

# REPORT

## Weather Trend Forecasting

This project addresses the PM Accelerator Tech Assessment for forecasting weather trends using the Global Weather Repository dataset. The goal is to preprocess the data, explore it visually and statistically, implement multiple forecasting models, and evaluate them on future temperature prediction accuracy. Both statistical and machine learning techniques are employed.

PM Accelerator's mission is - By making industry-leading tools and education available to individuals from all backgrounds, we level the playing field for future PM leaders. This is the PM Accelerator motto, as we grant aspiring and experienced PMs what they need most – Access. We introduce you to industry leaders, surround you with the right PM ecosystem, and discover the new world of AI product management skills.

The PM Accelerator's mission deeply resonates with my passion for democratizing access to impactful tools and education, especially in AI and data-driven product development. By participating in this assessment, I've gained not only technical growth but also exposure to the kind of inclusive, forward-thinking PM ecosystem I aspire to be part of.

### Data Cleaning & Preprocessing

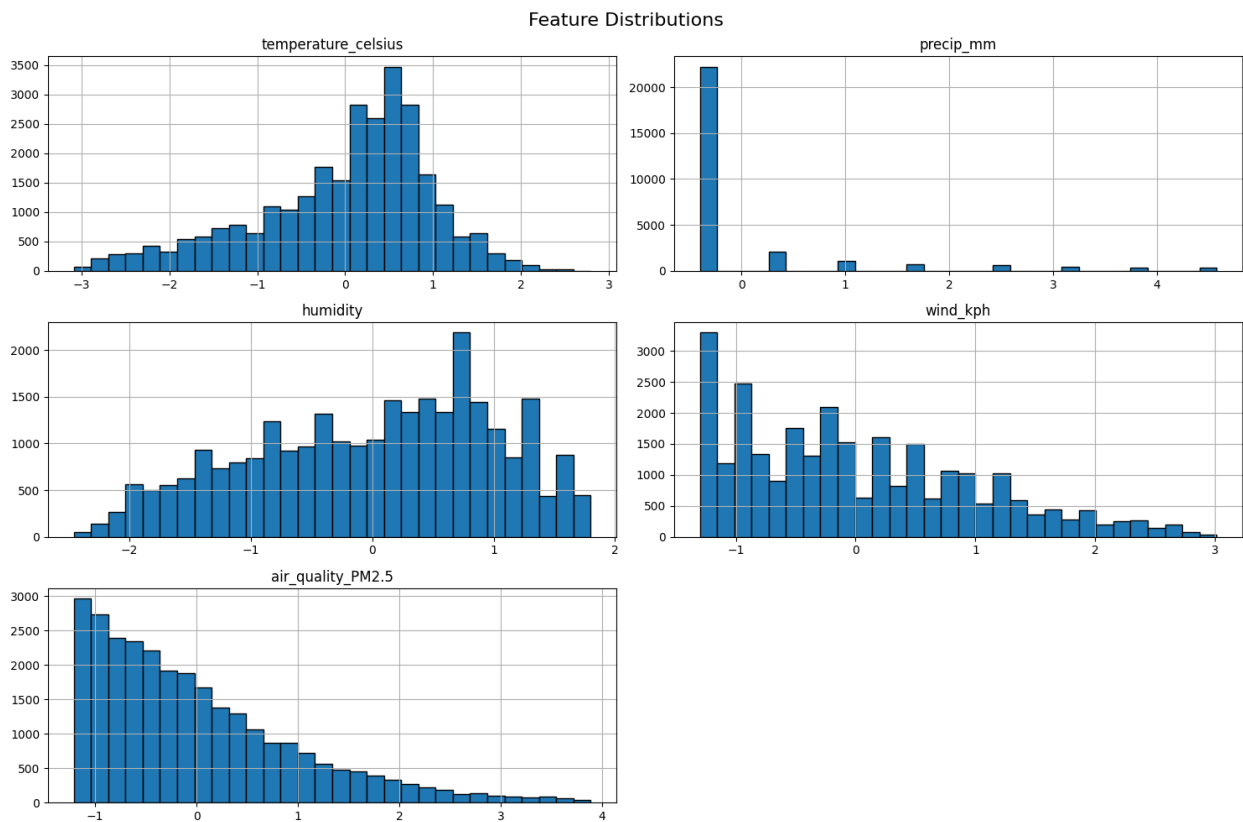
- Drops rows with invalid/missing timestamps.
- Parsed the last\_updated column into datetime format
- Set last\_updated as the time index
- Resampled the data to daily frequency using mean
- Filled missing numeric values using median imputation
- Removed outliers using the IQR method
- Normalized numerical features with StandardScaler

