

WEATHER DATA ANALYSIS

JEEVA JOSE

ARAVIND KRISHNA A V

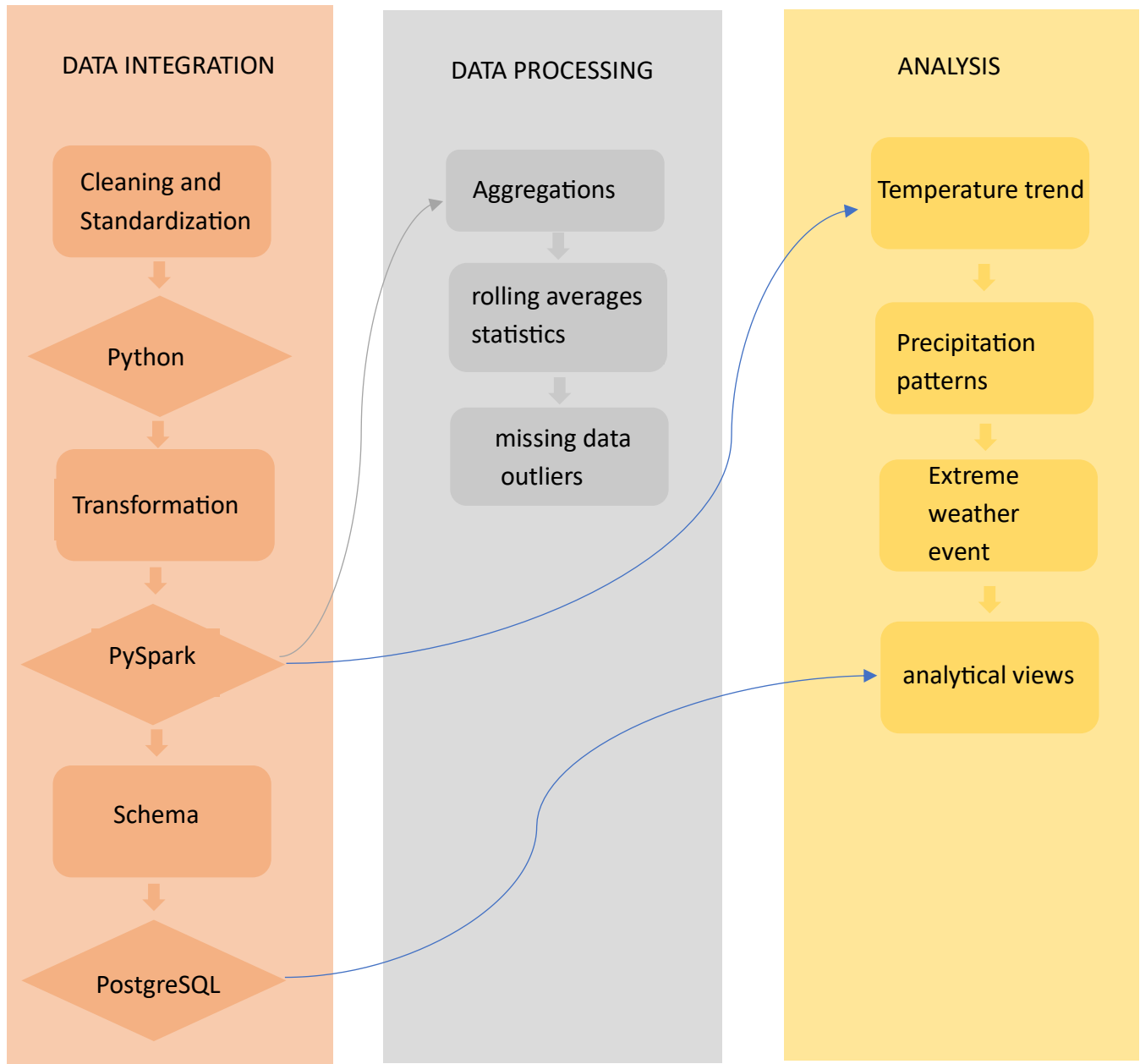
1.INTRODUCTION

Weather data analysis is essential for understanding climate patterns and environmental conditions that impact daily life. With the rise of large-scale datasets, efficient methods for cleaning, processing, and analyzing this data are crucial. This project leverages advanced tools like Python, pandas, PySpark, and PostgreSQL to uncover insights into temperature trends, precipitation patterns, and extreme weather events. These analyses benefit fields such as meteorology, agriculture, urban planning, and disaster management, providing valuable information for decision-making and preparedness.

