

01_data_exploration

November 29, 2024

```
[ ]: # Import required libraries
import warnings
from pathlib import Path

import matplotlib
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import plotly.express as px
import plotly.graph_objects as go
import seaborn as sns
from IPython.display import display, HTML

matplotlib.use('Agg') # Use non-interactive backend for PDF export

# Set up the output directory for saving figures
notebook_dir = Path().absolute()
project_root = notebook_dir.parent if notebook_dir.name == 'notebooks' else_
↳ notebook_dir
figures_dir = project_root / 'figures'
exploration_dir = figures_dir / 'exploration'
exploration_dir.mkdir(parents=True, exist_ok=True)

# Create directories
(figures_dir / 'exploration').mkdir(parents=True, exist_ok=True)
(figures_dir / 'feature_analysis').mkdir(parents=True, exist_ok=True)

# Suppress warnings
warnings.filterwarnings('ignore')

# Set plotting styles
plt.style.use('bmh') # Using a built-in style instead of seaborn
sns.set_palette("husl")
plt.rcParams['figure.figsize'] = [12, 6]

# Suppress warnings
warnings.filterwarnings('ignore')
```

```

# Set plotting styles
plt.style.use('bmh') # Using a built-in style instead of seaborn
sns.set_palette("husl")
plt.rcParams['figure.figsize'] = [12, 6]

# Load processed data
# Get the current notebook directory and construct the correct path
notebook_dir = Path().absolute()
project_root = notebook_dir.parent if notebook_dir.name == 'notebooks' else_
↳notebook_dir
processed_data_path = project_root / 'processed_data' / 'final_processed_data.
↳csv'

print(f"Looking for data file at: {processed_data_path}")
df = pd.read_csv(processed_data_path)

print("\nDataset Overview:")
print("=" * 80)
print(f"\nShape: {df.shape}")
print("\nFeatures:")
for col in df.columns:
    dtype = df[col].dtype
    missing = df[col].isnull().sum()
    print(f"- {col}: {dtype} (Missing: {missing})")

```

```

[ ]: # Load the datasets
def load_datasets():
    """Load all relevant datasets"""
    # Get the current notebook directory and construct the correct path
    notebook_dir = Path().absolute()
    project_root = notebook_dir.parent if notebook_dir.name == 'notebooks' else_
↳notebook_dir
    base_path = project_root / "data"

    print(f>Loading data from: {base_path}")

    # Global Energy Consumption & Renewable Generation
    global_energy_path = base_path / "Global Energy Consumption & Renewable_
↳Generation"
    print(f>Checking global energy path: {global_energy_path}")
    print(f>Path exists: {global_energy_path.exists()}")

    global_data = {
        'continent_consumption': pd.read_csv(global_energy_path /_
↳"Continent_Consumption_TWH.csv"),

```

```

        'country_consumption': pd.read_csv(global_energy_path /
↪ "Country_Consumption_TWH.csv"),
        'renewable_gen': pd.read_csv(global_energy_path /
↪ "renewablePowerGeneration97-17.csv"),
        'nonrenewable_gen': pd.read_csv(
            global_energy_path / "nonRenewablesTotalPowerGeneration.csv")
    }

    # Worldwide Renewable Data
    worldwide_path = base_path / "Renewable Energy World Wide 1965-2022"
    worldwide_data = {
        'renewable_share': pd.read_csv(worldwide_path / "01_
↪ renewable-share-energy.csv"),
        'renewable_consumption': pd.read_csv(
            worldwide_path / "02 modern-renewable-energy-consumption.csv"),
        'hydro_consumption': pd.read_csv(worldwide_path / "05_
↪ hydropower-consumption.csv"),
        'wind_generation': pd.read_csv(worldwide_path / "08 wind-generation.
↪ csv"),
        'solar_consumption': pd.read_csv(worldwide_path / "12_
↪ solar-energy-consumption.csv")
    }

    # Weather and US Data
    weather_data = pd.read_csv(base_path /
↪ "renewable_energy_and_weather_conditions.csv")
    us_data = pd.read_csv(base_path / "us_renewable_energy_consumption.csv")

    return global_data, worldwide_data, weather_data, us_data

# Print current working directory and verify paths
print("Current working directory:", Path().absolute())
print("\nTrying to load datasets...")
global_data, worldwide_data, weather_data, us_data = load_datasets()
print("\nDatasets loaded successfully!")

# Load datasets
global_data, worldwide_data, weather_data, us_data = load_datasets()