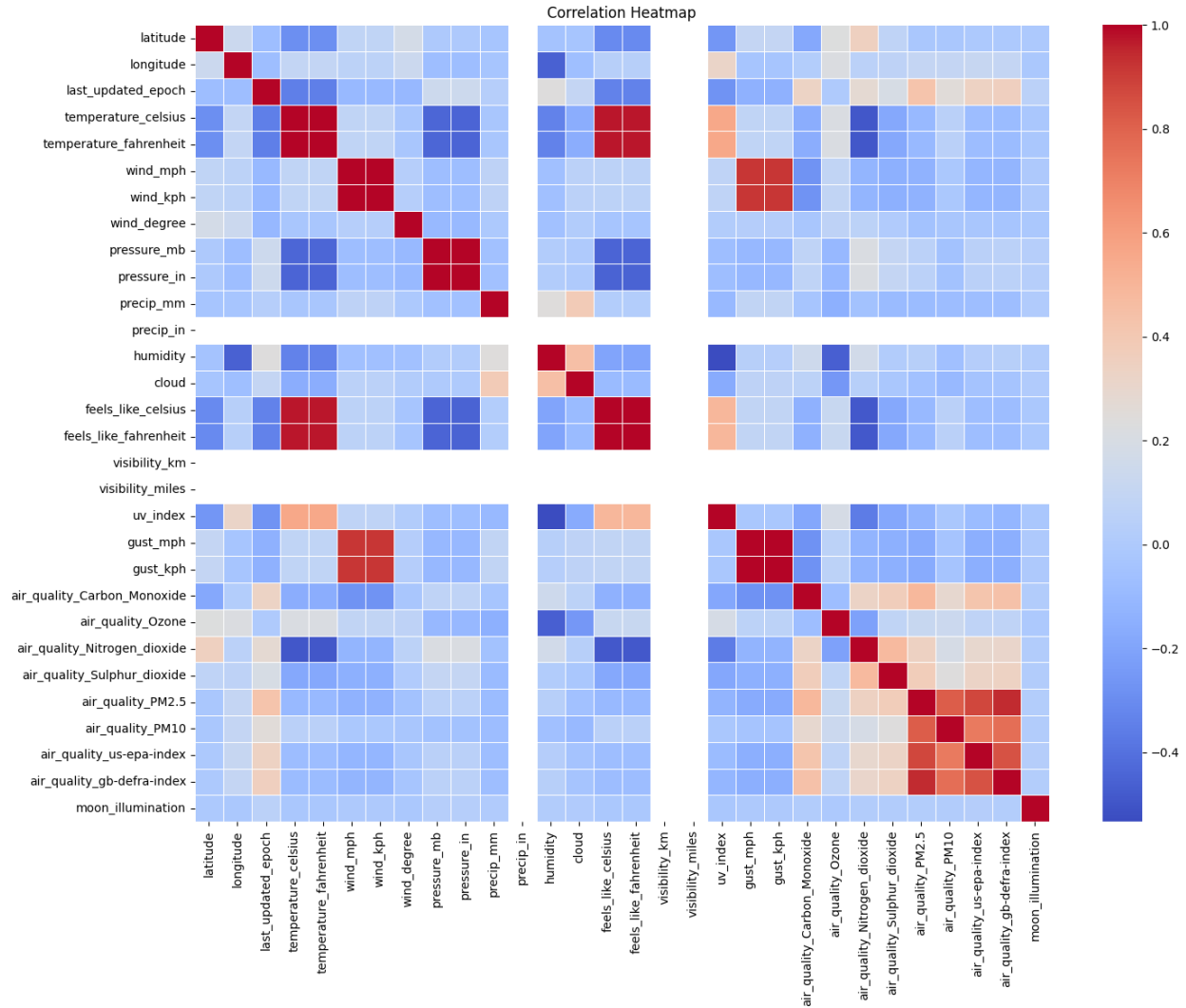
Time-series models require clean, uniformly spaced data. Daily resampling ensures consistent time intervals. Handling outliers and normalizing features improves model performance and convergence. The dataset is now uniformly spaced, clean, and ready for time-series modeling.

