

CHAPTER 3

RESEARCH METHODOLOGY

3.1 Introduction

This chapter will elaborate in detail on the methodology for establishing a prediction model of tropical cyclone landfall points. This research is centered on predicting the cyclone paths in the coastal regions of Malaysia. Given the high complexity and non - linearity inherent in the atmospheric system, traditional numerical physical models frequently encounter challenges in accurately predicting cyclone landfall locations.

Consequently, this study employs a data - driven approach grounded in the Random Forest (RF) algorithm for model construction. The overall methodological framework encompasses several crucial steps:

1. Comprehending the essence of the model employed.
2. Obtaining and pre - processing relevant meteorological datasets.

3. Formulating meaningful input features.

4. Training the model and assessing its performance.

Two publicly accessible meteorological datasets are utilized in this research: the International Best Track Archive for Climate Stewardship (IBTrACS) global cyclone path dataset and the ERA5 meteorological reanalysis dataset.

The primary objective of this chapter is to elucidate the process of converting raw meteorological data into model outputs capable of making predictions. Additionally, it will expound on the rationale behind the selection of the Random Forest model. In particular, its proficiency in dealing with non - linear relationships and its resilience against overfitting render the Random Forest model an optimal choice for spatial meteorological prediction tasks.

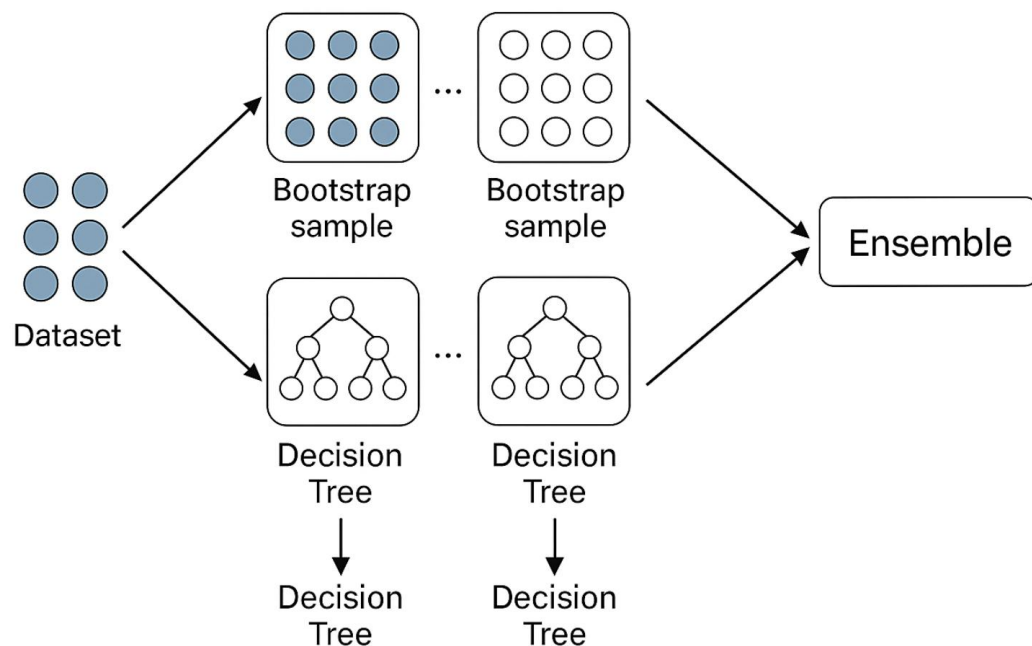
3.2 Model Overview: Random Forest

Random Forest (RF) is a supervised machine learning algorithm and an ensemble learning method that constructs multiple decision trees during the training phase. For prediction, it aggregates the outputs of all sub-models by averaging (for regression tasks) or majority voting (for classification tasks).