2.OBJECTIVES

1. **Integrate and Standardize Data:** Combine weather data from various sources and ensure consistency.
2. **Large-Scale Data Processing:** Utilize PySpark for efficient data transformation and PostgreSQL for storage.
3. **Aggregate and Analyze Data:** Calculate daily/monthly patterns, rolling averages, and handle outliers.
4. **Trend and Pattern Analysis:** Identify temperature trends, precipitation patterns, and seasonal variations.
5. **Detect Extreme Weather Events:** Identify and analyze the frequency and impact of extreme weather events.
6. **Visualize and Report:** Create graphs, charts, dashboards, and analytical views in PostgreSQL.
7. **Document Methodologies:** Provide comprehensive documentation for data integration, processing, and analysis.

3. TECHNICAL ARCHITECTURE



4.DATA INTEGRATION

4.1 DATA SOURCES

The Global Weather Repository is a comprehensive dataset that provides daily updated weather information from various locations worldwide. It includes data on temperature, precipitation, humidity, wind speed, and other relevant weather parameters.

4.2 DATA CLEANING

The data cleaning process involves handling missing values, checking for duplicates, and standardizing column names to ensure the dataset is accurate and consistent. Below are the steps taken to clean the weather data:

1. **Handling Missing Values:**

- The dataset is checked for missing values in each column to identify any gaps in the data.
- Used pandas methods for this..

2. **Checking for Duplicates:**

- The dataset is examined for duplicate rows to maintain data integrity.
- Any duplicate records are identified and removed to ensure uniqueness.

3. **Cleaning and Standardizing Column Names:**

- Column names are cleaned and standardized by stripping whitespace, converting to lowercase, and replacing spaces with underscores.
- This ensures consistency and readability in the dataset.

4. Converting Date Columns:

- Date columns, such as 'last_updated', are converted to the appropriate datetime format.
- This standardization allows for accurate time-based analysis and calculations.

4.3 DATA TRANSFORMATION

To handle large-scale weather data efficiently, we use PySpark for data transformation. Below are the steps involved in transforming the data:

1. Extract Date Information:

- Create a new column date_only by extracting the date from the last_updated column.
- Extract the year, month, and day from the date_only column and store them in separate columns.

2. Transformations:

- Utilize PySpark's withColumn method to perform the transformations.

4.4 DATA STORAGE

This PostgreSQL schema efficiently stores and analyzes weather data through seven interconnected tables.

Locations: holds geographic details

weather_observations :records real-time conditions.

weather_conditions: Additional data like(sun, moon, and condition text)

air_quality :(pollution levels) enhance insights

daily_weather: Aggregated data is stored

monthly_weather: for trend analysis.

extreme_weather_events: logs significant weather occurrences

5.DATA PROCESSING

5.1 DATA AGGREGATION

The data aggregation process involves grouping weather data by specific time periods and calculating statistical measures.

Daily Aggregation:

- Group data by year, month, and day.
- Calculate average temperature, total precipitation, average humidity, and average pressure for each day.

Monthly Aggregation:

- Group data by year and month.
- Calculate average temperature, total precipitation, average humidity, and average pressure for each month

5.2 ROLLING AVERAGES AND SUMS

- **Temperature:** Calculate a 3-day rolling average for temperature in Fahrenheit.
- **Precipitation:** Calculate a 3-day rolling average and rolling sum for precipitation in inches.
- **Humidity:** Calculate a 3-day rolling average for humidity in percentage.

5.3 STATISTICAL MEASURES

- Calculate the following statistics for temperature, precipitation, and humidity:
 - Mean, Sum, Count, Minimum, and Maximum values.

5.4 HANDLING OUTLIERS

The process of detecting and handling outliers is crucial for maintaining data quality and accuracy. Below are the key steps taken to detect and replace outliers in the weather data:

1. Detecting Outliers Using Z-Score:

- Calculate the mean and standard deviation of the target column.
- Compute the Z-score for each data point in the column.
- Identify outliers as data points with a Z-score greater than 3 (indicating they are more than three standard deviations away from the mean).

2. Replacing Outliers with Median Value:

- Calculate the median value of the target column.
- Replace the outliers with the median value to minimize their impact on the data.

3. Outlier Detection and Replacement:

- Apply the above steps to detect and replace outliers in the `temperature_fahrenheit` and `precip_in` columns.
- Ensure the cleaned data is saved for further analysis.

