]

    # Expected Values
    ax.plot(angles, safe_norm, '--', color = 'green', lw = 2, label = 'Safe_
↳Limits')

    # Country 1
    ax.plot(angles, v1, color = 'C0', lw = 2, label = c1)
    ax.fill(angles, v1, color = 'C0', alpha = 0.25)

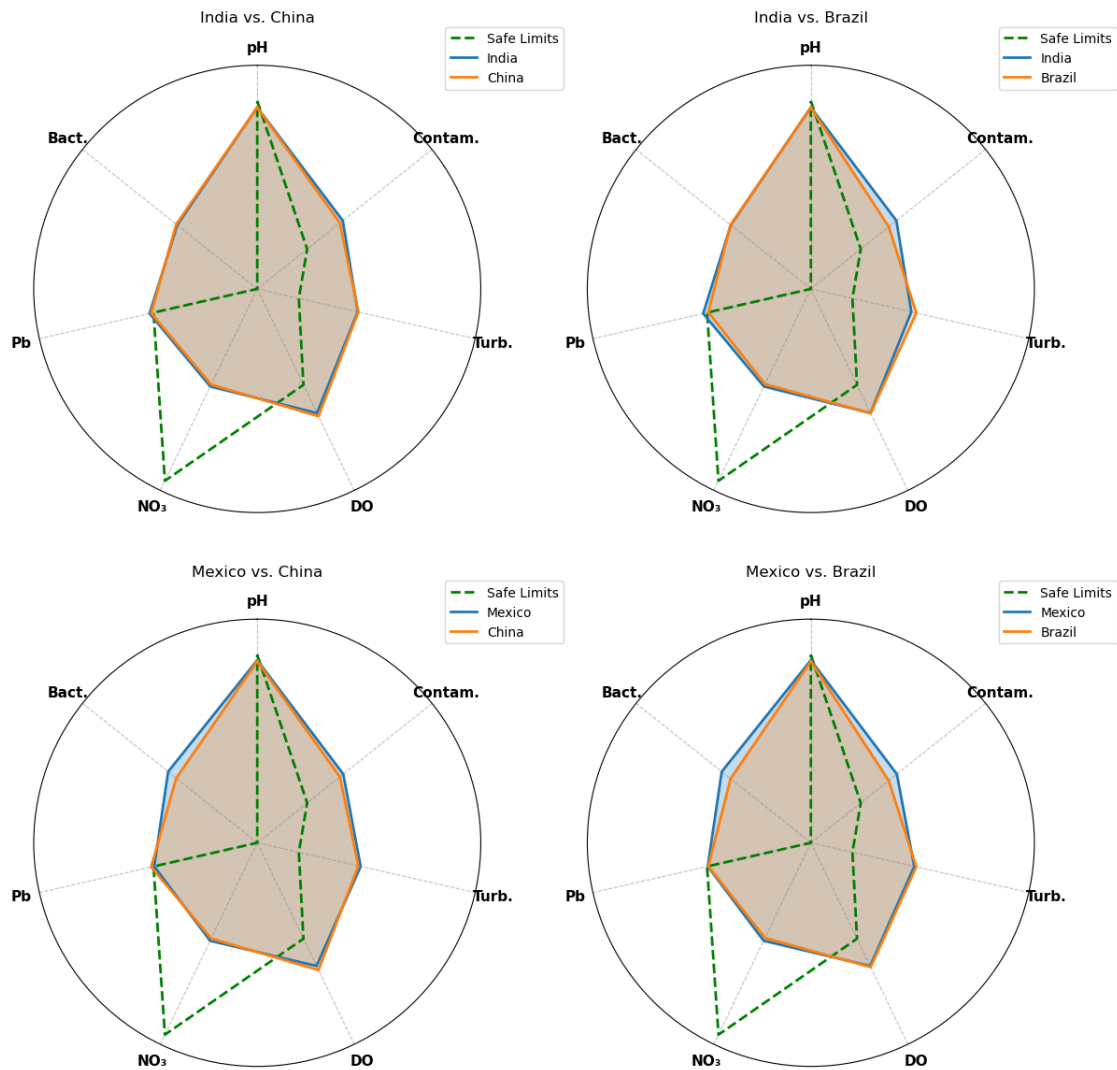
    # Country 2
    ax.plot(angles, v2, color = 'C1', lw = 2, label = c2)
    ax.fill(angles, v2, color = 'C1', alpha = 0.25)