In this study, we employ the Random Forest regression model to predict the landing points (latitude and longitude) of tropical cyclones. The Random Forest model demonstrates exceptional performance in handling nonlinear relationships between features, mitigating overfitting, and maintaining high predictive accuracy on medium and small-scale datasets, making it particularly suitable for this research.

This algorithm leverages a mechanism known as Bootstrap Aggregation (Bagging), where each decision tree is trained on a bootstrap sample of the dataset, and only a subset of features is considered at each split. This process enhances diversity among the trees, reduces overall model variance, and improves generalization capability.

Bootstrap Aggregation (Bagging)



For regression tasks, the final prediction of the Random Forest is obtained by averaging the outputs of all individual decision trees. The mathematical formulation is presented as follows:

$$\hat{y}_{RF}(x) = \frac{1}{T} \sum_{i=1}^T f_i(x)$$

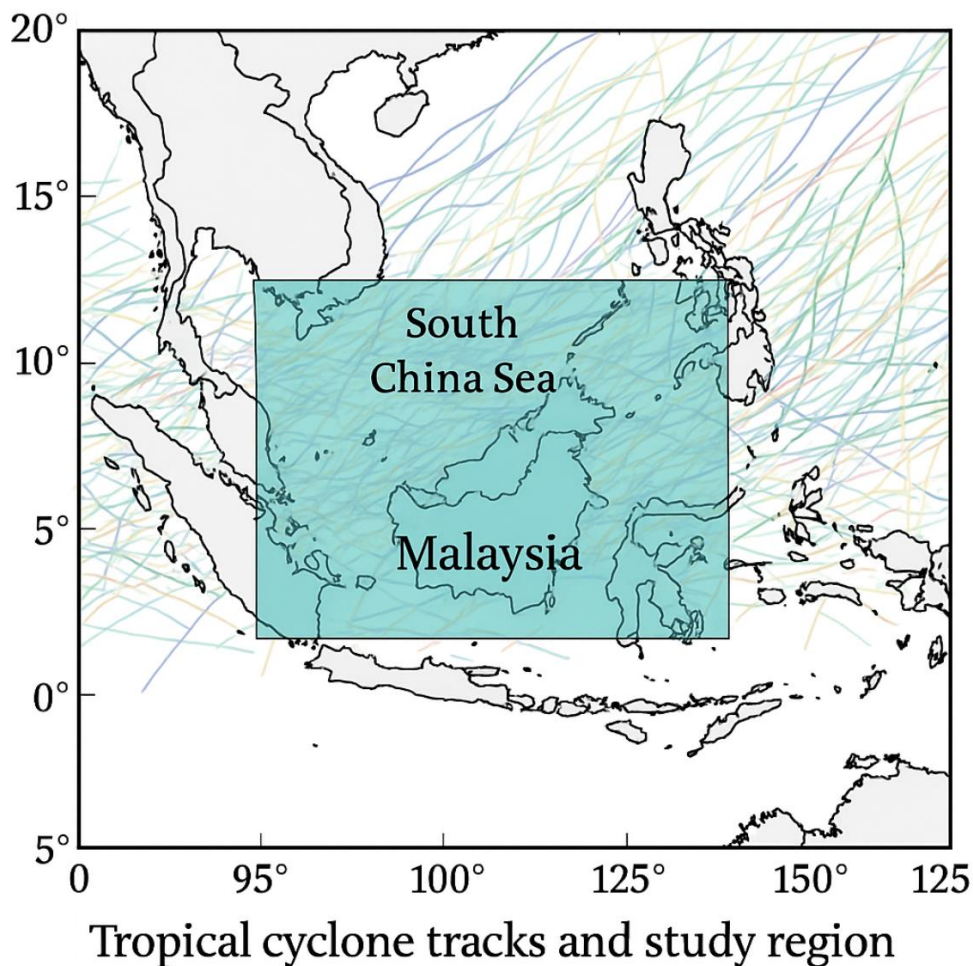
Among them:

- T denotes the total number of constructed decision trees
- $f_i(x)$ indicates the prediction outcome of the i tree for input x

The random forest model was selected because it has been extensively applied in meteorological prediction and geospatial modeling, achieving favorable results. In particular, under conditions where sample sizes are limited, feature dimensions are high, and data contain noise, the random forest model exhibits exceptional stability and robustness. Additionally, this algorithm can provide scores of feature importance, which aids in identifying the environmental variables that most significantly influence the prediction of tropical cyclone landfall locations.

