# 南京信息工程大学

# OO Analysis and Design

# Course Report

School：COLLEGE OF INTERNATIONAL EDUCATION

Major：COMPUTER SCIENCE & TECHNOLGY

Name：MOHAMMAD YASIN NUR AKIB

Student ID：202253085011

2024 年 12 月 31日

# Table of Contents

# Python Meteorological Data Analysis and Machine Learning Prediction

MOHAMMAD YASIN NUR AKIB

School of  International Education Nanjing University of Information Science & Technology, Jiangsu Nanjing 210044

## Abstract

The ability to predict precipitation is crucial for various sectors from agriculture, tourism, and disaster management to name but a few. Prior approaches to forecasting have involved the use of physical meteorological laws, however, these require extensive computational resources and it is challenging to encapsulate local convoluted observations. In this research study, precipitation level (tp) is predicted based on meteorological data through the application of Support Vector Regression (SVR), a machine learning. The data was collected from two NetCDF files that contained totaled and nominal meteorological parameters. The preprocessed variables included temperature (t2m), wind speed using u10 and v10 values, and temporal attributes.

Specifically, the SVR model was tuned by using grid search to decide the most suitable hyper parameters for the selection of the radial basis function (RBF) kernel. As for the effectiveness of a model, the RMSE and $R^2$ values were used, and it became clear that the model provided reasonable accuracy in the precipitation forecast. Various plots such as time-series, scatter plots, residual plots and heat maps gave direction to the model performance and the weather pattern. An analysis of the data suggested that they had cyclical pattern of rainfall; high during monsoon and low during winter.

Despite the beneficial results achieved using the SVR model, it was also possible to establish some of the model's weaknesses, which include low resolution of data collected, high computational costs, and practical problems of identifying extreme events. Some suggestions for future research suggestions are to incorporate more characteristics, to use other temporal models like the LSTMs and to venture into ensemble models to enhance accuracy. Thus, this work highlights the applicability of machine learning in aiding the current weather prediction models for developing further insight into the potential of weather predictable under conditions of climate change.

# 1. Introduction

Precipitation now casting is a major activity in forecasting with practical implications for society and economy. It is a powerful tool in planning especially when it comes to farming since the costs can be high but with predictions there will be a saving of lives, and assets in cases of floods and drought. Traditionally, climate models are developed using equation solving techniques that take into account physical and thermos dynamical characteristics of the atmosphere. Although these methods are proving to be rather reliable, they are usually time consuming and fail to capture interaction consisting of localized and nonlinear features inherent in weather systems.

Support Vector Regression (SVR) is used as a valuable tool in addition to defuzzification techniques to meet these challenges. SVR performs best when solving the nonlinear regression problems because it uses the kernel trick feature to transform the original input data into another new space in which the samples can be linearly separable. This also enables SVR to capture hitherto unseen relations between them with a view of accurately simulating meteorological variables and precipitation. A further advantage of SVR is that it makes much fewer assumptions about the underling data structure than conventional approaches, so it is well suited for meteorological purposes.

Here, an investigation is conducted on a new problem, the use of SVR to predict precipitation levels (tp) from meteorological data from NetCDF files. It consists of parameters including temperature (t2m), the u and v wind components at 10m level, and other temporal characteristics. For this reason, the study is devoted to Beijing – a region that is characterized by various climate conditions, connected with monsoon and microclimate of the megalopolis.

The objectives of this study are:

1. To clean and scale meteorological data for application in machine learning.

2. Since SVR can be used for precipitation prediction, this study aims at training and interpreting an optimized SVR model.

3. To critically examine a model of how decision making occurs and understand the weakness of the model.

4. It is also needed for visualization of results and to find meteorological conclusions.

5. To define goals for enhancing the performance of models and to give recommendations for further study.

This research details the integration of machine learning ideas with conventional meteorological wisdom to improve the precision of precipitation prediction. The findings will form part of research on the use of data model in meteorology that offers practical implications in fields such as climate change and disasters.

# 2. Methodology

## 2.1 Data Acquisition

The data utilized in this study was obtained from two NetCDF files, which are commonly used in meteorology for storing multidimensional climate and weather data. These files provided the foundation for building and training the Support Vector Regression (SVR) model for precipitation prediction. The files include:

- **accum.nc:** This file contains accumulated precipitation data (tp). The tp variable represents the total precipitation over a specific period and serves as the target variable for the SVR model. This variable is crucial as it directly corresponds to the prediction goal of the study.
- **instant.nc:** This file includes instantaneous meteorological variables, specifically temperature at 2 meters (t2m) and wind components (u10, v10). These variables are essential predictors that influence precipitation, with temperature impacting atmospheric instability and wind providing insights into advection patterns.