6.VISUALIZATION

6.1 YEARLY AND MONTHLY ANALYSIS

Yearly Average Temperature and Precipitation:

- Grouped data by year to calculate and visualize the average temperature and precipitation trends.

Monthly Average Temperature and Precipitation:

- Grouped data by year and month to calculate and visualize the average temperature and precipitation trends

6.2 EXTREME WEATHER EVENTS REPORT

Detection

1. **Data Preparation:** Ensure date format.
2. **Thresholds:** Temperature > 100°F, Precipitation > 1 inch.
3. **Identify Events:** Detect and combine extreme temperature and precipitation events. Add event type.

6.3 ANALYTICAL VIEWS

Avg_monthly_weather:

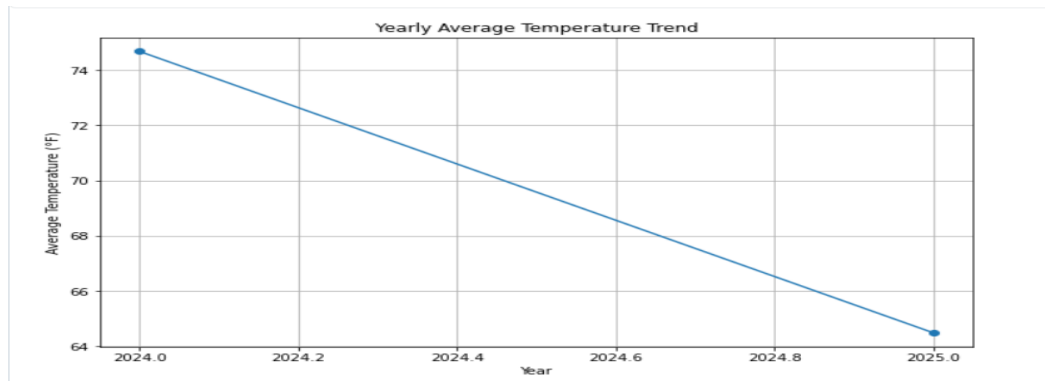
- Aggregates monthly weather data.
- Displays average temperature, precipitation, humidity, and wind speed for each location.

Temperature_trends:

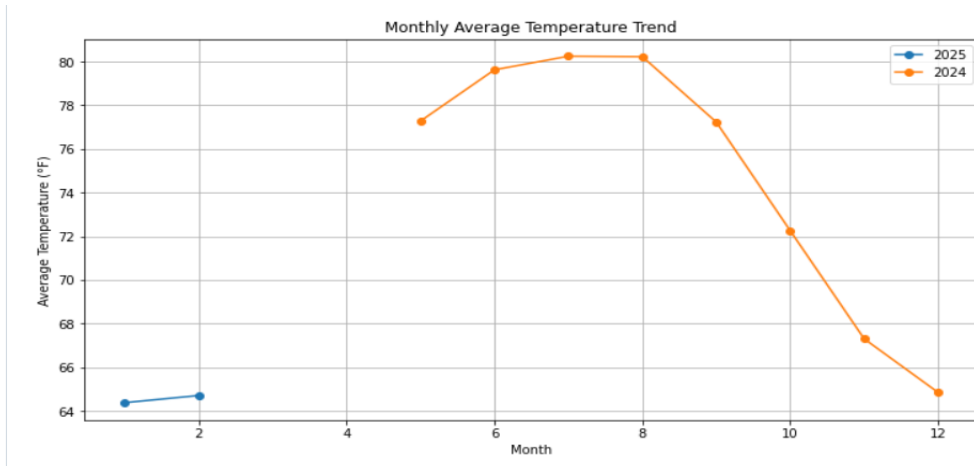
- Calculates a 7-day rolling average of temperatures.
- Helps analyze temperature fluctuations over time.