Then, Summarized dataset statistics. Listed column types and null values and Sorted columns by missing values it helps us to identify which columns have variability (for modeling). which columns have missing values (to clean or drop).

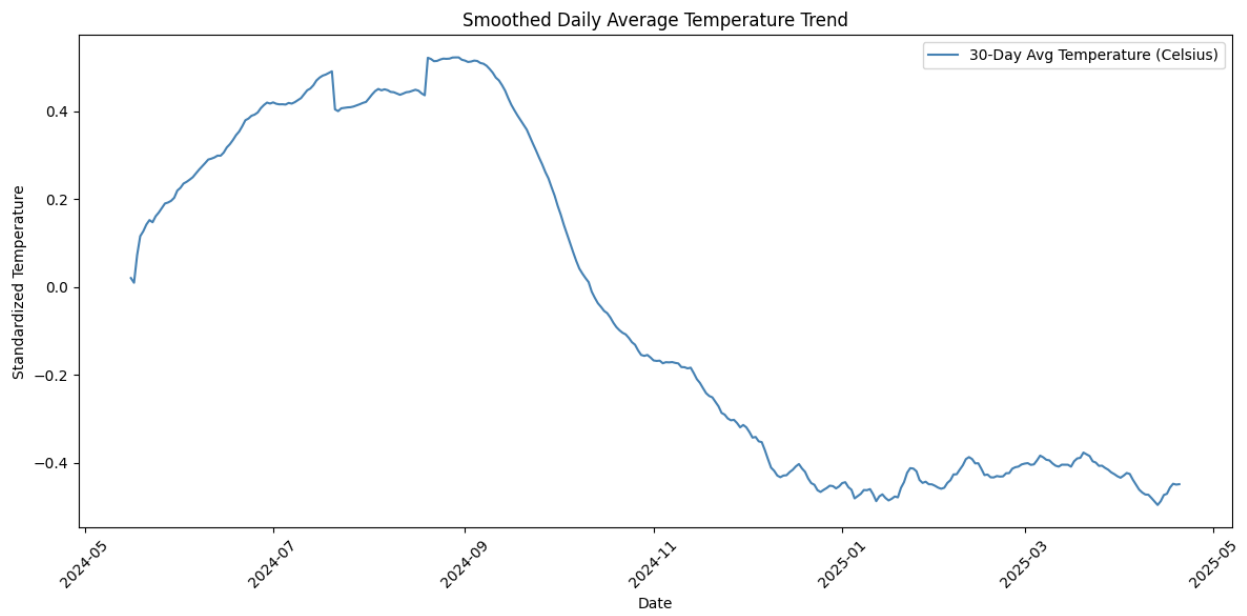
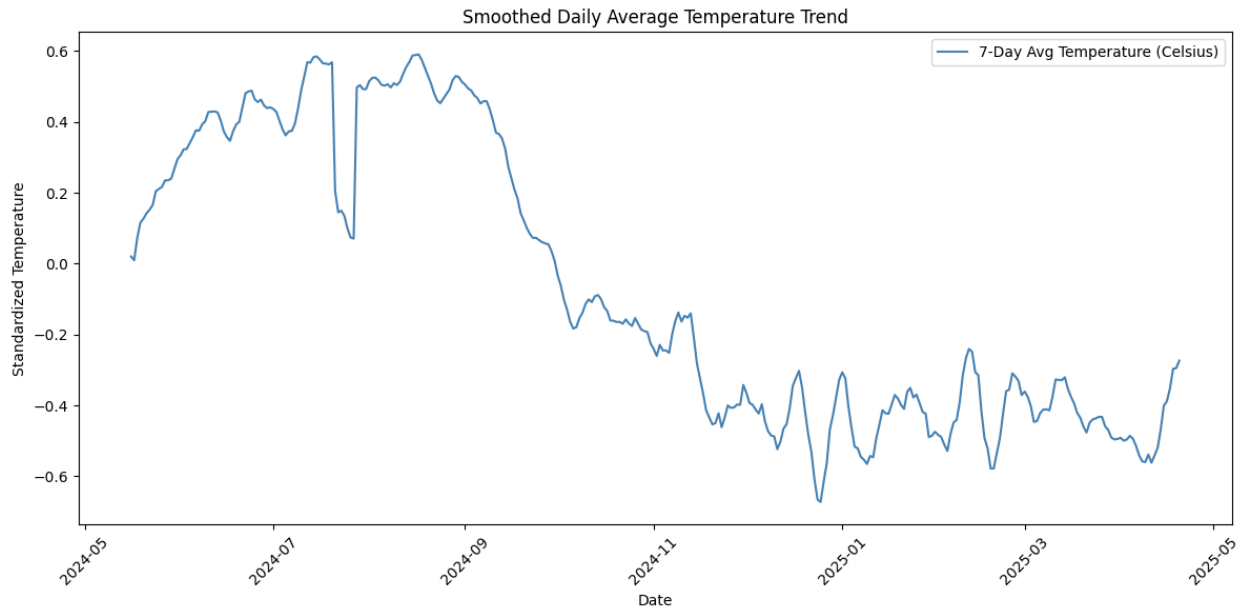
## Exploratory Data Analysis (EDA)