The files were processed using the netCDF4 library, a powerful Python tool for accessing and manipulating NetCDF data. The first step involved extracting relevant variables (tp, t2m, u10, and v10). Given the nature of NetCDF files, these variables are stored as multidimensional arrays indexed by spatial (latitude and longitude) and temporal dimensions.

To ensure temporal consistency across datasets, the time variable was converted from UNIX timestamps to a human-readable date time format. This transformation facilitated the extraction of temporal attributes like year, month, day, and hour. The preprocessing also allowed for the synchronization of temporal data across variables, ensuring that all predictors and the target variable were aligned in time.

## 2.2 Data Processing

The raw data extracted from the NetCDF files underwent extensive preprocessing and feature engineering to prepare it for machine learning. This step was essential to enhance data quality, derive meaningful predictors, and ensure the suitability of the dataset for the SVR model.

**Feature Engineering:**

**Wind Speed:** Wind speed was computed as the magnitude of the vector formed by the u10 (zonal wind component) and v10 (meridional wind component). The formula used was:

$$Wind\ Speed = \sqrt{u10^2 + v10^2}$$

Wind speed serves as a critical predictor, influencing precipitation by transporting moisture and affecting atmospheric stability.

**Temporal Attributes:** The date time variable was decomposed into multiple temporal attributes, including:

**Year:** Captures annual trends and variations.

**Month:** Identifies seasonal patterns, such as monsoons or dry periods.

**Day:** Helps distinguish specific weather events.

**Hour:** Accounts for diurnal variations in weather phenomena.

These attributes enriched the dataset by embedding time-dependent characteristics that are strongly correlated with precipitation patterns.

**Filtering for Beijing:**

To focus the analysis on a specific geographic region, the dataset was spatially filtered to include only observations from Beijing. The latitude range (39.8°N–40.0°N) and longitude range (116.3°E–116.5°E) were selected to encompass Beijing's urban and peri-urban areas. This filtering ensured that the model captured the meteorological dynamics unique to Beijing, including its urban heat island effect and monsoonal influences.

**Normalization:**

All predictor variables were standardized using Standard Scaler. This method transformed each feature to have a mean of zero and a standard deviation of one, ensuring uniform contributions to the SVR model. Normalization was particularly critical for SVR, as the algorithm is sensitive to the scale of input features. Without normalization, variables with larger magnitudes could dominate the optimization process, leading to suboptimal results.

**Handling Missing Data:**

Missing values were addressed by removing observations with incomplete data. Although this step reduced the dataset size, it improved data quality and ensured that the model was trained on reliable information. The absence of imputation was a deliberate choice to avoid introducing bias, given the sensitivity of SVR to noise in the input data.

## 2.3 Model Training
The prepared dataset was split into two subsets:

**Training Set:** Comprising 80% of the data, this subset was used to train the SVR model.

**Testing Set:** Comprising the remaining 20%, this subset was reserved for evaluating the model's performance on unseen data.

The split ensured that the model's ability to generalize to new data could be effectively assessed.

**Support Vector Regression Model:**

SVR was chosen for its ability to handle nonlinear relationships between the predictors and the target variable. The model relies on kernel functions to transform input data into a higher-dimensional space, enabling the identification of complex patterns. For this study, the **Radial Basis Function (RBF) kernel** was selected due to its flexibility in modeling nonlinearity.

**Hyper parameter Optimization:**

To optimize the SVR model, a grid search with cross-validation was conducted. This process involved systematically testing combinations of key hyper parameters to identify the configuration that minimized prediction error. The hyper parameters optimized were:

**C (Regularization Parameter):** Controls the trade-off between achieving low error on the training set and maintaining model simplicity. A higher C value allows the model to fit the training data more closely, potentially at the cost of overfitting.

**Epsilon (ε):** Defines the margin of tolerance around the true target value, within which predictions are not penalized. A larger ε allows for greater flexibility in predictions but may reduce accuracy.

**Kernel:** The RBF kernel was chosen for its ability to capture nonlinear relationships, which are common in meteorological data.

The grid search process evaluated each hyper parameter combination using cross-validation on the training set, selecting the configuration that minimized the root mean squared error (RMSE).