3.3 Data Sources

This study utilized two primary publicly accessible meteorological datasets: the International Best Track Archive for Climate Stewardship (IBTrACS) and the ERA5 reanalysis dataset. Both datasets are extensively employed in global meteorological research and are characterized by their high data quality and extensive spatiotemporal coverage.



IBTrACS is a globally comprehensive dataset of historical tropical cyclone tracks compiled and released by the National Centers for Environmental Information

(NCEI), which operates under the U.S. National Oceanic and Atmospheric Administration (NOAA). This dataset includes critical parameters such as the time, latitude and longitude positions, central pressure, and maximum sustained wind speed of cyclones. For this study, we focus on the Northwest Pacific sub-region (ibtracs.WP) to analyze tropical cyclone activities affecting Malaysia and its surrounding areas. The dataset spans from 1980 to the present, with a temporal resolution of six hours.

[Home](#) [Products](#) [International Best Track Archive for Climate Stewardship \(IBTrACS\)](#)

International Best Track Archive for Climate Stewardship (IBTrACS)

The International Best Track Archive for Climate Stewardship (IBTrACS) project is the most complete global collection of tropical cyclones available. It merges recent and historical tropical cyclone data from multiple agencies to create a unified, publicly available, best-track dataset that improves inter-agency comparisons. IBTrACS was developed collaboratively with all the World Meteorological Organization (WMO) Regional Specialized Meteorological Centres, as well as other organizations and individuals from around the world.

To help the project receive continued support, updates, and improvement, tell us how you use IBTrACS data by completing our optional User Registration Form.

[Optional User Registration](#)



ERA5 is a high-precision global reanalysis dataset provided by the European Centre for Medium-Range Weather Forecasts (ECMWF). It offers hourly estimates of various atmospheric and oceanic variables. Key variables relevant to this study include sea surface temperature (SST), surface pressure, 10-meter wind speed and direction, and mid-level relative humidity. These variables serve as essential factors influencing the path and intensity changes of tropical cyclones.

Climate Data Store
Datasets
Applications
User guide
Live
Background

Info
26 Sep 2024
Watch our [Forum](#) for Announcements, news and other discussed topics.

ERA5 hourly data on single levels from 1940 to present

Overview

Download

Documentation

ERA5 is the fifth generation ECMWF reanalysis for the global climate and weather for the past 8 decades. Data is available from 1940 onwards. ERA5 replaces the ERA-Interim reanalysis.

Reanalysis combines model data with observations from across the world into a globally complete and consistent dataset using the laws of physics. This principle, called data assimilation, is based on the method used by numerical weather prediction centres, where every so many hours (12 hours at ECMWF) a previous forecast is combined with newly available observations in an optimal way to produce a new best estimate of the state of the atmosphere, called analysis, from which an updated, improved forecast is issued. Reanalysis works in the same way, but at reduced resolution to allow for the provision of a dataset spanning

ERA5 2 metre temperature and Mean sea level pressure
1 January 2023 at 00:00 UTC

2 metre temperature (°C)

References

[Citation and attribution](#)

DOI: [10.24381/cds.adbb2d47](https://doi.org/10.24381/cds.adbb2d47)

Licence

[Licence to use Copernicus Products](#)

Publication date

2018-06-14

Update date

2025-06-05

cds.climate.copernicus.eu

By temporally aligning the cyclone trajectory information from IBTrACS with the environmental variables in ERA5, a supervised learning training dataset can be constructed. In this dataset, cyclone positions are treated as output variables, while the corresponding environmental features at each time step are used as input variables. The study focuses on the South China Sea and surrounding waters of Malaysia, covering a geographical range approximately between 0°–20°N latitude and 95°–125°E longitude.

