

RAINFALL PREDICTION IN BALIKPAPAN USING SUPPORT VECTOR REGRESSION

Annida Nur Rahmayanthi¹, M. Fathurahman^{2*}, Sri Wahyuningsih³

^{1,2,3}Statistics Study Program, Department of Mathematics, Faculty of Mathematics and Natural Sciences, Mulawarman University, Samarinda, Indonesia

e-mail: ¹nidamida17@gmail.com, ^{2}fathur@fmipa.unmul.ac.id, ³swahyuningsih@fmipa.unmul.ac.id

ABSTRACT

Support Vector Regression (SVR) is a widely used supervised machine learning technique for data mining predictions. SVR defines a hyperplane in feature space to identify support vectors and maximize the margin between data points. A common challenge in machine learning is overfitting, which refers to achieving near-perfect accuracy on training data but failing to generalize unseen observations. SVR mitigates overfitting by balancing model complexity against training error, thereby producing robust, reliable predictions. Balikpapan, a city in East Kalimantan Province, frequently faces landslides and floods triggered by intense rainfall. This study aims to identify the optimal SVR model for predicting monthly rainfall in Balikpapan City. We selected the appropriate kernel using the Terasvirta test, analyzed data from January 2014 through December 2023, and evaluated four train-test splits (60:40, 70:30, 80:20, and 90:10). The Terasvirta test indicated that a linear kernel is most suitable for this dataset. The best-performing SVR model used a 90:10 split with hyperparameters $\epsilon = 0.9$ and $C = 128$, yielding 24 support vectors and a bias term of 0.028. Model performance, measured by root mean square error (RMSE) was 0.159 on the training set and 0.15 on the testing set, demonstrating strong predictive accuracy.

Keywords: Support Vector Regression, Overfitting, Linear Kernel, Terasvirta Test, Rainfall Prediction.

Cited: Rahmayanthi, A. N., Fathurahman, M., & Wahyuningsih, S. (2025). Rainfall prediction in Balikpapan using support vector regression. *Parameter: Journal of Statistics*, 5(1), 27–34. <https://doi.org/10.22487/27765660.2025.v5.i1.17640>



Copyright © 2025 Rahmayanthi et al. This is an open-access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

INTRODUCTION

Indonesia, an archipelagic nation with diverse climatic conditions, is frequently influenced by the El Niño and La Niña phenomena (Fajri, Siregar, & Sahara, 2019). These events often result in above-normal rainfall intensities, leading to flooding and landslides. Balikpapan, a major city in East Kalimantan Province, has experienced both floods and landslides. Balikpapan, a major city in East Kalimantan Province, has experienced both floods and landslides whenever rainfall exceeded typical levels during La Niña episodes (Badan Penanggulangan Bencana Daerah Provinsi Kalimantan Timur, 2023).

Nicknamed the “Oil City,” Balikpapan is the city’s second-largest urban center after Samarinda and continues to expand annually in both population and infrastructure (Alaudin, Maslina, Melawardani, & Ryka, 2022). Geographically, it borders Kutai Kartanegara Regency to the north, Penajam Paser Utara Regency to the west, and the Makassar Strait to the east and south. The terrain is predominantly hilly, with sloping areas along rivers and coastlines; 60.9 percent of its land lies more than 20 m above sea level. Like much of Indonesia, Balikpapan has a tropical climate with year-round rainfall. In 2023, the highest monthly rainfall 274.3 mm occurred in April, while August recorded the lowest at 83.1 mm. March had the most rainy days (25), and that year saw 62 fallen-tree incidents and 54 landslides, underscoring the hazards posed by heavy rain in hilly terrain (Badan Pusat Statistik, 2024).

Rainfall is a complex meteorological variable that is inherently difficult to predict (Refonaa et al, 2019). Data mining techniques, which uncover patterns and relationships within large datasets, offer a promising approach to rainfall forecasting. The five primary data mining tasks estimation, prediction, classification, clustering, and association are central to knowledge discovery in databases and support informed decision-making (Larose & Larose, 2015).

Support Vector Regression (SVR), an adaptation of Support Vector Machines (SVM) for regression problems, is a robust supervised-learning method (Smola & Scholkopf, 2004). SVR has been widely applied to prediction tasks: it constructs a hyperplane in a high-dimensional feature space using kernel functions such as linear, polynomial, and radial basis function (RBF) to handle both linear and nonlinear relationships while controlling overfitting (Lagat, Waititu, & Wanjaya, 2018).