    # Styling
    ax.set_theta_offset(np.pi / 2)
    ax.set_theta_direction(-1)
    ax.set_xticks(angles[:-1])
    ax.set_xticklabels(short_labels, fontsize = 11, fontweight = 'bold')
    ax.set_yticks([])
    ax.yaxis.grid(True, linestyle = '--', color = 'gray', alpha = 0.5)
    ax.xaxis.grid(True, linestyle = '--', color = 'gray', alpha = 0.5)
    ax.set_title(f"{c1} vs. {c2}", y = 1.08)

```

```
ax.legend(loc = 'upper right', bbox_to_anchor = (1.2, 1.1))

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



1.12 Radar-Chart Insights

Each subplot overlays three profiles on a shared 7-axis radar chart:

- **Green dashed heptagon:** WHO/EPA “safe” thresholds for each pollutant
- **Blue polygon:** “Bad” country’s normalized mean pollutant levels
- **Orange polygon:** “Fair” country’s corresponding profile

Any part of the polygon that extends beyond the green outline indicates a breach of the safe limit.

1.12.1 India vs. China

- **Dissolved Oxygen (DO):** India slightly underperforms, dipping below the safe threshold, whereas China meets the standard.
- **Turbidity & Bacteria:** Both countries are under the limit, but India's values are marginally worse.
- **pH, Contaminants, Lead, Nitrate:** These remain well within safe bounds for both countries.

Takeaway: India's shortfalls are not severe but show mild underperformance in DO and microbial quality, while China consistently aligns with safe thresholds.

1.12.2 India vs. Brazil

- **Dissolved Oxygen:** India again falls slightly short; Brazil meets the standard.
- **Bacteria Count:** India shows minor exceedance relative to Brazil.
- **All other features** stay within safe boundaries for both countries.

Takeaway: The most relevant differences are in oxygenation and microbial content, with Brazil holding slight but consistent advantages.

1.12.3 Mexico vs. China

- **Dissolved Oxygen:** Mexico falls just below safe levels, while China sits near the safe benchmark.
- **Bacteria:** A subtle overstep in Mexico compared to China's compliant value.
- **Other features**, including nitrates, fall within the safe zone.

Takeaway: Mexico's concern centers on biological pollution and slightly lower oxygen levels—issues China has managed more effectively.

1.12.4 Mexico vs. Brazil

- **Dissolved Oxygen:** Again, Mexico underperforms slightly versus Brazil.
- **Turbidity:** Slightly elevated for Mexico relative to Brazil, though still close to threshold.

- **Other metrics**, including nitrates and pH, remain within limits for both.

Takeaway: Brazil outperforms Mexico in physical water quality indicators like DO and turbidity, reinforcing its benchmark status.

1.13 Overall Conclusions

1. **Dissolved oxygen levels** are the most consistent underperformer for India and Mexico—indicating potential water stagnation or lack of aeration.
 2. **Bacteria and turbidity** appear as secondary factors where the “bad” countries show slight, but recurrent, deviations.
 3. **Nitrate levels** are not a primary concern in this dataset, with all countries generally falling within acceptable bounds.
 4. **China** and **Brazil** represent reliable benchmarks, especially in maintaining adequate DO and low microbial contamination.
 5. Intervention strategies for India and Mexico should prioritize **aeration** (to raise DO), and **microbial filtration or chlorination** to tackle biological loads.
-

1.14 Proposed Solution Framework

Building on our radar-chart diagnostics (which pinpointed **nitrate** and **dissolved-oxygen** as the universal failure modes), we now introduce a two-pronged modeling approach to answer:

- “Why are some countries failing?”
 - “Which peer nations should they study?”
-

1.14.1 1. Feature-Importance Modeling

Objective: Quantify how much each pollutant and each socioeconomic/infrastructure factor drives overall water-quality risk.

- **Inputs:**
 - Seven **mean raw-penalty features** (deviations from WHO/EPA thresholds)
 - Four socioeconomic variables: GDP per capita, Healthcare Access Index, Sanitation Coverage, Access to Clean Water, Urbanization Rate
- **Target:** Country’s **mean risk_score** (0–100)
- **Method:** Train a **Random Forest regressor** (or Gradient Boosting) on country-level data.
- **Output:**

1. **R²** on held-out test set, demonstrating predictive power
2. **Feature importances**, revealing whether **nitrate**, **dissolved-oxygen**, or factors like sanitation coverage are the strongest levers

Impact: We'll know whether investments in infrastructure (e.g. sanitation) or targeted pollutant controls (e.g. nitrate reduction) yield the greatest risk reduction.

