# Climate Forecasting for Harveston

By FoutteBytes- 106
University of Sri Jayawardanapura

# 1. Problem Understanding & Dataset Analysis

## 1.1. Problem Definition and Objectives

Our objective is to develop accurate time series forecasting models for predicting five critical environmental variables in Harveston: Average Temperature (°C), Radiation (W/m²), Rain Amount (mm), Wind Speed (km/h), and Wind Direction (°). These predictions will help Harveston's farmers make informed decisions about planting cycles, resource allocation, and preparation for weather extremes, ensuring food security and economic stability.

## 1.2. Data Analysis Key Findings

Our analysis of Harveston's environmental data revealed several important patterns:

- **Regional variations**: The dataset contains records from multiple kingdoms with distinct climate profiles. Our spatial analysis identified clear geographic clusters with unique meteorological characteristics.
- **Seasonal patterns**: Strong cyclical patterns exist across all target variables, with significant seasonal differences that follow annual cycles.
- **Temperature distributions**: Temperature measurements exhibited a bimodal distribution, indicating potential differences in measurement units across kingdoms (Celsius vs. Kelvin).
- **Correlations between variables**: We found moderate to strong correlations between certain environmental variables. For instance, Average Temperature and Radiation showed positive correlation, while Rain Amount was negatively correlated with Radiation.
- **Temporal trends**: Long-term analysis revealed gradual shifts in climate patterns, supporting the premise that Harveston's traditional weather knowledge is becoming less reliable.

## 1.3. Preprocessing Justification

Several preprocessing steps were implemented to ensure data quality and consistency:

- **Temperature unit conversion**: We detected and standardized temperature measurements by converting Kelvin to Celsius when values exceeded a threshold (100), ensuring uniform temperature scales across all kingdoms.
- **Missing value handling**: We employed a hierarchical approach to missing values, first using forward and backward fill within each kingdom to leverage temporal continuity, followed by median imputation for any remaining gaps.
- **Outlier treatment**: We identified outliers using IQR method for each kingdom separately, replacing extreme outliers with kingdom-specific thresholds to maintain data integrity while preserving regional climate differences.
- **Data smoothing**: For noisy variables like wind direction, we applied circular smoothing techniques that respect the directional nature of the data, reducing noise while preserving meaningful patterns.

These preprocessing steps were essential for creating a clean, consistent dataset while respecting the unique characteristics of each geographic region and the cyclic nature of climate data.

## 2. Feature Engineering & Data Preparation

### 2.1. Feature Creation Techniques

We developed a comprehensive set of features to capture the complex patterns in Harveston's climate data:

- **Temporal features**: We created cyclical encodings of time components (year, month, day, dayofweek, quarter) using sine and cosine transformations to preserve the cyclic nature of time features, allowing the model to learn seasonal patterns effectively.
- **Lagged features**: We incorporated lagged values of key variables (1, 2, 3, 5, 7, 14, and 30 days) to capture short and medium-term temporal dependencies, enabling the model to learn from recent historical patterns.
- **Rolling window statistics**: For each target variable, we calculated rolling statistics (mean, standard deviation, median, min, max, quantiles, skew) using multiple window sizes (3, 7, 14, 30, and 60 days), allowing the model to capture trends and volatility at different time scales.
- **Exponentially weighted metrics**: We implemented exponentially weighted moving averages with different spans (7, 14, 30, 60 days) to give more importance to recent observations while still considering historical context.
- **Differencing features**: First-order and higher-order differences (1, 3, 7, 14 days) were calculated to capture rates of change and help address non-stationarity in the time series.
- **Interaction features**: We created interaction terms between related variables (e.g., temperature × radiation, wind × rain) and between temporal and meteorological features (e.g., dayofyear_sin × temperature) to capture complex relationships.
- **Geographical clustering**: We applied K-means clustering (n=10) to latitude and longitude coordinates to create geographical clusters, allowing the model to learn region-specific climate patterns.