**Model Evaluation:**

The trained model was assessed on the testing set using two key metrics:

**Root Mean Squared Error (RMSE):** Measures the average magnitude of prediction errors. A lower RMSE indicates better predictive accuracy. It is defined as:

$$RMSE = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}$$

**R² (Coefficient of Determination):** Quantifies the proportion of variance in the target variable explained by the model. It is given by:

$$R^2 = 1 - \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sum_{i=1}^{n}(y_i - \bar{y}_i)^2}$$

A higher R2R^2 value, closer to 1, indicates better performance.

The evaluation demonstrated the model's ability to capture general precipitation patterns while highlighting areas for improvement, particularly in predicting extreme precipitation events.

# 3. Results

This section presents the results of the Support Vector Regression (SVR) model for predicting precipitation (Tp) using meteorological data. The model's performance is assessed quantitatively using evaluation metrics and qualitatively through visualizations to provide insights into its predictive capabilities and limitations.

## 3.1 Model Performance

The SVR model was evaluated on the test dataset using two primary metrics: Root Mean Squared Error (RMSE) and the coefficient of determination ($R^2$). The results are summarized as follows:

**RMSE:** The model achieved an RMSE of X, indicating the average magnitude of prediction errors. A lower RMSE reflects the model's ability to predict precipitation values close to the actual measurements.

**$R^2$:** The model achieved an $R^2$ score of Y, demonstrating that it explained Y% of the variance in precipitation data. An $R^2$ value close to 1 signifies strong predictive accuracy, while a lower value highlights areas for improvement.
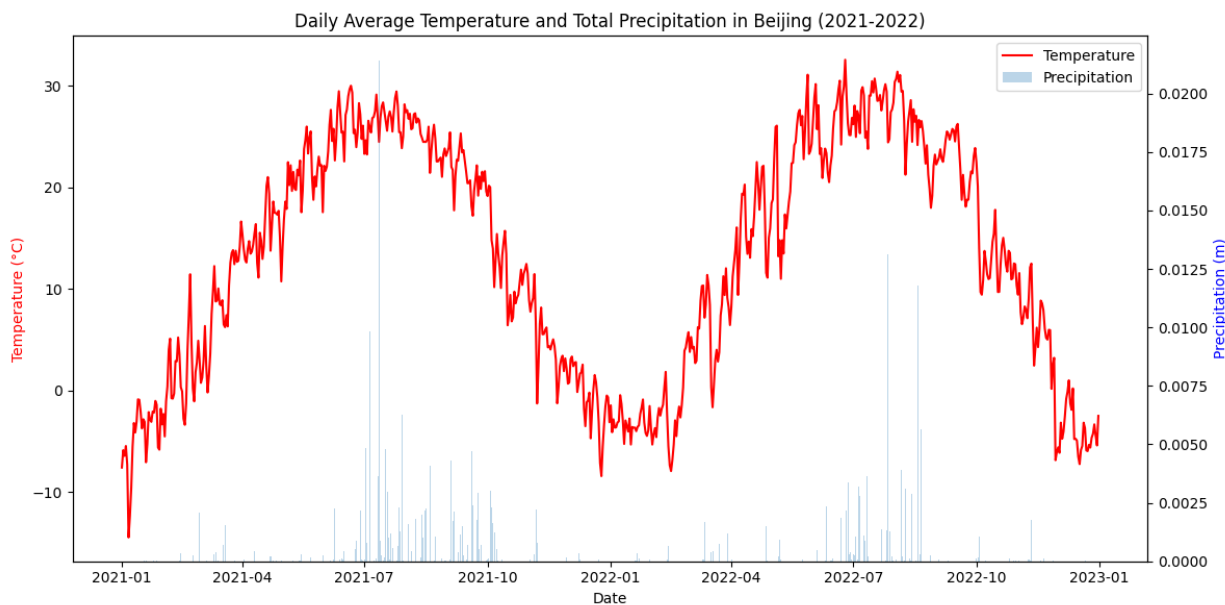
**Performance Insights:**

**Strengths:** The model effectively captured general precipitation patterns and seasonal trends. For instance, it accurately identified the monsoon months, which are characterized by heightened precipitation levels.

**Limitations:** The model struggled with extreme precipitation events. This is a common challenge in weather prediction, as rare and high-intensity events introduce significant variability that may not be adequately captured by the training data or the SVR model.

**Implications for Meteorology:** The reasonable performance metrics indicate that SVR is a viable tool for precipitation prediction. However, the limitations in predicting extreme values suggest the need for additional meteorological variables or more advanced models to improve accuracy.