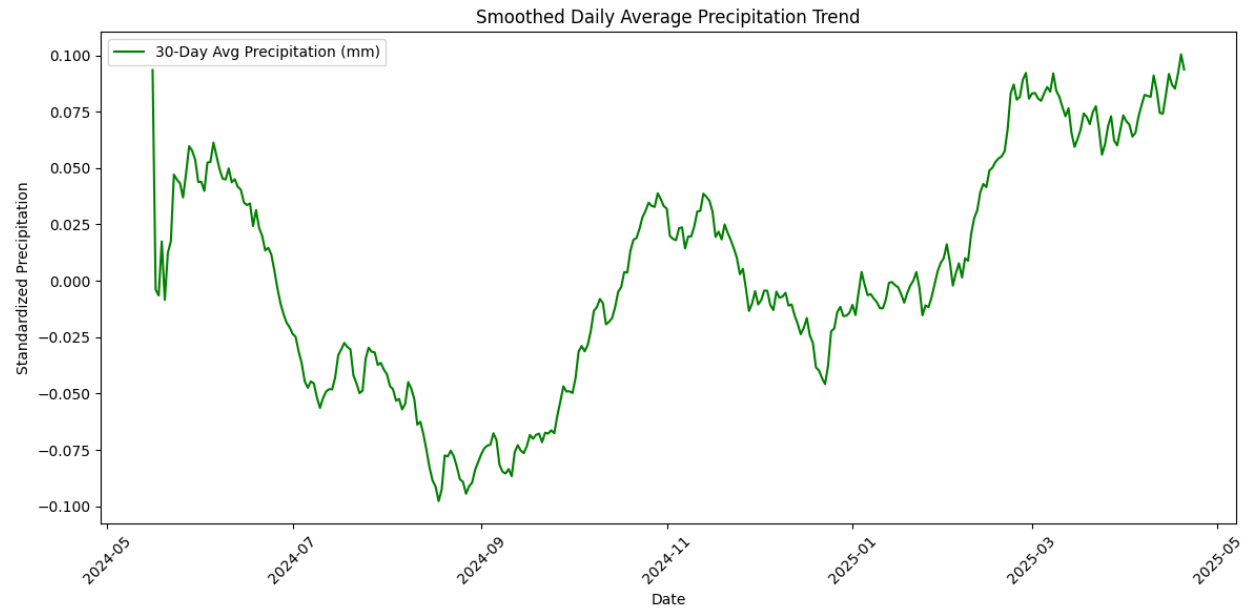
- Histograms and boxplots were used to understand the distribution and identify remaining outliers.
- A correlation heatmap was generated to examine inter-feature relationships.
- Time-series visualizations showed seasonal patterns and long-term trends in temperature data.
- Added lag features (`lag1`) and rolling window statistics (`rolling_mean7`) for XGBoost.