7.RESULTS

Temperature trend analysis

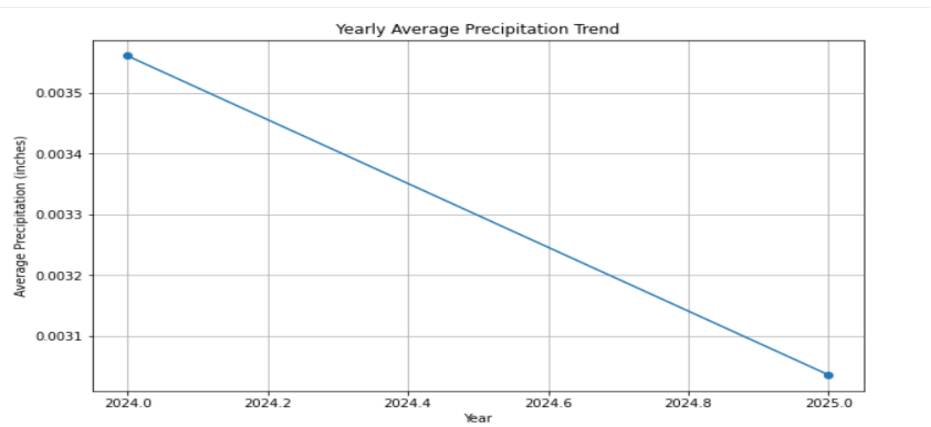


The line graph titled "**Yearly Average Temperature Trend**" shows a decrease in average temperature from 74°F in 2024 to 64°F in 2025. This indicates a downward trend, which could reflect changing climate patterns.

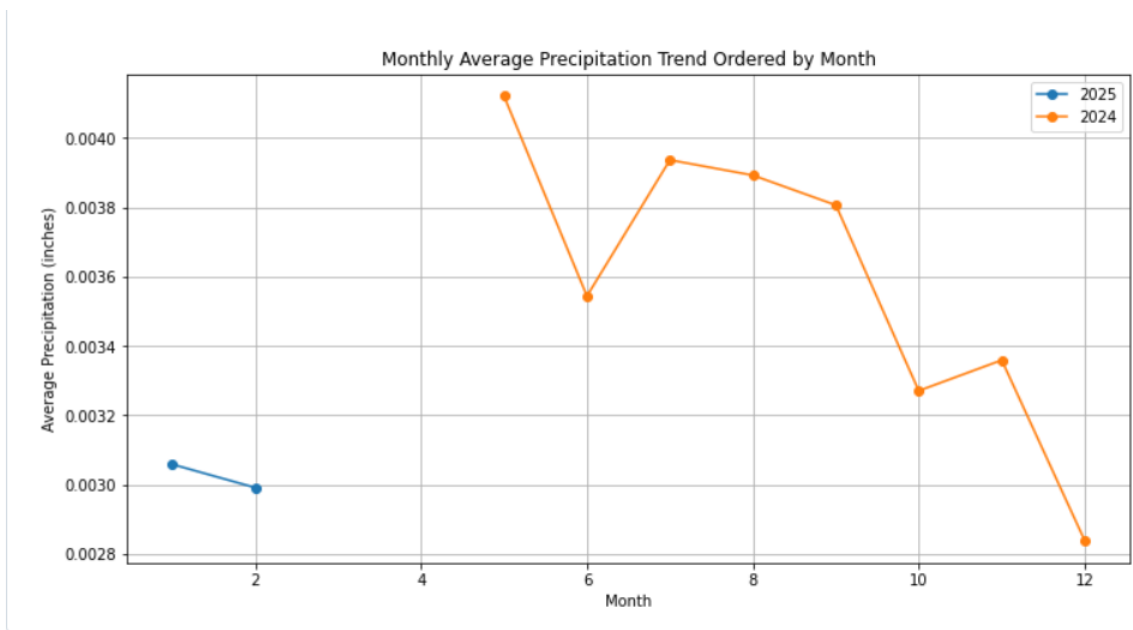


It is a graph that compares **monthly average temperatures** for the years 2024 and 2025. The graph highlights how temperatures fluctuate seasonally, with peaks around July and dips in the winter months