Several studies have successfully applied SVR to rainfall prediction. Prahutama and Yasin (2015) compared linear, polynomial, and RBF kernels for weekly rainfall forecasting in Semarang, finding that the polynomial kernel yielded the highest R^2 (71.61%) with $p = 2$ and $C = 2$. In Malang Regency, Yulianto, Mahmudy, and Soebroto (2020) evaluated SVR optimized via Particle Swarm Optimization (SVR-PSO) using linear, RBF, and analysis of variance (Anova) RBF kernels; the ANOVA RBF produced the lowest RMSE (2.193) compared to the linear kernel (7.998) and the RBF kernel (27.172) for a ten-day forecast in January 2019. Siregar (2022) predicts rainfall in Medan City using SVR with an RBF kernel ($C = 0.0001$, $\gamma = 0.0005$, $\varepsilon = 1$), achieving an RMSE of 0.0388. Wati, Adriyansyah, and Sulistiana (2024) compared SVR and seasonal autoregressive integrated moving average (SARIMA) for rainfall forecasting in South Bangka. The SVR model with an RBF kernel ($C = 1,000$, $\gamma = 235$, $\varepsilon = 0.0001$) outperformed SARIMA, yielding MAPE value of 0.015.

This study aims to identify the optimal SVR configuration and apply it to predict rainfall in Balikpapan. To date, no research has specifically addressed rainfall prediction for Balikpapan, and our findings will provide a valuable decision-support tool for local agencies to mitigate flood and landslide risks.

MATERIALS AND METHODS

Support Vector Regression

Support Vector Regression (SVR) is the regression counterpart of Support Vector Machine (SVM), producing continuous real-valued outputs. Introduced by Drucker to address regression problems, SVR effectively mitigates overfitting and delivers strong predictive performance (Akbar et al., 2022).

While SVM seeks to separate data into two classes, SVR fits all data points within a “tube” around the regression function, minimizing deviations beyond a threshold ε . Given training inputs $\mathbf{x} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \subseteq R$ and corresponding outputs $y = \{y_1, y_2, \dots, y_n\} \subseteq R$, SVR approximates the target function by

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b \quad (1)$$

where $\mathbf{w}^T = [\mathbf{w}_1 \mathbf{w}_2 \dots \mathbf{w}_q]$ is the weight vector, \mathbf{x} is the vector of predictor variables, b is the bias term, and q is the number of predictor variables (Purnama & Setianingsih, 2020).

The parameters \mathbf{w} and b in Equation (1) are found by minimizing the regularized risk functional:

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n E_{\varepsilon}(y_i - f(\mathbf{x}_i)) \quad (2)$$

where:

$$E_{\varepsilon}(y_i - f(\mathbf{x}_i)) = \begin{cases} |y_i - f(\mathbf{x}_i)| - \varepsilon, & |y_i - f(\mathbf{x}_i)| > \varepsilon \\ 0, & |y_i - f(\mathbf{x}_i)| \leq \varepsilon \end{cases} \quad (3)$$

and $C > 0$ controls the trade-off between model flatness (small $\|\mathbf{w}\|$) and tolerated deviations beyond ε .

To handle points outside the ε -tube, slack variables $\xi_i, \xi_i^* \geq 0$ are introduced, yielding the primal problem:

$$\min_{\mathbf{w}, \xi_i, \xi_i^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n (\xi_i, \xi_i^*) \quad (4)$$

subject to

$$\begin{cases} y_i - \mathbf{w}^T \mathbf{x}_i - b \leq \varepsilon + \xi_i, \\ \mathbf{w}^T \mathbf{x}_i + b - y_i \leq \varepsilon + \xi_i^*, \\ \xi_i, \xi_i^* \geq 0. \end{cases}$$

The dual formulation introduces Lagrange multipliers α_i, α_i^* , resulting in

$$\max_{\alpha_i, \alpha_i^*} \frac{1}{2} \sum_{i,j=1}^n (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) \mathbf{x}_i^T \mathbf{x}_j - \varepsilon \sum_{i=1}^n (\alpha_i + \alpha_i^*) + \sum_{i=1}^n y_i (\alpha_i^* - \alpha_i) \quad (5)$$

subject to

$$\sum_{i=1}^n (\alpha_i - \alpha_i^*) = 0 \text{ and } 0 \leq \alpha_i, \alpha_i^* \leq C.$$

Support vectors are those \mathbf{x}_i for which $0 < \alpha_i < C$ or $0 < \alpha_i^* < C$. The bias b can be computed from any support vector (\mathbf{x}_i, y_i) as

$$b = y_i - \sum_{i=1}^n (\alpha_i - \alpha_i^*) \mathbf{x}_i^T \mathbf{x}_i \pm \varepsilon, \quad (6)$$

with “ $+\varepsilon$ ” if $\alpha_i^* \in (0, C)$ and “ $-\varepsilon$ ” if $\alpha_i \in (0, C)$.