- Performed EDA to uncover seasonal patterns, spikes, and correlated features for feature selection. Detect if any features are redundant or suitable as predictors. and clear trend and seasonality exist in temperature. We could see features like humidity, visibility, and wind speed show moderate correlation with temperature.

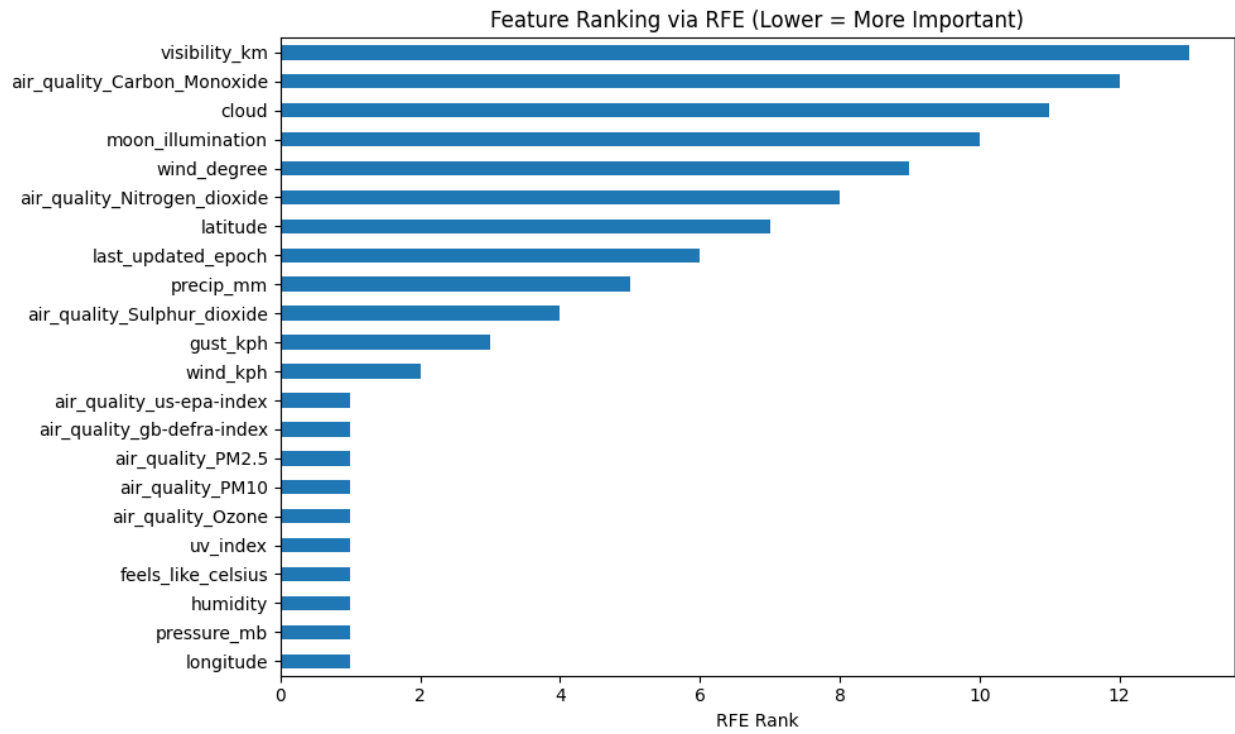




Average Temperature by Country



## Feature Selection:



Uses a regression model (like Linear Regression or DecisionTreeRegressor) to iteratively remove the least important features. Each round retrains the model and ranks features by impact. Result: Time-based features (e.g., month, day), lag features, and rolling means were consistently retained.