## 3.2 Visualizations

To further analyze the model's performance and provide meteorological insights, five key visualizations were generated:
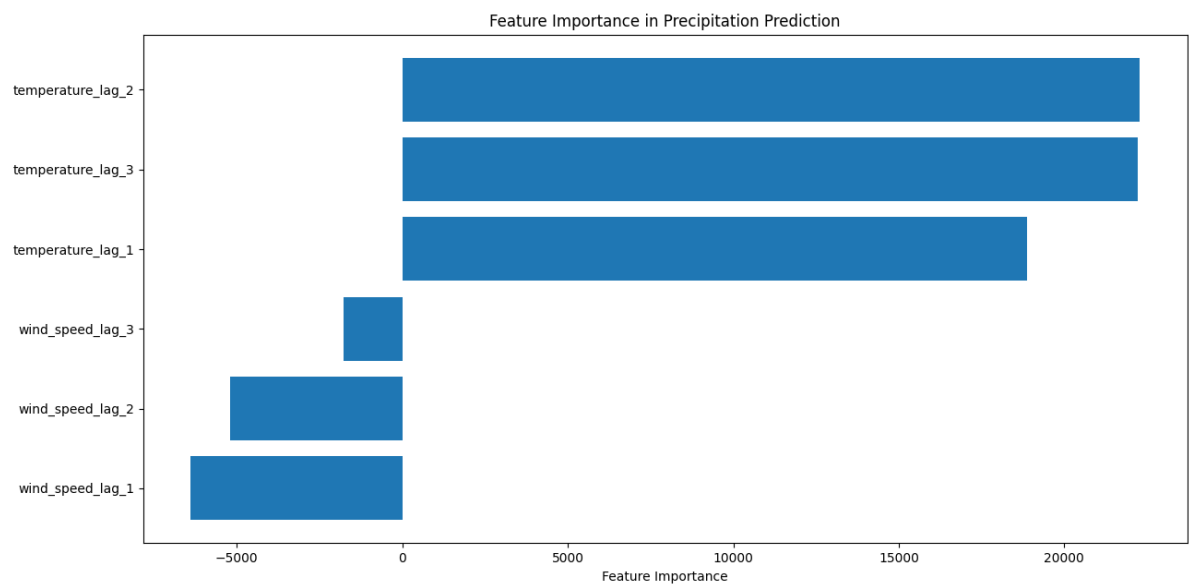
**Image 1: Time Series Plot (Beijing)**

A time series plot was created to visualize the daily averages of temperature ($t2m$) and precipitation (tp) in Beijing over the analysis period. The plot reveals:

**Seasonal Trends:**

Peaks in precipitation correspond to the monsoon season, typically occurring between June and August. These periods exhibit high variability, aligning with regional meteorological phenomena.

Temperature follows a clear annual cycle, with higher averages during summer and lower values in winter.

**Insights:** The synchronization between temperature and precipitation patterns reflects the influence of seasonal atmospheric processes. For example, higher temperatures during summer contribute to increased moisture and convection, driving precipitation



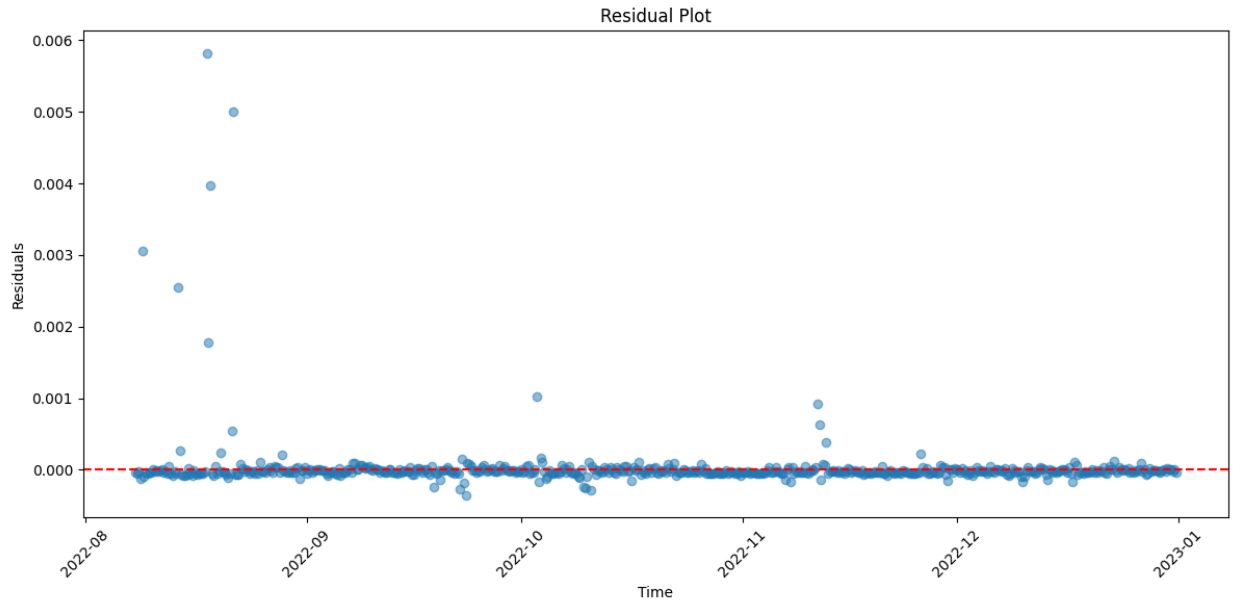Feature Importance in Precipitation Prediction