3.4 Data Preprocessing

Before model training, a series of preprocessing steps must be performed on the IBTrACS and ERA5 datasets to ensure a consistent data structure and complete information, thereby making them suitable for supervised learning tasks. These steps help reduce noise and improve the reliability of the model inputs.

The IBTrACS dataset can be obtained from its official website (<https://www.ncei.noaa.gov/products/international-best-track-archive>).

A1																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																				
----	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--

Given the substantial size of the IBTrACS dataset files and the inclusion of unnecessary data for our study, we will conduct a filtering process on the IBTrACS dataset in this step. Our filtering criteria are as follows:

- Time range: from 1980 to the present
- Geographical range: covering Malaysia and its surrounding seas (Longitude: approximately 95°E to 120°E; Latitude: approximately 0°N to 15°N)

3.4.2 Downloading and Filtering of the ERA5 Dataset

The ERA5 dataset can be obtained from its official website (<https://cds.climate.copernicus.eu/datasets/reanalysis-era5-single-levels?tab=overview>).

Given the substantial size of the ERA5 dataset files and the inclusion of unnecessary data for our study, we will conduct a filtering process on the ERA5 dataset in this step. Our filtering criteria are as follows:

- Time range: from 1980 to the present

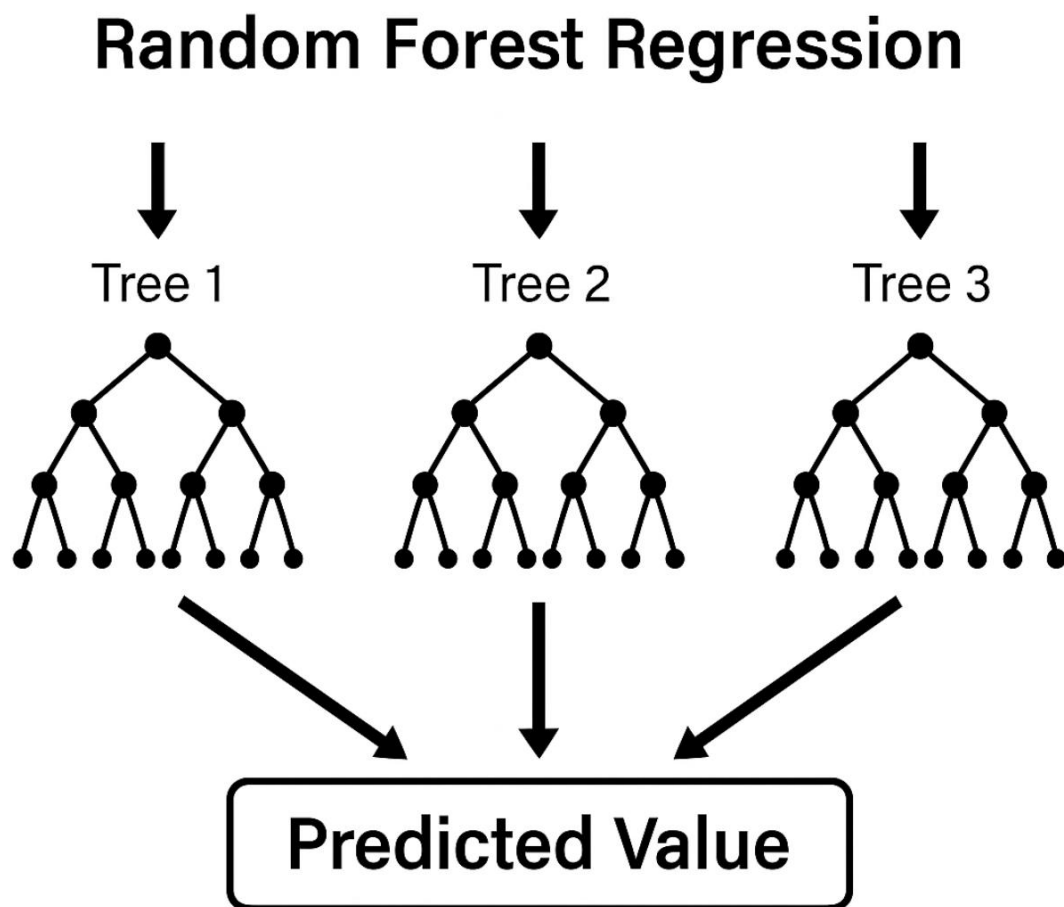
- Temporal resolution: every 6 hours

- Spatial coverage: Malaysia and its surrounding seas (Longitude: approximately 95°E to 120°E; Latitude: approximately 0°N to 15°N)

- Environmental variables for each cyclone at each time step: sea surface temperature, wind speed, air pressure, humidity, etc.

3.5 Model Structure and Justification

The machine learning model selected for this study is the Random Forest Regression (RFR) algorithm. Random Forest is an ensemble learning method that constructs multiple decision trees during training and aggregates their predictions through averaging, thereby reducing overfitting and enhancing the model's generalization ability. Each decision tree is trained on a random subset of both data samples and features, introducing randomness and diversity into the model. This structure enables Random Forest to effectively capture the complex nonlinear relationships between environmental variables (e.g., sea surface temperature, wind shear, humidity) and the landing locations of tropical cyclones (latitude and longitude).



The Random Forest model constructed in this study incorporates the following key structural parameters:

- Number of trees: Determines the total number of decision trees to be built. A higher number of trees generally improves performance but increases computational cost.