Precipitation patterns

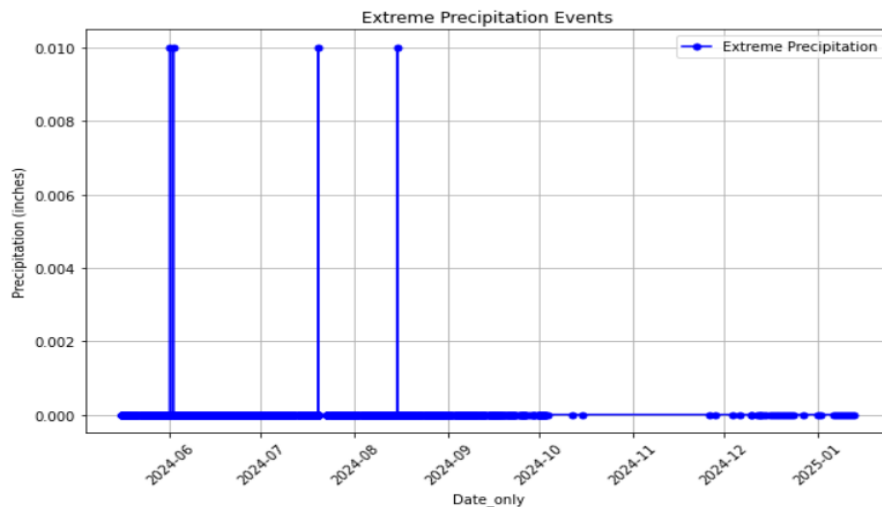


The line graph titled "**Yearly Average Precipitation Trend**" shows a significant decrease in average precipitation from approximately 0.00355 inches in 2024 to about 0.0031 inches in 2025. This downward trend could suggest a shift in weather patterns or a drier overall climate between these years.

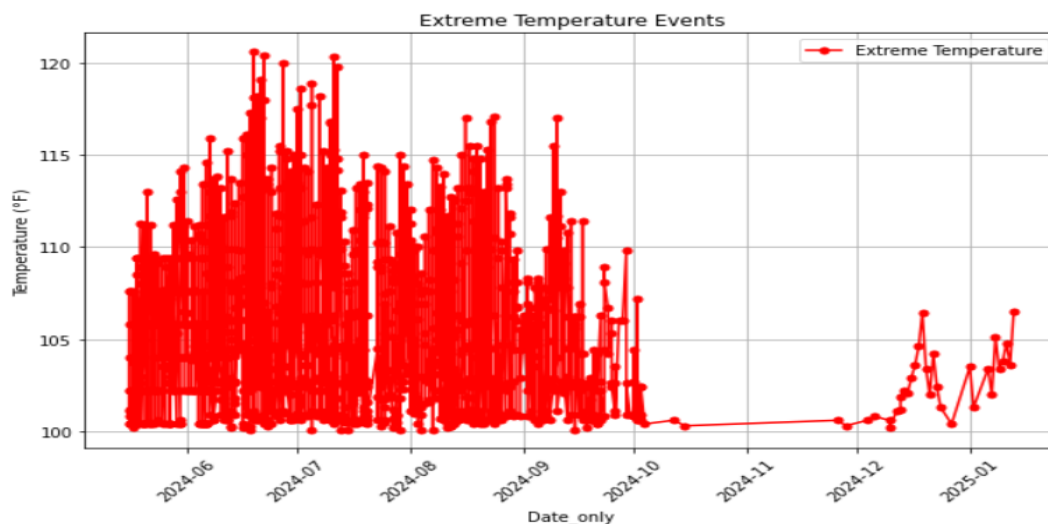


The graph compares **monthly precipitation trends** for 2024 and early 2025, showing fluctuations and suggesting incomplete data for 2025. It highlights seasonal changes and differences between the two years.

Extreme weather event detection



The graph "**Extreme Precipitation Events**" tracks precipitation from June 2024 to January 2025, showing notable spikes in rainfall around mid-2024 and smaller, frequent events later. The blue lines and dots highlight the intensity and timing of each occurrence.



The graph "**Extreme Temperature Events**" depicts intense heatwaves from June to September 2024, with frequent spikes above 110°F. Extreme events become rarer and less intense after September, apart from slight increases in December 2024 and January 2025.

8.CONCLUSION

The Weather Data Analysis project has successfully implemented a comprehensive data pipeline to collect, clean, standardize, and analyze weather data. Key achievements include:

Data Integration and Cleaning: Weather data was efficiently integrated from multiple sources, cleaned, and standardized to ensure consistency and accuracy.

Data Processing and Aggregation: The data was aggregated on a daily and monthly basis, enabling the calculation of relevant statistical measures such as rolling averages and sums.

Outlier Detection and Handling: Outliers were detected and replaced with median values to maintain data integrity.

Trend and Pattern Analysis: Temperature and precipitation trends were analyzed on both yearly and monthly scales, providing valuable insights into weather patterns.

Extreme Weather Events Detection: Extreme temperature and precipitation events were identified and visualized, highlighting patterns and potential impacts.

-Visualization and Reporting: Clear and informative visualizations were created to effectively communicate findings and insights.

9.FUTURE SCOPE

The project scope includes collecting, cleaning, and standardizing weather data, aggregating data, calculating statistics, handling outliers, analyzing trends and extreme weather events, designing PostgreSQL schemas, and creating visualizations with documentation. Real-time weather forecasting, detailed climate change analysis, and developing new data collection tools are out of scope, ensuring focused and efficient execution.