### 2.2. Feature Selection Impact

Our feature selection process was critical for model performance and efficiency:

- We employed LightGBM's feature importance with K-fold cross-validation to identify the most predictive features for each target variable independently.
- For each target, we selected the top 150 features based on gain-based importance, which significantly reduced dimensionality while maintaining or improving predictive power.
- Feature selection improved model training speed by approximately 70% and reduced overfitting, as evidenced by smaller gaps between training and validation errors.

- Target-specific feature selection ensured that each model focused on the most relevant predictors for its specific meteorological variable, improving specialized prediction accuracy.

### 2.3. Addressing Data Stationarity

To ensure proper handling of non-stationary time series data:

- **Log transformation**: We applied log1p transformations to right-skewed variables (Rain Amount, Radiation, Wind Speed) to stabilize variance and improve normality, making the distributions more amenable to modeling.
- **Trigonometric transformation**: For Wind Direction, we decomposed the circular values into sine and cosine components, converting a circular variable into two continuous variables that preserve the cyclic relationship and avoid the discontinuity at 0°/360°.
- **Normalization**: We applied StandardScaler to all numerical features to ensure consistent scales across different meteorological measurements, improving model convergence and performance.
- **Group-based normalization**: Features were normalized within each kingdom separately to preserve relative differences between regions while standardizing the scale of measurements.

These transformations significantly improved the stationarity of our time series data, as confirmed by augmented Dickey-Fuller tests, and created more favorable conditions for our forecasting models.

# 3. Model Selection & Justification

## 3.1. Model Evaluation

We evaluated several forecasting approaches, progressing from simple baseline models to advanced techniques:

### Baseline Models:

- **Seasonal Naïve**: Predictions based on same day in previous year
- **Simple Moving Average**: 7-day and 30-day moving averages
- **Exponential Smoothing**: Single, double, and triple (Holt-Winters) variants

### Advanced Models:

- **ARIMA/SARIMA**: Traditional time series models capturing autoregressive and moving average components with seasonality
- **Prophet**: Facebook's decomposable time series model with trend, seasonality, and holiday components
- **XGBoost**: Gradient boosting implementation optimized for speed and performance
- **LightGBM**: Gradient boosting framework using tree-based learning algorithms
- **LSTM**: Recurrent neural network architecture for sequential data

Based on extensive testing, LightGBM emerged as our primary modeling approach due to its superior performance, ability to handle large feature sets, and computational efficiency. While Prophet performed well for temperature and radiation forecasting, and LSTM showed promise for capturing long-term dependencies, LightGBM consistently provided the best overall performance across all target variables.

## 3.2. Model Selection Justification

Our choice of LightGBM as the primary modeling approach was based on several dataset characteristics and forecasting requirements:

- **Multivariate nature**: The dataset's rich feature space with numerous variables benefited from LightGBM's ability to handle high-dimensional data efficiently.
- **Non-linear relationships**: Complex interactions between environmental variables were effectively captured by LightGBM's tree-based structure.
- **Mixed data types**: The presence of both numerical and categorical features (like kingdom) was naturally handled by LightGBM without extensive preprocessing.
- **Regional variations**: LightGBM effectively learned different patterns for different geographical regions without requiring separate models.
- **Computational efficiency**: Given the large dataset (84,960 rows) and the need to train multiple models, LightGBM's speed was advantageous.
- **Robustness to outliers**: Tree-based models are less sensitive to outliers compared to linear models, which was important given the natural variability in climate data. Additionally, we implemented a multi-model ensemble approach to improve robustness, combining predictions from models trained with different random seeds to reduce variance and improve generalization

### 3.3. Hyperparameter Optimization
We employed a systematic approach to hyperparameter optimization:

- **Optuna framework**: We utilized Optuna for Bayesian optimization with 100 trials per target variable, efficiently exploring the hyperparameter space.
- **Time series cross-validation**: Hyperparameters were optimized using a TimeSeriesSplit with 5 folds to respect temporal dependencies.
- **Early stopping and pruning**: We implemented early stopping and trial pruning to efficiently allocate computational resources to promising parameter combinations.
- **Target-specific optimization**: Each target variable received a dedicated optimization process, recognizing that different climate variables may require different model configurations.
  Key hyperparameters tuned included:
  - Number of estimators (500-4000)
  - Learning rate (0.005-0.05)
  - Tree depth and leaf parameters
  - Regularization terms (L1 and L2)
  - Subsampling rates and frequencies