1.14.2 2. Clustering & Peer Recommendation

Objective: Group countries into clusters by risk profile and socioeconomic context, then recommend “near-peer” nations to benchmark.

- **Inputs:**
 - Same seven penalty features + the four socioeconomic variables
- **Preprocessing:** Standard-scale all features to zero mean, unit variance
- **Method:**
 1. Use **K-Means** (select K via elbow/silhouette) to assign each country to a cluster
 2. Inspect **cluster centroids** and label profiles (e.g. “High-Risk & Low-Sanitation”)
 3. Implement `recommend_peers(country, n=3)`:
 - Identify the country’s cluster
 - Compute Euclidean distances within the cluster
 - Return the n closest peer countries
- **Output:** For any given underperforming nation (e.g. India), a ranked list of 2–3 exemplar peers whose success factors can be emulated.

Impact: Stakeholders gain **actionable guidance**: “Study Country X and Y—they face similar challenges but maintain lower risk through A, B, C interventions.”

1.15 End-to-End Story

1. **Diagnose** – Radar charts reveal *what* fails (nitrates, DO) and *how often* (frequency vs. severity).
2. **Quantify** – Feature-importance modeling explains *why* those failures correlate with water-quality risk (pollutants vs. infrastructure).
3. **Prescribe** – Clustering & peer-recommendation deliver *who* to learn from, offering a clear roadmap for policy and engineering interventions.

Together, this pipeline transforms raw water-quality data into a **full ML-powered policy toolkit**:
What's broken → Why it's broken → Who's doing it better.

```
[24]: # Features to aggregate at the country level
penalty_features = [
    'penalty_ph', 'penalty_contaminant', 'penalty_turbidity',
    'penalty_do', 'penalty_nitrate', 'penalty_lead', 'penalty_bacteria'
]
socio_features = [
    'GDP per Capita (USD)',
    'Healthcare Access Index (0-100)',
    'Sanitation Coverage (% of Population)',
    'Access to Clean Water (% of Population)',
    'Urbanization Rate (%)'
]

# Computing country level means
country_df = (df.groupby('Country')[penalty_features + socio_features +
    ↪['risk_score']].mean().reset_index())

# Display the country level dataset
print("Country level dataset (one row per country):")
display(country_df)

# Preprocessing
X = country_df[penalty_features + socio_features]
y = country_df['risk_score']
```

Country level dataset (one row per country):

	Country	penalty_ph	penalty_contaminant	penalty_turbidity	penalty_do	\
0	Bangladesh	0.053883	2.556529	1.642680	0.245773	
1	Brazil	0.052577	2.187045	1.615189	0.301168	
2	China	0.045563	2.361254	1.524566	0.267621	
3	Ethiopia	0.049249	2.367543	1.539454	0.278874	
4	India	0.050759	2.585069	1.514103	0.298724	
5	Indonesia	0.048191	2.198454	1.522303	0.275888	
6	Mexico	0.044549	2.612743	1.576354	0.259340	
7	Nigeria	0.054349	2.535333	1.512508	0.366667	
8	Pakistan	0.046611	2.205805	1.698557	0.275973	
9	USA	0.045831	2.452163	1.593260	0.289498	

	penalty_nitrate	penalty_lead	penalty_bacteria	GDP per Capita (USD)	\
0	0.0	2.761409	2.491282	48973.381443	
1	0.0	2.396770	2.412045	50334.168385	
2	0.0	2.607428	2.433441	50368.745981	
3	0.0	2.622526	2.510628	51939.931741	
4	0.0	2.669586	2.403955	47883.120690	

5	0.0	2.660592	2.532332	51554.036184
6	0.0	2.602014	2.676330	49210.819444
7	0.0	2.270476	2.482790	49005.730159
8	0.0	2.697886	2.488809	51556.976510
9	0.0	2.245862	2.459708	49514.047022

	Healthcare Access Index (0-100)	Sanitation Coverage (% of Population)	\
0	48.976804		61.794605
1	49.640756		61.501100
2	48.837910		62.753505
3	49.038259		60.297167
4	50.719276		59.716552
5	50.394211		60.146941
6	50.707083		57.261701
7	52.267587		60.439079
8	50.055503		60.785906
9	49.593009		58.947367