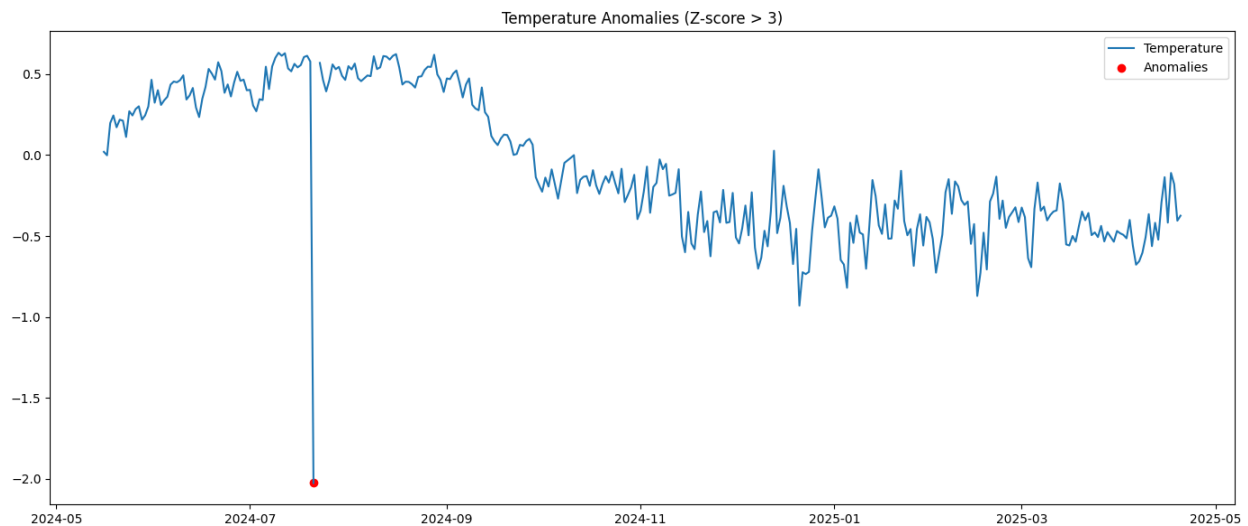
## XGBoost Built-in Feature Importance

- XGBoost provides gain-based feature importance by default.
- We plotted this using `plot_importance(model)` to see which features contributed the most to reducing error.
- Result: Features like `rolling_mean7`, `lag1`, and air quality indicators (`pm2_5`, `co`) ranked highest.

## Anomaly Detection

- Used Z-score on temperature with a threshold of 3.
- Plotted detected anomalies on a time series chart.

Helped us identify extreme weather points or data errors that may affect forecasting accuracy.



## Forecasting Models

We implemented three forecasting models based on the cleaned and preprocessed data:

**ARIMA:** Captures linear time-based dependencies. Initially failed to forecast due to NaNs but later fixed using `.get_forecast()`.