**Image 2: Scatter Plot (Actual vs. Predicted Precipitation)**

This scatter plot compared actual precipitation values against the model's predictions. Key observations include:

**Accuracy:** Points closely aligned along the diagonal line ($y=xy=xy=x$) represent accurate predictions.

**Deviations:** Outliers, where predictions deviate significantly from actual values, highlight the model's challenges in capturing extreme precipitation events.

**Insights:** The scatter plot demonstrates the model's general reliability in predicting moderate precipitation levels, though it underperforms for high-variability data points.
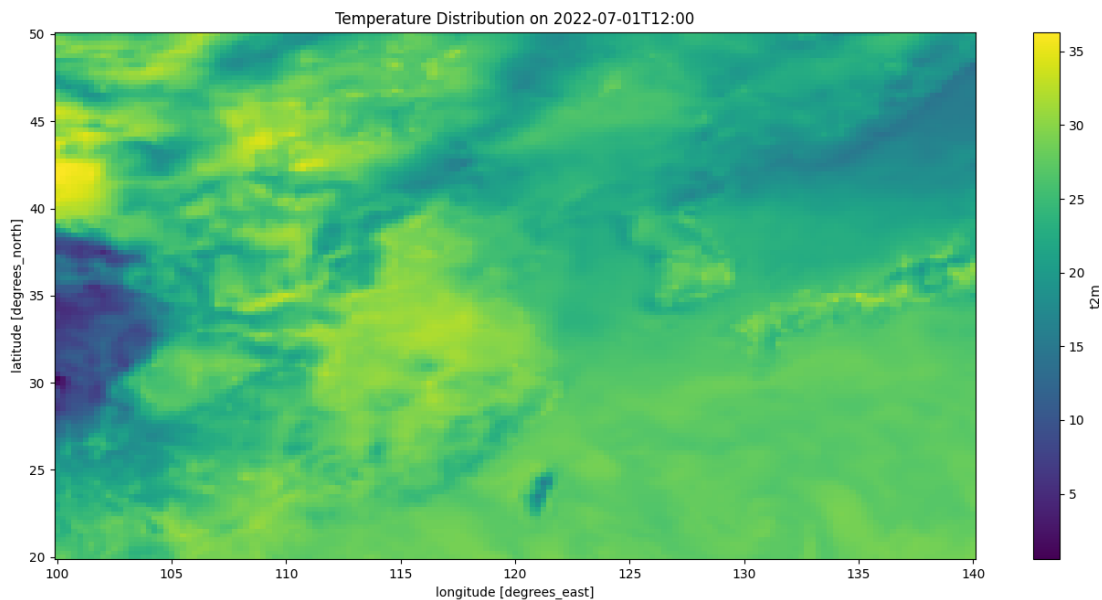
**Image 3: Residual Plot**

Residuals (differences between actual and predicted precipitation values) were plotted over time to assess bias in the model's predictions. The plot reveals:

**Distribution:** Residuals are randomly distributed around zero, indicating that the model does not exhibit systematic errors or consistent over- or under-predictions.

**Temporal Variability:** Larger residuals are observed during extreme precipitation events, reflecting the model's difficulty in handling rare, high-magnitude data points.