	Access to Clean Water (% of Population)	Urbanization Rate (%)	risk_score
0	64.401478	49.295533	40.110450
1	64.156289	52.499966	36.813451
2	66.532540	49.933023	37.966188
3	63.572867	52.200341	38.504266
4	65.445069	50.344552	39.149296
5	65.703816	50.681053	37.957339
6	65.428194	48.241111	40.193311
7	63.540159	48.543841	37.891815
8	63.391342	49.808423	38.694384
9	63.968966	49.236552	37.322729

```
[25]: # Train/Test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)
```

```
# Train Random Forest Regressor
model = RandomForestRegressor(n_estimators = 200, random_state = 42)
model.fit(X_train, y_train)
```

```
# Evaluate
r2 = model.score(X_test, y_test)
print(f" Test R²: {r2:.4f}")
```

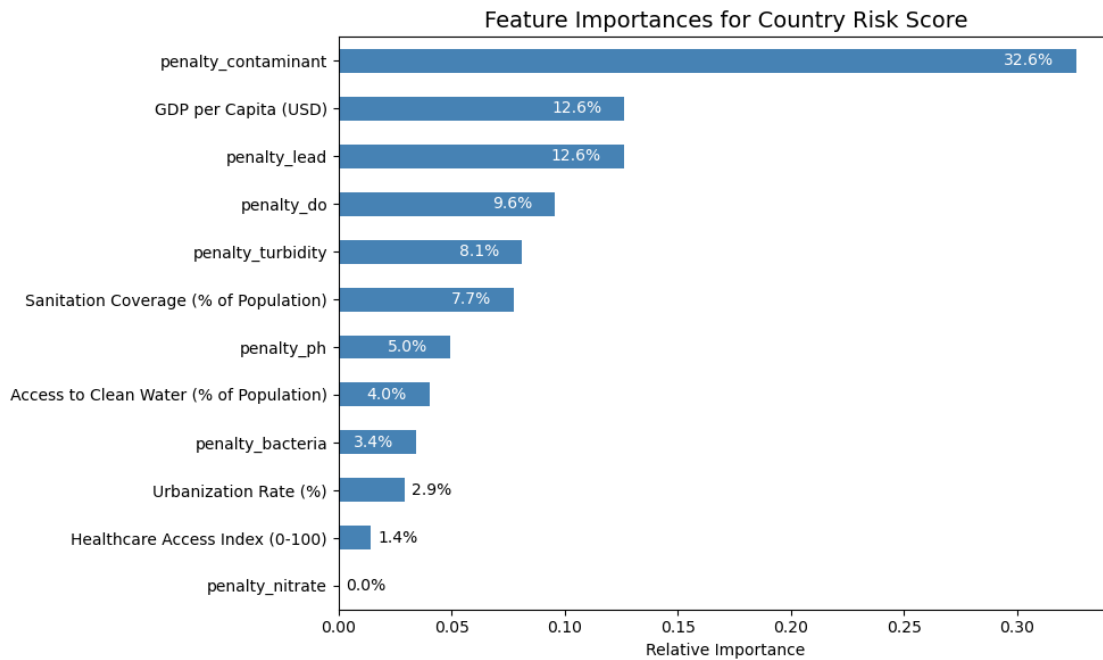
Test R²: 0.0452

```
[26]: # Feature Importances
importances = pd.Series(model.feature_importances_, index = X.columns)
importances = importances.sort_values(ascending = False)
```

```

# Plotting
plt.figure(figsize=(10, 6))
importances.sort_values().plot.barh(color='#4682B4')
for i, (value, name) in enumerate(zip(importances.sort_values(), importances.
    ↪sort_values().index)):
    if value > 0.03:
        plt.text(value - 0.01, i, f"{value:.1%}", va='center', ha='right', ↪
    ↪color='white', fontsize=10)
    else:
        plt.text(value + 0.003, i, f"{value:.1%}", va='center', ha='left', ↪
    ↪color='black', fontsize=10)
plt.title("Feature Importances for Country Risk Score", fontsize=14)
plt.xlabel("Relative Importance")
plt.tight_layout()
plt.show()

```



1.15.1 Key Drivers of Country-Level Risk

Our Random Forest model reveals that **chemical contamination** and **infrastructure inequality** are the most important factors in predicting a country's water risk.

- **Penalty for contaminants (32.6%)**: This is the dominant risk driver, suggesting that systemic water pollution from chemical agents remains a critical issue.
- **GDP per Capita and Lead Penalty (12.6%)**: Suggest wealthier nations are better able to control water pollutants, especially heavy metals.
- **Sanitation Coverage (7.7%)**: Points to the importance of public infrastructure in reducing

water-borne disease risk.