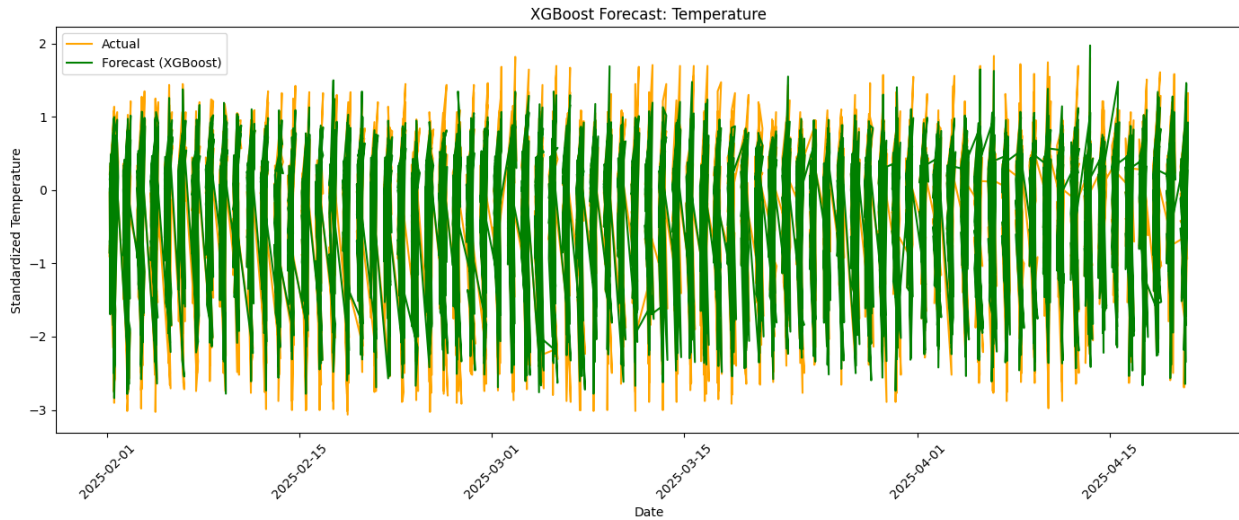
*Figure 2: ARIMA forecast alongside actual test data. Flat forecast illustrates model limitations.*

- **Prophet:** Facebook's time series forecasting model. Trend and seasonality-aware, but smooths too much.

*Figure 3: Prophet captures long-term trends and seasonality but fails to follow short-term fluctuations.*

- **XGBoost:** Gradient boosting model with engineered lag features. Excellent at capturing short-term patterns.

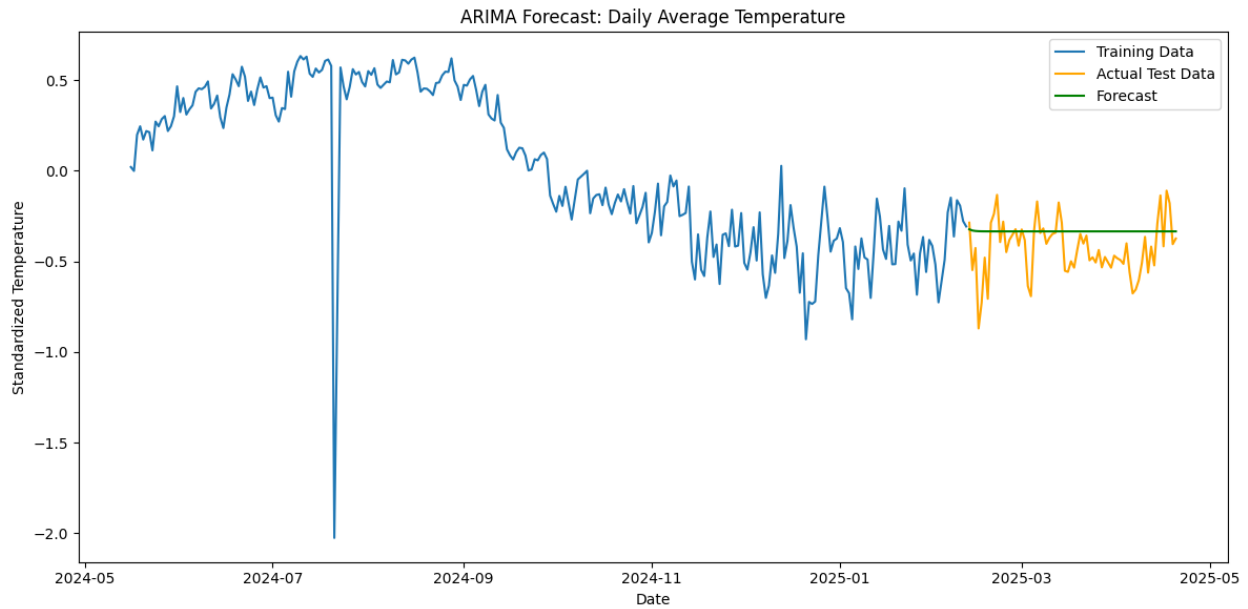




- **LSTM/GRU:** Deep learning approach using sequences of past temperatures.