```

```

[ ]: # Initial Data Overview
def display_dataset_info(data_dict, title):
    """Display basic information about datasets"""
    print(f"\n{title}")
    print("=" * 80)
    for name, df in data_dict.items():

```

```

print(f"\nDataset: {name}")
print(f"Shape: {df.shape}")
print("\nColumns:")
for col in df.columns:
    dtype = df[col].dtype
    missing = df[col].isnull().sum()
    print(f"- {col}: {dtype} (Missing: {missing})")
print("-" * 40)

# Display information for each dataset group
display_dataset_info(global_data, "Global Energy Consumption & Renewable_
↳Generation Datasets")
display_dataset_info(worldwide_data, "Worldwide Renewable Energy Datasets")
print("\nWeather Conditions Dataset")
print("=" * 80)
display(weather_data.info())
print("\nUS Renewable Energy Dataset")
print("=" * 80)
display(us_data.info()) # Cell 3: Initial Data Overview

def display_dataset_info(data_dict, title):
    """Display basic information about datasets"""
    print(f"\n{title}")
    print("=" * 80)
    for name, df in data_dict.items():
        print(f"\nDataset: {name}")
        print(f"Shape: {df.shape}")
        print("\nColumns:")
        for col in df.columns:
            dtype = df[col].dtype
            missing = df[col].isnull().sum()
            print(f"- {col}: {dtype} (Missing: {missing})")
        print("-" * 40)

# Display information for each dataset group
display_dataset_info(global_data, "Global Energy Consumption & Renewable_
↳Generation Datasets")
display_dataset_info(worldwide_data, "Worldwide Renewable Energy Datasets")
print("\nWeather Conditions Dataset")
print("=" * 80)
display(weather_data.info())
print("\nUS Renewable Energy Dataset")
print("=" * 80)
display(us_data.info())

```

```
[ ]: # Data Quality Assessment
def assess_data_quality(data_dict, title):
    """Assess data quality for each dataset"""
    print(f"\n{title}")
    print("=" * 80)

    for name, df in data_dict.items():
        print(f"\nDataset: {name}")

        # Missing values
        missing = df.isnull().sum()
        if missing.any():
            print("\nMissing Values:")
            print(missing[missing > 0])

        # Duplicates
        duplicates = df.duplicated().sum()
        print(f"\nDuplicate Rows: {duplicates}")

        # Basic statistics
        print("\nNumerical Columns Statistics:")
        print(df.describe().round(2))

        print("-" * 40)

    # Assess data quality for each dataset group
    assess_data_quality(global_data, "Global Energy Data Quality Assessment")
    assess_data_quality(worldwide_data, "Worldwide Renewable Data Quality_
↪Assessment")

    print("\nWeather Data Quality Assessment")
    print("=" * 80)
    display(weather_data.describe())
    print("\nUS Data Quality Assessment")
    print("=" * 80)
    display(us_data.describe())

[ ]: def plot_time_series(df, x_col, y_col, title, hue=None):
    """
    Create time series plot using plotly with case-insensitive column matching

    Args:
        df: DataFrame to plot
        x_col: Name of x-axis column
        y_col: Name of y-axis column
        title: Plot title
```

```

        hue: Column name for color grouping
    """
    # Print available columns for debugging
    print(f"Available columns: {list(df.columns)}")

    # Create a copy to avoid modifying original
    plot_df = df.copy()

    # Find actual column names (case-insensitive)
    x_col_actual = next((col for col in df.columns if col.lower() == x_col.
↪lower()), None)
    y_col_actual = next((col for col in df.columns if col.lower() == y_col.
↪lower()), None)
    hue_actual = next((col for col in df.columns if col and col.lower() == hue.
↪lower()),
                       None) if hue else None

    if not x_col_actual:
        raise ValueError(f"Column '{x_col}' not found. Available columns:␣
↪{list(df.columns)}")
    if not y_col_actual:
        raise ValueError(f"Column '{y_col}' not found. Available columns:␣
↪{list(df.columns)}")
    if hue and not hue_actual:
        raise ValueError(f"Column '{hue}' not found. Available columns:␣
↪{list(df.columns)}")

    # Create the plot
    fig = px.line(plot_df,
                  x=x_col_actual,
                  y=y_col_actual,
                  title=title,
                  color=hue_actual if hue else None)

    fig.update_layout(
        xaxis_title=x_col,
        yaxis_title=y_col,
        template='plotly_white'
    )

    filename = f"{title.lower().replace(' ', '_').replace('(', '').replace(')', '␣
↪')}.png"
    fig.write_image(str(exploration_dir / filename))
    fig.show()