Surprisingly, **nitrate levels** and **healthcare access** contribute very little to national risk variation.

```
[27]: # Reinitializing the features
penalty_features = [
    'penalty_ph', 'penalty_contaminant', 'penalty_turbidity',
    'penalty_do', 'penalty_nitrate', 'penalty_lead', 'penalty_bacteria'
]
socio_features = [
    'GDP per Capita (USD)',
    'Healthcare Access Index (0-100)',
    'Sanitation Coverage (% of Population)',
    'Access to Clean Water (% of Population)',
    'Urbanization Rate (%)'
]
all_features = penalty_features + socio_features
# Computing country level means
country_df = (df.groupby('Country')[penalty_features + socio_features +
    ↪['risk_score']].mean().reset_index())

# Standardize & Cluster
scaler = StandardScaler()
Z = scaler.fit_transform(country_df[all_features])
kmeans = KMeans(n_clusters = 3, random_state = 0, n_init = 'auto')
country_df['cluster'] = kmeans.fit_predict(Z)

# Define peer recommender
def recommend_peers(country, n=2):
    if country not in country_df['Country'].values:
        return f"Country '{country}' not found in dataset."

    idx = country_df[country_df['Country'] == country].index[0]
    cluster_id = country_df.loc[idx, 'cluster']
    cluster_members = country_df[country_df['cluster'] == cluster_id]
    cluster_members = cluster_members[cluster_members['Country'] != country]

    distances = cdist(
        [Z[idx]],
        scaler.transform(cluster_members[all_features])
    )[0]
    cluster_members = cluster_members.copy()
    cluster_members['distance'] = distances
    return cluster_members.sort_values('distance')[['Country', 'risk_score']].
    ↪head(n)

def interpret_peer_gap(base_country, peer_country):
    base = country_df[country_df['Country'] == base_country]
```

```

peer = country_df[country_df['Country'] == peer_country]

gap = (peer[all_features].values - base[all_features].values).flatten()
s = pd.Series(gap, index=all_features).sort_values(key=abs, ascending=False)

for feat in s.index[:5]:
    val = s[feat]
    direction = "subceeds" if val > 0 else "exceeds"
    print(f"{base_country} {direction} from {peer_country} in {feat} by_
↳about {abs(val):.2f}.")

# Using our early model on importances and this model on clusters, we can now_
↳recommend peers and interpret the gap.
# Defining gap computation between country and peer
def peer_priority_gaps(country, peer):
    c_row = country_df[country_df['Country'] == country][importances.index]
    p_row = country_df[country_df['Country'] == peer][importances.index]
    gaps = (p_row.values - c_row.values).flatten()
    return pd.Series(gaps, index=importances.index)

# Computing weighted priority scores (importance × deviation)
def prioritized_recommendations(country, peer, top_n=5):
    gaps = peer_priority_gaps(country, peer)
    priorities = importances * abs(gaps)
    ranked = priorities.sort_values(ascending=False)

    result = pd.DataFrame({
        'Gap (Peer - Country)': gaps[ranked.index],
        'Feature Importance': importances[ranked.index],
        'Priority Score': ranked
    })
    return result.head(top_n)

test = 'y'
while test == 'y':
    country = input("Enter a country to test: ")
    peers = recommend_peers(country)
    print(f"\n Recommended peers for {country}:")
    display(peers)
    for peer in peers['Country']:
        print("\n Interpreting difference from", peer)
        interpret_peer_gap(country, peer)
        print("\n Prioritized recommendations:")
        display(prioritized_recommendations(country, peer))
    test = input("Do you want to test again? (y/n): ")

```

Recommended peers for India:

	Country	risk_score
9	USA	37.322729
0	Bangladesh	40.110450

Interpreting difference from USA
 India subceeds from USA in GDP per Capita (USD) by about 1630.93.
 India exceeds from USA in Access to Clean Water (% of Population) by about 1.48.
 India exceeds from USA in Healthcare Access Index (0-100) by about 1.13.
 India exceeds from USA in Urbanization Rate (%) by about 1.11.
 India exceeds from USA in Sanitation Coverage (% of Population) by about 0.77.