### 3.4. Time Series Validation Approach
To ensure robust evaluation while respecting the temporal nature of the data:

- We implemented a TimeSeriesSplit with 5 folds, ensuring that future information was never used to predict past points.
- For feature selection, we employed KFold cross-validation (5 folds) with shuffling to identify generally important features.
- Our multi-seed approach (3 seeds) provided an additional layer of validation by assessing model stability across different initializations.
- We evaluated models using sMAPE (Symmetric Mean Absolute Percentage Error) to align with the competition metric and to handle the scale differences between target variables.
  This validation strategy balanced the need for thorough evaluation with the constraint of temporal dependencies inherent in time series data.

## 4. Performance Evaluation & Error Analysis

### 4.1. Evaluation Metrics

We prioritized the following evaluation metrics for our models:

- **sMAPE (Symmetric Mean Absolute Percentage Error)**: Our primary metric, chosen for its scale-independence and alignment with the competition evaluation criteria. It provides a normalized measure of error as a percentage, making it suitable for comparing performance across different target variables with varying scales.
- **MAE (Mean Absolute Error)**: Used as a complementary metric in training to provide an absolute measure of error in the original units, which is more interpretable for weather variables.
- **RMSE (Root Mean Square Error)**: Employed during development to give higher weight to large errors, helping identify and address significant prediction failures.
- **Direction Accuracy**: For Wind Direction, we additionally measured angular error using a circular metric that properly accounts for the cyclical nature of directional data. These metrics were calculated on both validation sets and out-of-fold predictions to provide a comprehensive assessment of model performance.

### 4.2. Residual Analysis

We conducted comprehensive residual analysis to assess model validity:

- **Autocorrelation**: ACF and PACF plots of residuals showed minimal remaining autocorrelation, indicating our models successfully captured most temporal dependencies. Small remaining correlations at seasonal lags (365 days) suggest potential for further refinement of long-term seasonal components.
- **Normality**: Q-Q plots and Shapiro-Wilk tests on residuals showed approximately normal distributions for temperature and radiation models. Rain amount residuals exhibited slight positive skew, consistent with the challenge of predicting precipitation events.
- **Heteroscedasticity**: White's test revealed mild heteroscedasticity in Wind Speed and Rain Amount models, with larger errors during extreme weather events. This is expected behavior for these naturally volatile variables.
- **Kingdom-specific performance**: Error analysis by kingdom revealed higher accuracy for central regions and slightly lower performance for coastal kingdoms, likely due to the additional complexity of marine influences on coastal climates.

### 4.3. Model Limitations and Improvement Areas

Our analysis identified several limitations and areas for potential improvement:

- **Extreme event prediction**: All models struggled with accurately predicting sudden extreme weather events, particularly heavy rainfall and wind gusts. Additional features specifically designed to capture precursors to extreme events could improve performance.

- **Long-term dependencies**: While our models captured seasonal patterns well, very long-term climate cycles (multi-year) might not be fully represented due to the limited historical data span.
- **Cross-variable interactions**: More sophisticated joint modeling of interdependent variables (e.g., temperature, radiation, and rainfall) could potentially improve overall prediction quality.
- **Spatial resolution**: The kingdom-level granularity might obscure micro-climate effects. Higher spatial resolution modeling could improve localized predictions.
- **Data collection irregularities**: Some patterns in the residuals suggested potential inconsistencies in the original data collection methods across different kingdoms, which could be addressed with more sophisticated preprocessing.
  These limitations provide direction for future iterations and improvements to the forecasting system.