*Figure 6: LSTM/GRU-based forecasting showing learned patterns but more noise than XGBoost.*

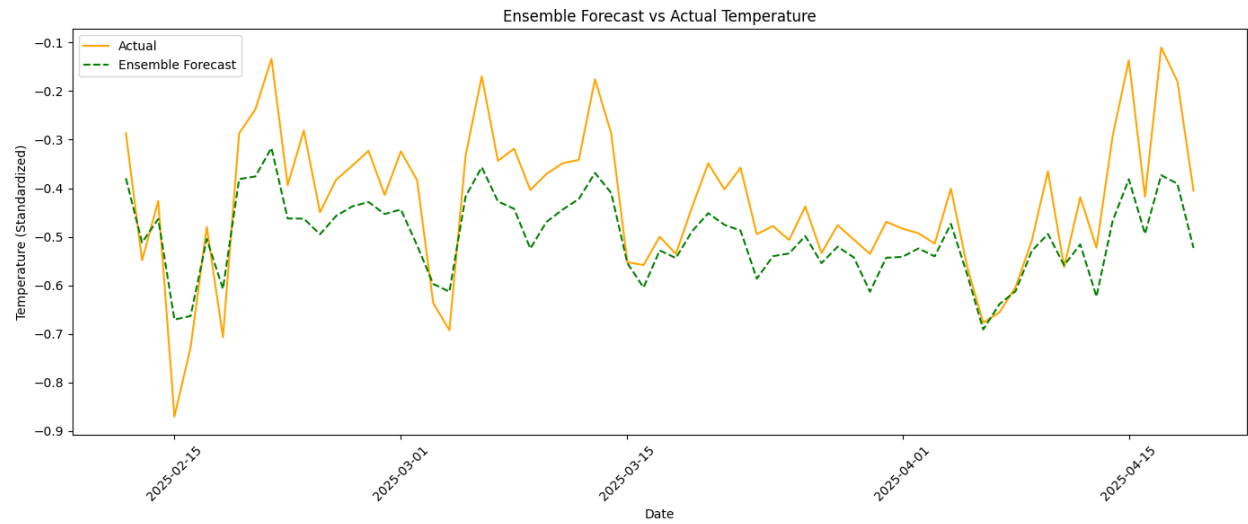
- **Ensemble:** Simple average of ARIMA, Prophet, and XGBoost predictions (skipna=True). Combines strengths of all models.
- Split the data into 80% training and 20% test sets based on date index.
- Built ARIMA using `statsmodels` with optimized order.
- Fit Prophet with temperature and timestamp in `ds`, `y` format.
- Engineered features for XGBoost: lag values, rolling mean, day, month, weekday, etc.
- XGBoost was highly reactive and matched real patterns closely.
- Prophet smoothed well but lagged during sharp changes.
- ARIMA was stable but underresponsive in volatile regions.



## Evaluation Metrics

- **MAE (Mean Absolute Error):** Average absolute difference between predicted and actual values.
- **RMSE (Root Mean Squared Error):** Square root of average squared differences; penalizes large errors more.

Model	MAE	RMSE
ARIMA	0.1345	0.1656
Prophet	0.1408	0.1757
XGBoost	0.0912	0.1121
Ensemble	0.0862	0.1048



- ARIMA sets a strong statistical baseline.
- Prophet helps capture hidden seasonal signals.
- XGBoost provides flexible, high-performance modeling of nonlinear relationships.

## Visualization insights