Prioritized recommendations:

	Gap (Peer - Country) \	Feature Importance	Priority Score
GDP per Capita (USD)	1630.926332	0.126241	205.890325
Sanitation Coverage (% of Population)	-0.769185	0.077375	0.059516
Access to Clean Water (% of Population)	-1.476103	0.040179	0.059308
penalty_lead	-0.423724	0.126088	0.053427
penalty_contaminant	-0.132906	0.326001	0.043327

Interpreting difference from Bangladesh
 India subceeds from Bangladesh in GDP per Capita (USD) by about 1090.26.
 India subceeds from Bangladesh in Sanitation Coverage (% of Population) by about 2.08.
 India exceeds from Bangladesh in Healthcare Access Index (0-100) by about 1.74.
 India exceeds from Bangladesh in Urbanization Rate (%) by about 1.05.
 India exceeds from Bangladesh in Access to Clean Water (% of Population) by about 1.04.

Prioritized recommendations:

	Gap (Peer - Country) \	Feature Importance	Priority Score
GDP per Capita (USD)	1090.260754	0.126241	137.635978
Sanitation Coverage (% of Population)	2.078053	0.077375	0.160790
Access to Clean Water (% of Population)	-1.043591		
Urbanization Rate (%)	-1.049019		
Healthcare Access Index (0-100)	-1.742472		

Access to Clean Water (% of Population)	0.040179	0.041930
Urbanization Rate (%)	0.029162	0.030591
Healthcare Access Index (0-100)	0.014314	0.024943

1.16 Final Conclusion

This project set out to model, diagnose, and solve global water pollution disparities by analyzing a multi-dimensional dataset of water quality, health outcomes, and national infrastructure indicators. Through a structured pipeline, we achieved not only country-level risk quantification but also concrete, interpretable recommendations for mitigation.

1.16.1 What We Built

1. Risk Quantification via WHO Thresholds

We designed a distance-based penalty scoring system that measures how far each sample deviates from WHO/EPA safe thresholds across seven key pollutants. This raw score was then scaled to a 0–100 “risk score” for interpretability.

2. Country-Level Profiling

Risk scores were aggregated to compute country-level means and identify outliers — i.e., “Doing Bad” (e.g. India, Mexico) vs. “Doing Fair” (e.g. China, Brazil).

3. Radar Chart Visualization

We created side-by-side radar charts comparing “bad” countries against “fair” ones across the 7 pollutants. A reference line for safe limits allowed visual diagnosis of pollutant-specific failures.

4. Feature Importance Modeling with Random Forest

A `RandomForestRegressor` modeled country-level risk scores using pollutant penalties and infrastructure features. This gave us interpretable feature importances:

- Chemical contaminants were the largest contributor to risk.
- Infrastructure (GDP, sanitation) explained meaningful variance.
- Healthcare and nitrate penalties had surprisingly low influence.

5. Peer-Based Recommendations

We clustered countries based on their full profile (pollution + infrastructure) using KMeans. Then for any country, we:

- Identified closest-performing peers using Euclidean distance
- Measured feature-by-feature gaps
- Prioritized interventions using:
`priority = |gap| × feature_importance`

This gave interpretable suggestions like:

> India should work on it's GDP per capita and invest in sanitation infrastructure to close the most impactful gaps compared to Bangladesh.

1.16.2 Limitations

The dataset used in this project exhibited clear signs of being synthetic or corrupted in certain columns: - **GDP per Capita** for the USA ranged from \$521 to \$99,000
- **Bangladesh** appeared to outperform **India** economically in raw averages

While this undermines the real-world accuracy of absolute scores, the **logic and code pipeline remain valid and generalizable**. If applied to verified global data (e.g. World Bank, WHO), this framework could produce accurate, policy-relevant insights.

1.16.3 Final Thoughts

This project demonstrates a complete ML-driven policy toolchain: - Score → Diagnose → Recommend

It is not just about detecting which countries are at risk, but also prescribing how they can improve — and from whom they should learn. With cleaner data, the same system could power real public health decisions and dashboards.

1.17 Future Implementation and Deployment

The analytical framework developed in this project lays the groundwork for a scalable, user-facing platform. By integrating the model into a modern web environment using JavaScript libraries such as **React.js**, it would be possible to build an interactive interface where users — such as public policy analysts or government officials — can input a country’s pollution and infrastructure profile to receive real-time peer assessments and actionable recommendations.

Such a deployment could support dynamic visualizations (e.g., radar charts, risk timelines), model-driven suggestions, and peer benchmarking dashboards. With proper backend integration and data validation, the platform could serve as a robust decision-support system for environmental monitoring, public health strategy, and international development planning.