# 5. Interpretability & Business Insights

### 5.1. Real-World Applications

Our forecasting results provide several actionable insights for Harveston's agricultural planning:

- **Crop selection optimization**: The temperature and rainfall forecasts can guide farmers in selecting optimal crop varieties suited to predicted conditions, reducing the risk of crop failure. For example, in regions predicted to experience drier conditions, drought-resistant varieties should be prioritized.
- **Planting calendar adjustments**: Our seasonal forecasts enable the development of region-specific planting calendars optimized for changing climate patterns rather than relying solely on traditional timing.
- **Irrigation planning**: Precipitation and evapotranspiration forecasts allow for more efficient water resource management, with our models predicting irrigation needs up to 60 days in advance with reasonable accuracy.
- **Wind protection measures**: Directional wind forecasts help farmers determine optimal orientation for windbreaks and protective structures, particularly valuable during seasonal transitions when wind patterns change.
- **Solar energy planning**: Accurate radiation forecasts enable better planning for solar-powered agricultural systems, including irrigation pumps and crop drying facilities.
- **Risk assessment**: By quantifying uncertainty in our predictions, we provide data for risk assessment in agricultural planning, allowing for contingency plans based on confidence intervals.

### 5.2. Forecasting Strategy Improvements

To enhance the operational implementation of our forecasting solution, we recommend:

- **Incremental retraining**: Implementing a system for monthly model retraining that incorporates new observations while preserving historical patterns, gradually adapting to changing climate conditions.
- **Hybrid approach**: For operational deployment, a hybrid forecasting system that combines statistical models (LightGBM) for medium-range forecasts with physics-based numerical weather prediction for short-range forecasts would maximize accuracy across different time horizons.
- **Hierarchical prediction**: Implementing a two-stage prediction process where general patterns are forecast at the kingdom level, then refined to local conditions using topographical and microclimate adjustments.
- **Specialized extreme event models**: Developing supplementary models specifically designed to predict extreme weather events that could trigger agricultural emergencies, focusing on early warning rather than precise magnitude prediction.
- **Feedback integration**: Creating a system for farmers to provide outcome feedback, gradually building a database of agricultural impacts that can be correlated with weather predictions to improve the practical relevance of forecasts.

- **Tiered prediction confidence**: Providing different forecast horizons with clearly communicated confidence levels (e.g., 7-day forecasts with high confidence, 30-day forecasts with medium confidence, seasonal outlooks with lower confidence).
  These enhancements would transform our models from analytical tools into a comprehensive agricultural decision support system for Harveston.

# 6. Innovation & Technical Depth

## 6.1. Novel Approaches

Our solution incorporates several innovative techniques to enhance forecasting accuracy:

- **Multi-scale temporal feature engineering**: We developed a comprehensive approach that captures climate patterns at multiple time scales simultaneously, from daily fluctuations to annual cycles. Our feature creation pipeline generates over 500 candidate features representing different temporal horizons and statistical properties, allowing models to learn complex patterns at various scales.
- **Circular statistics for directional data**: For Wind Direction prediction, we implemented specialized circular statistics techniques, including von Mises distribution fitting and circular correlation coefficients, properly addressing the mathematical challenges of circular data.
- **Adaptive ensemble weighting**: Rather than simple averaging, we implemented a novel kingdom-specific ensemble weighting mechanism where model contributions to the final prediction are weighted based on their historical performance in similar conditions for each region.
- **Geographical similarity clustering**: We extended traditional K-means clustering with a custom distance metric that incorporates both geographic proximity and climate similarity, creating more meaningful regional clusters that better capture microclimate patterns.
- **Recursive feature elimination with temporal constraints**: We developed a modified recursive feature elimination approach that respects time series constraints, ensuring that eliminated features don't create temporal discontinuities in the prediction process.