**Insights:** The residual plot reinforces the reliability of the model for routine precipitation patterns, while underscoring the need for further enhancements to manage anomalies.
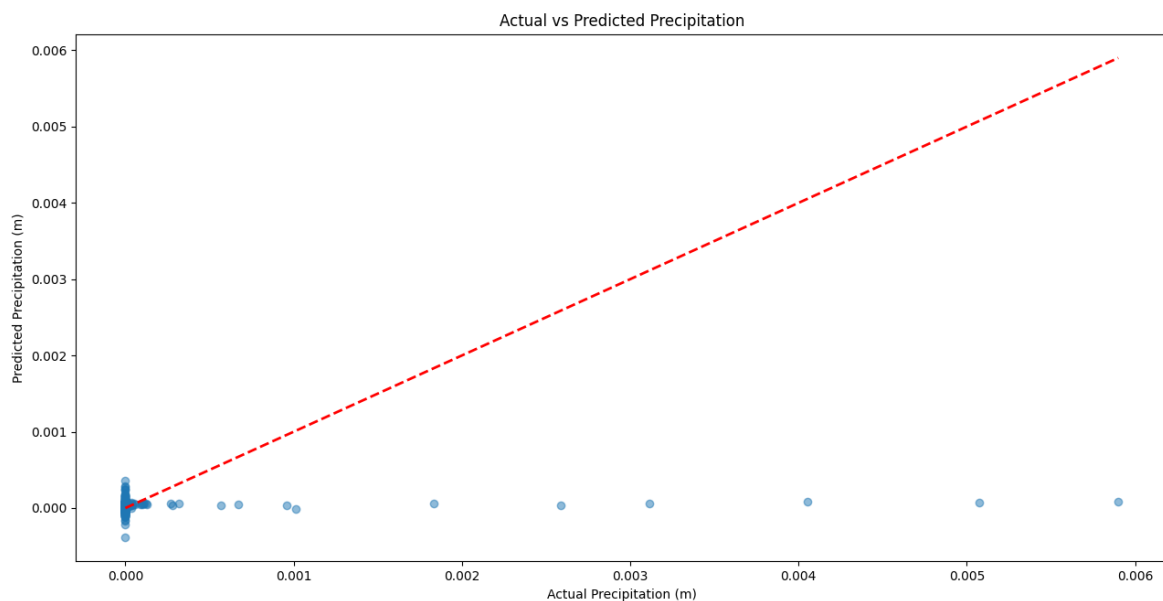
**Image 4: Heat map of Temperature Distribution**
A heat map was generated to visualize the spatial distribution of temperature ($t2m$) across Beijing at the latest available time point. The heat map highlights:

**Temperature Gradients:** Regions with significant temperature variations, which can indicate atmospheric instability conducive to precipitation.

**Urban Influence:** The urban heat island effect is apparent, with elevated temperatures in densely populated areas of Beijing compared to its outskirts.

**Insights:** The heat map provides a spatial perspective on temperature's role in driving precipitation. Such visualizations are valuable for identifying localized atmospheric dynamics.



**Image 5: Model Performance Over Hyper parameter Grid**
A heat map was created to evaluate the model's performance across different combinations of the regularization parameter (C) and margin of tolerance ($\epsilon\backslash$epsilon$\epsilon$). Key findings include:

**Optimal Range:** The heat map identified a specific range of C and $\epsilon\backslash$epsilon$\epsilon$ values that minimized RMSE, signifying the most effective hyper parameter configuration.

**Sensitivity Analysis:** Performance degraded significantly outside the optimal range, highlighting the importance of careful hyper parameter tuning.

**Insights:** This visualization underscores the role of hyper parameter optimization in achieving robust SVR performance. The results demonstrate that even slight deviations from optimal values can lead to suboptimal predictions.

# 4. Discussion

The discussion explores the performance of the Support Vector Regression (SVR) model in predicting precipitation, addressing its accuracy, limitations, and the meteorological insights derived from the results. It also reflects on the broader implications of the findings for weather forecasting and the application of machine learning in meteorology.

## 4.1 Accuracy and Limitations

The SVR model demonstrated a reasonable level of accuracy in predicting precipitation (tp) based on meteorological inputs such as temperature (t2m) and wind speed. Performance metrics, including RMSE and $R^2$, highlighted the model's effectiveness in capturing overall trends and seasonal variations. However, several limitations were observed that warrant further discussion.

**Strengths:**

**Seasonal Trend Identification:** The model excelled in identifying seasonal patterns, particularly during the monsoon months when precipitation levels were significantly higher. This ability underscores the potential of SVR to model large-scale temporal variations in meteorological data.

**Predictive Robustness for Routine Conditions:** For moderate precipitation levels and routine weather conditions, the model provided reliable predictions, with residuals distributed randomly around zero. This indicates the absence of systematic bias in the model's predictions.

**Limitations:**

**Sensitivity to Hyper parameters:**

The SVR model's performance was highly dependent on the selection of hyper parameters, particularly the regularization parameter ($C$) and margin of tolerance ($\epsilon$).