```

```

# Print data information before plotting
print("\nGlobal Data - Renewable Generation:")
print(global_data['renewable_gen'].head())
print("\nColumns:", list(global_data['renewable_gen'].columns))

print("\nWorldwide Data - Renewable Share:")
print(worldwide_data['renewable_share'].head())
print("\nColumns:", list(worldwide_data['renewable_share'].columns))

# Plot renewable generation trends
print("\nPlotting renewable generation trends...")
plot_time_series(
    global_data['renewable_gen'],
    'Year', # Changed from 'year' to 'Year'
    'Hydro(TWh)', # Using an actual column name
    'Renewable Power Generation Trends (1997-2017)'
)

# Plot renewable share evolution
print("\nPlotting renewable share evolution...")
plot_time_series(
    worldwide_data['renewable_share'],
    'Year',
    'Renewables (% equivalent primary energy)',
    'Evolution of Renewable Energy Share (1965-2022)',
    hue='Entity'
)

# Create some additional plots to show different aspects of the data
print("\nPlotting solar and wind generation trends...")
if 'Solar PV (TWh)' in global_data['renewable_gen'].columns:
    plot_time_series(
        global_data['renewable_gen'],
        'Year',
        'Solar PV (TWh)',
        'Solar Power Generation Trends (1997-2017)'
    )

if 'wind_generation' in worldwide_data:
    plot_time_series(
        worldwide_data['wind_generation'],
        'Year',
        'Electricity from wind (TWh)',
        'Wind Power Generation Trends',
        hue='Entity'
    )

```

```
[ ]: # Geographic Distribution Analysis
def plot_choropleth(df, color_col, title):
    """
    Create choropleth map using plotly

    Args:
        df: DataFrame containing the data
        color_col: Column containing values to plot
        title: Plot title
    """
    # Print data info for debugging
    print(f"\nCreating choropleth for {color_col}")
    print(f"Available columns: {list(df.columns)}")
    print(f"Sample data:\n{df.head()}")

    # Melt the dataframe to get country-wise data
    # Convert wide format (countries as columns) to long format
    melted_df = df.melt(
        id_vars=['Year'],
        var_name='Country',
        value_name='Generation' # Use a generic name instead of the column name
    )

    print(f"\nMelted data sample:\n{melted_df.head()}")

    # Filter to only the data we want to plot
    plot_data = melted_df[melted_df['Country'] == color_col].copy()

    # Create the choropleth map
    fig = px.choropleth(
        plot_data,
        locations='Country',
        locationmode='country names',
        color='Generation',
        hover_name='Country',
        title=title,
        color_continuous_scale='Viridis'
    )

    # Update layout
    fig.update_layout(
        template='plotly_white',
        title_x=0.5, # Center the title
        margin=dict(l=0, r=0, t=30, b=0)
    )

    filename = f"choropleth_{title.lower().replace(' ', '_')}.png"
```



```

fig.write_image(str(exploration_dir / filename))
fig.show()

# Print information about the renewable generation data
print("Renewable Generation Data Info:")
print("\nColums:", list(global_data['renewable_gen'].columns))
print("\nSample Data:")
print(global_data['renewable_gen'].head())

# Get the latest year data
latest_year = global_data['renewable_gen']['Year'].max()
print(f"\nLatest year in data: {latest_year}")

# Get renewable energy columns (exclude 'Year' column)
renewable_cols = [col for col in global_data['renewable_gen'].columns if col !=
↳ 'Year']

# Create summary dataframe for the latest year
latest_data = global_data['renewable_gen'][
    global_data['renewable_gen']['Year'] == latest_year].copy()

# Create bar chart showing total generation by type
generation_by_type = latest_data[renewable_cols].sum()
fig = px.bar(
    x=generation_by_type.index,
    y=generation_by_type.values,
    title=f'Total Renewable Energy Generation by Type ({latest_year})'
)
fig.update_layout(xaxis_tickangle=-45, showlegend=False)
fig.write_image(str(exploration_dir / 'total_generation_by_type.png'))
fig.show()

# Create pie chart showing energy mix
fig = px.pie(
    values=generation_by_type.values,
    names=generation_by_type.index,
    title=f'Global Renewable Energy Mix ({latest_year})'
)
fig.write_image(str(exploration_dir / 'global_renewable_mix.png'))
fig.show()

# Create bar chart showing generation over time
yearly_totals = global_data['renewable_gen'].groupby('Year')[renewable_cols].
↳ sum()
fig = px.line(
    yearly_totals,