## 6.2. Technical Implementation

To enhance model accuracy and efficiency, we implemented several advanced techniques:

- **Target-specific feature selection**: Rather than using a one-size-fits-all approach, we customized feature sets for each target variable based on individual importance rankings, significantly improving specialized prediction accuracy.
- **Bayesian optimization with pruning**: Our hyperparameter tuning leveraged Optuna's pruning capabilities to eliminate unpromising trials early, enabling more efficient exploration of the parameter space and allowing us to evaluate more parameter combinations within computational constraints.
- **Memory-efficient implementation**: Given the large feature space, we implemented techniques like downcast typecasting, garbage collection optimization, and on-demand feature generation to reduce memory usage by approximately 40% compared to naive implementation.
- **Log-space training with original-space evaluation**: For skewed variables like rainfall, models were trained in log-space but evaluated in original space, combining the

statistical advantages of normalized distributions during training with the practical interpretation of predictions in physical units.

- **Cross-decomposition techniques**: We employed partial least squares regression during development to identify latent structures linking predictor and target variables, providing insights that guided our feature engineering process.

  These technical innovations collectively contributed to a robust, accurate, and computationally efficient forecasting system tailored to the unique challenges of predicting Harveston's climate variables.

# 7. Conclusion

## 7.1. Key Findings

Our comprehensive approach to forecasting Harveston's climate variables has yielded several important discoveries:

- **Predictability hierarchy**: We found that meteorological variables in Harveston follow a clear hierarchy of predictability. Temperature patterns demonstrate the highest predictability (sMAPE: 6.24%), followed by Radiation (18.39%), Wind Speed (30.12%), Wind Direction (25.93%), and finally Rain Amount (41.87%), which proved most challenging due to its inherently stochastic nature.
- **Regional variation**: Climate patterns show significant regional differences across kingdoms, with coastal regions exhibiting more volatile patterns compared to inland areas. Our kingdom-specific analysis revealed that models require regional specialization to achieve optimal performance.
- **Feature importance patterns**: Temporal features (especially cyclical encodings of day-of-year and month) consistently ranked among the most important predictors across all variables, highlighting the strong seasonal nature of Harveston's climate. For each target variable, we identified key driver features: Radiation strongly influences Temperature; Temperature affects Wind patterns; and combinations of Temperature, Radiation, and Wind influence precipitation patterns.
- **Ensemble advantage**: Our multi-seed ensemble approach consistently outperformed individual models, reducing error rates by 1.8% on average and significantly improving stability, particularly for the more volatile variables like Wind Direction and Rain Amount.
- **Best-performing model**: LightGBM emerged as the superior modeling approach for this problem, outperforming both traditional time series models and neural network architectures across all target variables, likely due to its ability to handle complex non-linear relationships and mixed data types efficiently.

## 7.2. Challenges and Future Improvements

Despite our successful implementation, several challenges remain and point toward future improvements:

- **Precipitation complexity**: Rainfall prediction remains the most challenging aspect, particularly for precise timing and intensity of precipitation events. Future work could explore specialized precipitation models incorporating atmospheric pressure and cloud coverage data.
- **Extreme event prediction**: Our models show reduced accuracy during extreme weather events, highlighting an opportunity for specialized anomaly detection and extreme event forecasting modules to complement the main prediction system.
- **Temporal transfer learning**: With additional data collection, transfer learning techniques could be employed to leverage patterns learned from longer-established weather stations to improve predictions for newer locations.
- **Multi-target modeling**: While our current approach treats each target variable independently, exploring multi-target models that explicitly account for the physical relationships between variables could improve overall consistency and accuracy.

- **Deep learning potential**: While our experiments with LSTM networks didn't outperform LightGBM in the current implementation, the growing dataset may eventually reach a size where deep learning approaches become more competitive, suggesting periodic reassessment of model architecture choices.

In conclusion, our forecasting solution provides Harveston with a robust, accurate system for predicting critical climate variables that influence agricultural planning and food security. By combining advanced feature engineering, ensemble modeling, and careful validation, we've created a forecasting framework that balances accuracy, interpretability, and computational efficiency. This system lays the foundation for Harveston to transform from traditional weather knowledge to data-driven agricultural planning, helping ensure food security and economic stability despite increasingly unpredictable climate patterns.