- Maximum tree depth: Controls the maximum depth of each tree. Shallower trees help mitigate overfitting by limiting the complexity of individual trees.

- Number of features to consider at each split: Restricts the number of features evaluated at each split, promoting diversity among trees and reducing correlation between them.

- Bootstrap sampling: Each tree is trained using a dataset sampled with replacement, ensuring variability in the training data for each tree.

Compared with other machine learning methods (e.g., Support Vector Regression, LSTM, Linear Regression), the Random Forest model was chosen for the following reasons:

1. Strong nonlinear modeling capability: The formation and trajectory of tropical cyclones are governed by the nonlinear interactions of various meteorological and oceanic factors. Random Forest excels at capturing such complex relationships due to its inherent flexibility.

2. Robustness against overfitting: As an ensemble method, Random Forest reduces variance by averaging predictions across multiple trees, enabling it to generalize well even in scenarios with limited data.

3. High fault tolerance: Random Forest is relatively insensitive to noise and missing data, making it well-suited for environmental datasets that often contain incomplete or noisy observations.

4. Interpretability of variable importance: Random Forest provides a mechanism to calculate feature importance, allowing researchers to identify which environmental variables most significantly influence the predicted landing locations of tropical cyclones. This enhances the interpretability and transparency of the model.

In conclusion, Random Forest is particularly well-suited for predicting tropical cyclone landing points in Malaysia and its surrounding regions, where the problem is characterized by high complexity and sparse data.

3.6 Summary

This chapter presents the methodology employed in this study for predicting the landing points of tropical cyclones in the Malaysian region. First, it elaborates on the random forest regression algorithm chosen for this study, emphasizing its strengths in nonlinear modeling, stability, and interpretability for geographical space prediction problems.

Subsequently, the chapter details the two primary data sources utilized: the IBTrACS cyclone trajectory dataset and the ERA5 environmental variable dataset. It also outlines the data acquisition and preprocessing procedures. Furthermore, the chapter examines the structural composition of the random forest model and provides justification for its selection, highlighting its effectiveness in addressing multi-variable, nonlinear, and incomplete data challenges.

The methodological framework established in this chapter serves as a foundation for the model training and evaluation presented in the subsequent chapter.