```

```

        title='Renewable Energy Generation Over Time'
    )
    fig.update_layout(
        xaxis_title='Year',
        yaxis_title='Generation (TWh)',
        showlegend=True
    )
    fig.write_image(str(exploration_dir / 'generation_over_time.png'))
    fig.show()

    print("\nVisualization Summary:")
    print(f"- Data covers years from {global_data['renewable_gen']['Year'].min()}_
    to {latest_year}")
    print(f"- Total types of renewable energy tracked: {len(renewable_cols)}")
    print("- Energy types:", renewable_cols)

```

```

[ ]: # Weather Impact Analysis
def analyze_weather_impact():
    """Analyze the impact of weather conditions on renewable energy"""
    # First, let's examine the data
    print("Weather Data Info:")
    print("\nColumns:", list(weather_data.columns))
    print("\nData Types:")
    print(weather_data.dtypes)

    # Convert Time column to datetime if it isn't already
    weather_df = weather_data.copy()
    weather_df['Time'] = pd.to_datetime(weather_df['Time'])

    # Select only numeric columns for correlation analysis
    numeric_cols = weather_df.select_dtypes(include=[np.number]).columns
    print("\nNumeric columns for analysis:", list(numeric_cols))

    # Calculate correlations for numeric columns
    weather_corr = weather_df[numeric_cols].corr()

    # Plot correlation heatmap
    plt.figure(figsize=(15, 12))
    sns.heatmap(weather_corr,
                annot=True,
                cmap='coolwarm',
                center=0,
                fmt='.2f',
                square=True)
    plt.title('Correlation between Weather Variables')
    plt.xticks(rotation=45, ha='right')
    plt.yticks(rotation=0)

```

```

plt.tight_layout()
# Save correlation heatmap
plt.figure(figsize=(15, 12))
sns.heatmap(weather_corr, annot=True, cmap='coolwarm', center=0, fmt='.2f',
↪square=True)
plt.title('Correlation between Weather Variables')
plt.xticks(rotation=45, ha='right')
plt.yticks(rotation=0)
plt.tight_layout()
plt.savefig(exploration_dir / 'weather_correlation.png', dpi=300,
↪bbox_inches='tight')
plt.show()

# Select key variables for scatter matrix
key_vars = ['temp', 'wind_speed', 'GHI'] # Adjust these based on actual
↪column names
if 'Energy delta[Wh]' in weather_df.columns:
    key_vars.append('Energy delta[Wh]')

print("\nCreating scatter matrix for variables:", key_vars)

# Create scatter matrix for key relationships
fig = px.scatter_matrix(
    weather_df,
    dimensions=key_vars,
    title='Relationships between Key Weather Variables'
)
fig.update_layout(
    title_x=0.5,
    title_y=0.95
)
# Save scatter matrix
fig = px.scatter_matrix(
    weather_df,
    dimensions=key_vars,
    title='Relationships between Key Weather Variables'
)
fig.write_image(str(exploration_dir / 'weather_relationships.png'))
fig.show()

# Time series analysis
# Group by hour of day to see daily patterns
weather_df['hour'] = weather_df['Time'].dt.hour
hourly_avg = weather_df.groupby('hour')[key_vars].mean()

# Plot daily patterns
fig = go.Figure()

```

```

    for col in key_vars:
        fig.add_trace(go.Scatter(x=hourly_avg.index, y=hourly_avg[col],
↪name=col))
    fig.update_layout(
        title='Average Daily Patterns of Weather Variables',
        xaxis_title='Hour of Day',
        yaxis_title='Value',
        hovermode='x'
    )
    fig.write_image(str(exploration_dir / 'daily_weather_patterns.png'))
    fig.show()

# Monthly patterns
    weather_df['month'] = weather_df['Time'].dt.month
    monthly_avg = weather_df.groupby('month')[key_vars].mean()

    fig = go.Figure()
    for col in key_vars:
        fig.add_trace(go.Scatter(x=monthly_avg.index, y=monthly_avg[col],
↪name=col))
    fig.update_layout(
        title='Average Monthly Patterns of Weather Variables',
        xaxis_title='Month',
        yaxis_title='Value',
        hovermode='x'
    )
    fig.write_image(str(exploration_dir / 'monthly_weather_patterns.png'))
    fig.show()

# Print summary statistics
    print("\nSummary Statistics:")
    print(weather_df[key_vars].describe())

# Calculate and print key findings
    print("\nKey Findings:")
    for var1 in key_vars:
        for var2 in key_vars:
            if var1 < var2: # Avoid duplicate combinations
                corr = weather_df[var1].corr(weather_df[var2])
                print(f"Correlation between {var1} and {var2}: {corr:.2f}")

# Run the analysis
    print("Starting weather impact analysis...")
    analyze_weather_impact()