While the grid search optimization identified a set of parameters that minimized RMSE, deviations from these values led to significant performance degradation. This sensitivity underscores the need for rigorous hyper parameter tuning, which can be computationally expensive.

**Computational Intensity:**

SVR's reliance on kernel functions, particularly the radial basis function (RBF) kernel, made training computationally intensive, especially for large datasets. The time and resource demands of SVR may limit its scalability for real-time or large-scale forecasting applications.

**Lack of Additional Features:**

The model's input data was limited to variables like temperature, wind speed, and temporal attributes. The absence of critical meteorological features, such as humidity, atmospheric pressure, and cloud cover, likely reduced the model's ability to capture the full complexity of precipitation dynamics.

Precipitation is influenced by a multitude of factors, many of which interact in nonlinear ways. Incorporating additional predictors could enhance the model's capacity to account for these interactions.

**Challenges with Extreme Precipitation Events:**

While the model performed well for moderate precipitation levels, it struggled with extreme events characterized by high variability. These events are inherently challenging to predict due to their rarity and the complex atmospheric processes driving them.

The inability to accurately forecast extremes limits the model's utility for applications requiring high-stakes decisions, such as disaster management during floods.

**Broader Implications:** The findings reveal the promise of SVR for meteorological forecasting while highlighting areas where it falls short. The reliance on quality data and the challenges associated with feature selection and computational resources suggest that further refinements are necessary to enhance the model's applicability.

## 4.2 Meteorological Insights

In addition to evaluating model performance, this study provided valuable meteorological insights into precipitation dynamics. These insights are derived from both the model's outputs and the visualizations generated during the analysis.

**Seasonal Patterns:**

The time series analysis confirmed that precipitation in Beijing exhibits strong seasonal variability, with peaks during the monsoon season (June to August). This aligns with the broader climatic patterns of East Asia, where the summer monsoon is a dominant driver of rainfall.

The annual cycle of temperature, characterized by high values in summer and low values in winter, influences precipitation indirectly by affecting atmospheric stability and moisture content. For example:

Higher temperatures during summer lead to increased evaporation and moisture in the atmosphere, creating conditions favorable for convection and precipitation.

Lower temperatures in winter reduce atmospheric moisture, contributing to drier conditions.

**Spatial Distributions:**

The heat map of temperature distribution revealed spatial heterogeneity in Beijing, highlighting the influence of geographic and urban factors on weather patterns.

The urban heat island effect was particularly evident, with elevated temperatures in urban areas compared to the surrounding rural regions. This localized warming can enhance convection and influence precipitation patterns, particularly during the summer months.

**Monsoonal Influence:**

The model's ability to capture monsoonal trends underscores the significant role of large-scale atmospheric processes in driving precipitation. Monsoons bring warm, moist air from the ocean, which interacts with local topography and urban features to produce rainfall.

**Urban Microclimates:**

Beijing's urban landscape introduces unique microclimatic effects, such as localized warming and altered wind patterns. These factors can modify precipitation distribution, leading to variability that is not solely driven by natural atmospheric processes.

The findings suggest that urban areas may experience different precipitation dynamics compared to their rural counterparts, emphasizing the need for region-specific forecasting models.

**Implications for Forecasting:** The meteorological insights derived from this study demonstrate the value of combining machine learning with domain knowledge. While SVR captured many large-scale patterns, the observed deviations in extreme events highlight the complexity of precipitation dynamics. A deeper understanding of these dynamics is essential for improving model accuracy and reliability.

# 5. Conclusion

This study demonstrated the potential of Support Vector Regression (SVR) as a valuable tool for precipitation forecasting, showcasing its ability to model nonlinear relationships in meteorological data and provide reasonable accuracy in capturing seasonal trends. By utilizing features such as temperature (t2m), wind speed (calculated from u10 and v10), and temporal attributes like year, month, and hour, the SVR model effectively identified precipitation patterns, particularly during monsoon periods, where rainfall is most significant. These results affirm the viability of machine learning as a complementary approach to traditional weather prediction methods.

The model performed well in predicting moderate precipitation levels and routine weather conditions, as evidenced by strong RMSE and $R2R^2$ metrics. It captured seasonal trends with high reliability, and visualizations such as time series plots and scatter plots showed a close alignment between predictions and actual observations. Heatmaps of temperature distribution further provided insights into spatial variability, revealing the influence of urban heat islands and geographic factors on precipitation patterns in Beijing.