```

```
[ ]: # Energy Mix Analysis
def analyze_energy_mix():
    """Analyze the composition of energy sources"""
    # First, let's examine the data structure
    print("Renewable Generation Data Columns:")
    print(global_data['renewable_gen'].columns)
    print("\nNon-renewable Generation Data Columns:")
    print(global_data['nonrenewable_gen'].columns)
    print("\nRenewable Consumption Data Columns:")
    print(worldwide_data['renewable_consumption'].columns)

    # Calculate total renewable generation (sum all TWh columns)
    renewable_cols = [col for col in global_data['renewable_gen'].columns if
↳ 'TWh' in col]
    renewable_total = global_data['renewable_gen'][renewable_cols].sum().sum()

    # Get non-renewable total
    if 'Contribution (TWh)' in global_data['nonrenewable_gen'].columns:
        nonrenewable_total = global_data['nonrenewable_gen']['Contribution_
↳ (TWh)'].sum()
    else:
        print("\nWarning: Could not find non-renewable generation column")
        nonrenewable_total = 0

    print(f"\nTotal Renewable Generation: {renewable_total:.2f} TWh")
    print(f"Total Non-renewable Generation: {nonrenewable_total:.2f} TWh")

    # Create pie chart for total energy mix
    fig = go.Figure(data=[go.Pie(
        labels=['Renewable', 'Non-Renewable'],
        values=[renewable_total, nonrenewable_total],
        hole=0.4
    )])
    fig.update_layout(title='Global Energy Mix')
    fig.write_image(str(exploration_dir / 'global_energy_mix.png'))
    fig.show()

    # Analyze renewable energy composition
    print("\nAnalyzing renewable energy composition...")

    # Create a year-by-year analysis of renewable sources
    yearly_renewable = global_data['renewable_gen'].
↳ groupby('Year')[renewable_cols].sum()

    # Create a stacked area chart for renewable composition

    # Create pie chart for renewable mix in latest year
```

```

latest_year = yearly_renewable.index.max()
latest_mix = yearly_renewable.loc[latest_year]

fig = px.area(
    yearly_renewable,
    title='Evolution of Renewable Energy Composition'
)
fig.write_image(str(exploration_dir / 'renewable_composition_evolution.
↳png'))
fig.show()

# Save renewable mix pie chart
fig = go.Figure(data=[go.Pie(
    labels=latest_mix.index,
    values=latest_mix.values,
    hole=0.4
)])
fig.write_image(str(exploration_dir / f'renewable_mix_{latest_year}.png'))
fig.show()

# Calculate and display summary statistics
print(f"\nRenewable Energy Mix Analysis for {latest_year}:")
for source in latest_mix.index:
    percentage = (latest_mix[source] / latest_mix.sum()) * 100
    print(f"{source}: {latest_mix[source]:.0f} TWh ({percentage:.1f}%)")

# Calculate growth rates
growth_rates = yearly_renewable.pct_change().mean() * 100
print("\nAverage Annual Growth Rates:")
for source in growth_rates.index:
    print(f"{source}: {growth_rates[source]:.1f}% per year")

# Run the analysis
print("Starting energy mix analysis...")
analyze_energy_mix()

```

```

[ ]: # Statistical Analysis
def perform_statistical_analysis():
    """Perform statistical analysis on the datasets"""
    # First, let's examine the data structure
    print("Renewable Generation Data Structure:")
    print("\nColumns:", list(global_data['renewable_gen'].columns))
    print("\nSample data:")
    print(global_data['renewable_gen'].head())

    # Get renewable energy columns

```

```

renewable_cols = [col for col in global_data['renewable_gen'].columns if
↳ 'TWh' in col]
print("\nAnalyzing columns:", renewable_cols)

# Time series analysis for each type
yearly_data = global_data['renewable_gen'].copy()

# Growth rates analysis
growth_rates = pd.DataFrame()
for col in renewable_cols:
    growth_rates[col] = yearly_data[col].pct_change() * 100

print("\nGrowth Rates Statistics (%):")
print(growth_rates.describe().round(2))

# Variance analysis
variance_analysis = pd.DataFrame({
    'mean': yearly_data[renewable_cols].mean(),
    'std': yearly_data[renewable_cols].std(),
    'var': yearly_data[renewable_cols].var(),
    'cv': yearly_data[renewable_cols].std() / yearly_data[renewable_cols].
↳ mean() * 100
    # Coefficient of variation
}).sort_values('var', ascending=False)

print("\nVariance Analysis:")
display(variance_analysis)

# Distribution analysis
plt.figure(figsize=(15, 10))

# Create subplots for each renewable type
rows = (len(renewable_cols) + 1) // 2 # Calculate number of rows needed
fig, axes = plt.subplots(rows, 2, figsize=(15, 5 * rows))
axes = axes.flatten() # Flatten axes array for easier indexing

for idx, col in enumerate(renewable_cols):
    if idx < len(axes):
        sns.histplot(data=yearly_data, x=col, ax=axes[idx])
        axes[idx].set_title(f'Distribution of {col}')
        axes[idx].set_xlabel('Generation (TWh)')
        axes[idx].tick_params(axis='x', rotation=45)

# Remove any empty subplots
for idx in range(len(renewable_cols), len(axes)):
    fig.delaxes(axes[idx])

```

```

plt.tight_layout()
plt.show()

# Time series analysis
plt.figure(figsize=(15, 8))
for col in renewable_cols:
    plt.plot(yearly_data['Year'], yearly_data[col], label=col)
plt.title('Renewable Energy Generation Over Time')
plt.xlabel('Year')
plt.ylabel('Generation (TWh)')
plt.legend(bbox_to_anchor=(1.05, 1), loc='upper left')
plt.tight_layout()
plt.show()

# Calculate summary statistics
print("\nSummary Statistics:")
total_generation = yearly_data[renewable_cols].sum().sum()
print(f"Total Generation: {total_generation:.2f} TWh")

latest_year = yearly_data['Year'].max()
print(f"\nLatest Year ({latest_year}) Generation Mix:")
latest_data = yearly_data[yearly_data['Year'] == ↵
↵latest_year][renewable_cols].iloc[0]
for col in renewable_cols:
    percentage = (latest_data[col] / latest_data.sum()) * 100
    print(f"{col}: {latest_data[col]:.2f} TWh ({percentage:.1f}%)")

# Calculate compound annual growth rate (CAGR)
print("\nCompound Annual Growth Rate (CAGR):")
years = latest_year - yearly_data['Year'].min()
for col in renewable_cols:
    initial_value = yearly_data[yearly_data['Year'] == ↵
↵min()][col].iloc[0]
    final_value = latest_data[col]
    if initial_value > 0: # Avoid division by zero
        cagr = (pow(final_value / initial_value, 1 / years) - 1) * 100
        print(f"{col}: {cagr:.1f}%")

# Run the analysis
print("Starting statistical analysis...")
perform_statistical_analysis()

```

```

[ ]: # Summary and Insights
def generate_summary():
    """Generate summary of key findings"""
    summary = ""

```


Key Findings from Data Exploration:

1. Data Quality:

- Minimal missing values in core variables
- No significant data quality issues
- Some outliers present in renewable generation data

2. Temporal Patterns:

- Clear upward trend in renewable energy adoption
- Significant seasonal variations in generation
- Acceleration in growth rates post-2010

3. Geographic Distribution:

- High concentration in developed countries
- Significant regional variations
- Emerging markets showing rapid growth

4. Weather Impact:

- Strong correlation with solar radiation
- Moderate wind speed dependency
- Temperature effects vary by region

5. Energy Mix:

- Increasing share of renewables
- Hydro and wind dominate renewable sources
- Solar showing fastest growth rate

Next Steps:

1. Feature Engineering:

- Create weather-based features
- Calculate growth rates and trends
- Generate regional indicators

2. Preprocessing:

- Handle outliers in generation data
- Normalize weather variables
- Create consistent time series format

"""

```
display(HTML(f"<pre>{summary}</pre>"))
```

```
generate_summary()
```