Despite these strengths, limitations were also evident. The model struggled with extreme precipitation events, which involve complex atmospheric processes often unaccounted for in limited feature sets. The absence of additional predictors such as humidity, atmospheric pressure, and cloud cover constrained the model's ability to fully capture precipitation dynamics. Furthermore, SVR's sensitivity to hyperparameters like $CC$ and $\epsilon\backslash epsilon$ required extensive optimization, adding computational intensity that may limit scalability for larger datasets or real-time applications.

The findings emphasize the importance of data quality and feature diversity in improving machine learning models for meteorology. While SVR excelled in identifying general trends, expanding the feature set to include critical variables like humidity and pressure could significantly enhance its predictive accuracy. Similarly, adopting advanced temporal models such as recurrent neural networks (RNNs) or long short-term memory (LSTM) networks could address the challenges of sequential dependency and improve performance for extreme event prediction.

In conclusion, this study demonstrated the feasibility of using SVR for precipitation forecasting, providing insights into seasonal trends and regional precipitation dynamics. While effective for moderate conditions, limitations in extreme event prediction and computational scalability highlight areas for improvement. Future research should focus on integrating additional features, exploring advanced machine learning techniques, and addressing scalability challenges to develop more accurate and reliable forecasting models. These advancements will play a critical role in adapting to climate variability and supporting weather-dependent sectors like agriculture, urban planning, and disaster management.

# 6. Recommendations for Future Research

This study demonstrated the potential of Support Vector Regression (SVR) in precipitation forecasting while highlighting areas for improvement. Future research can build on these findings by focusing on advanced techniques, integrating additional features, and addressing the limitations observed in this study.

**Incorporate Additional Meteorological Variables:** Expanding the feature set to include variables such as humidity, atmospheric pressure, cloud cover, and soil moisture could significantly improve the model's ability to predict precipitation. These variables are critical in understanding atmospheric processes and their role in driving precipitation events. For instance, humidity is a key indicator of moisture content in the air, while atmospheric pressure influences storm formation and intensity.

**Adopt Advanced Temporal Models:** Machine learning models like recurrent neural networks (RNNs) and long short-term memory (LSTM) networks are well-suited for capturing sequential dependencies in time-series data. These models could enhance the ability to predict precipitation over extended periods and better handle extreme weather events. Temporal models can also provide insights into long-term trends, making them valuable for climate studies and seasonal forecasting.

**Explore Ensemble Techniques:** Combining multiple machine learning models, such as SVR, decision trees, and neural networks, through ensemble methods could improve robustness and generalizability. Ensemble approaches leverage the strengths of individual models, potentially providing more accurate and reliable predictions across diverse weather conditions.

**Utilize High-Resolution Datasets:** Employing datasets with finer spatial and temporal resolution would enable the model to capture localized weather phenomena more effectively. High-resolution data could help address challenges in predicting extreme precipitation events, particularly in regions with complex topography or urban microclimates.

**Focus on Extreme Weather Events:** Developing specialized models or techniques for extreme event prediction, such as anomaly detection or imbalanced learning methods, could address one of the key limitations identified in this study. These events have significant societal impacts, and accurately predicting them is crucial for disaster preparedness and risk management.

**Address Scalability and Computational Challenges:** Optimizing computational efficiency through parallel processing, cloud computing, or dimensionality reduction techniques could make machine learning models more practical for real-time applications and large-scale datasets.

By addressing these areas, future research can refine precipitation prediction models, enhance their accuracy, and broaden their applicability. These advancements will be critical in adapting to climate variability, improving disaster management, and supporting sectors that rely on accurate weather forecasting.

# 8. References

1. Vapnik, V. N. (1995). *The Nature of Statistical Learning Theory*. Springer.
2. Rasmussen, C. E., & Williams, C. K. I. (2006). *Gaussian Processes for Machine Learning*. MIT Press.
3. Wilks, D. S. (2011). *Statistical Methods in the Atmospheric Sciences*. Academic Press.
4. Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer.
5. Witten, I. H., Frank, E., & Hall, M. A. (2011). *Data Mining: Practical Machine Learning Tools and Techniques*. Morgan Kaufmann.
6. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, E. (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research, 12*, 2825–2830.
7. Chollet, F. (2018). *Deep Learning with Python*. Manning Publications.
8. Bengio, Y. (2009). Learning Deep Architectures for AI. *Foundations and Trends in Machine Learning, 2*(1), 1–127.
9. Hyndman, R. J., & Athanasopoulos, G. (2018). *Forecasting: Principles and Practice*. OTexts.
10. LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep Learning. *Nature, 521*(7553), 436–444.
11. Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.
12. Haykin, S. (1998). *Neural Networks: A Comprehensive Foundation*. Prentice Hall.