The resulting linear SVR predictor is

$$f(\mathbf{x}) = \sum_{i=1}^n (\alpha_i - \alpha_i^*) \mathbf{x}_i^T \mathbf{x}_j + b. \quad (7)$$

For nonlinear relationships, inputs are mapped via $\varphi: \mathbf{x} \rightarrow \varphi(\mathbf{x})$ into a higher-dimensional feature space. The kernel trick replaces inner products $\varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}_j)$ with a kernel function $K(\mathbf{x}_i, \mathbf{x}_j)$. Common kernels include:

1. Linear

$$K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j. \quad (8)$$

2. Polynomial

$$K(\mathbf{x}_i, \mathbf{x}_j) = (\gamma \mathbf{x}_i^T \mathbf{x}_j + r)^d, \quad (9)$$

where:

$\gamma > 0$ is a scaling parameter,

$r \geq 0$ is a free coefficient (also called “offset”),

d is the polynomial degree.

3. RBF

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2), \quad (10)$$

with parameters C, γ , and ε .

4. Sigmoid

$$K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\gamma \mathbf{x}_i^T \mathbf{x}_j + r), \quad (11)$$

where:

$\gamma > 0$ controls the slope of the sigmoid,

r is the offset.

The general nonlinear SVR prediction then becomes:

$$f(\mathbf{x}) = \sum_{i=1}^n (\alpha_i - \alpha_i^*) K(\mathbf{x}_i, \mathbf{x}) + b, \quad (12)$$

where K may be any of the above kernels (linear, polynomial, RBF, or sigmoid) chosen to best capture the underlying data relationships.

Terasvirta Test

The choice between a linear and a non-linear kernel is determined by testing the model's linearity using the Terasvirta test (Terasvirta, Linc, & Granger, 1993). The hypothesis are:

$H_0 : f(\mathbf{x}_i)$ is a linear function of \mathbf{x} (model is linear)

$H_1 : f(\mathbf{x}_i)$ is a non-linear function of \mathbf{x} (model is non-linear)

The test statistic is defined as:

$$\chi^2 = nR^2 \quad (14)$$

where n is the sample size and R^2 is the coefficient of determination. Under the null hypothesis, χ^2 follows a chi-square distribution. The null hypothesis (H_0) is rejected when the $\chi^2 > \chi^2_{(\alpha, df)}$ or if the p -value is less than α .

Min-Max Normalization

The min-max normalization technique rescales each variable to the $[0, 1]$ range. The normalized value x_{ig}^* is given by (Permana & Salisah, 2022):

$$x_{ig}^* = \frac{x_{ig} - \min(x_g)}{\max(x_g) - \min(x_g)} \quad (15)$$

To recover the original value from its normalized form, use:

$$x_{ig} = x_{ig}^* (\max(x_g) - \min(x_g)) + \min(x_g) \quad (16)$$

where:

x_{ig}^* is the normalized value of the i -th observation for variable g ,

x_{ig} is the original value,

$\min(x_g)$ and $\max(x_g)$ are the minimum and maximum of variable g , respectively.

Data Sources and Research Variables

Secondary data were obtained from the National Bureau of Statistics of the Republik of Indonesia (BPS). Monthly rainfall data for Balikpapan spanning January 2014 to December 2023 were used as the response variable. The predictor variables comprise air temperature, air humidity, air pressure, and wind speed. Four train test splits: 60:40, 70:30, 80:20, and 90:10 were evaluated. The research variables are summarized in Table 1.

Table 1. Research Variables

Symbols	Variable	Description	Units
Y	Rainfall	Total rainfall in each area over a specified time.	mm
X₁	Air temperature	Measure of atmospheric heat or cold.	°C
X₂	Air humidity	The proportion of water vapor is present in the air.	%
X₃	Air pressure	Force exerted by the weight of air per unit area.	mb
X₄	Wind velocity	Horizontal speed of air movement.	m/sec

Data Analysis Steps

The data analysis for this study comprised the following steps:

1. Conduct descriptive statistical analysis to characterize rainfall, air temperature, relative humidity, air pressure, and wind speed in Balikpapan from January 2014 to December 2023.
2. Perform a linearity test using the Terasvirta procedure to identify the appropriate kernel function for the SVR model.

3. Normalize all variables using the min-max normalization method.
4. Split the dataset into training data and testing subsets using four different ratios: 60:40, 70:30, 80:20, and 90:10.
5. Develop an SVR model for rainfall prediction in Balikpapan employing the linear kernel determined by the Terasvirta test.
6. Use the tuned SVR model with optimal hyperparameters to predict monthly rainfall in Balikpapan.
7. Evaluate model performance and select the best-performing configuration based on the root mean square error (RMSE).
8. Synthesize findings and draw conclusions regarding rainfall modeling in Balikpapan.

RESULTS AND DISCUSSION

Descriptive Statistical Analysis

Descriptive statistics were used to summarize the characteristics of the dataset and provide an overall view of the research variables, as shown in Table 2.

Table 2. Descriptive Statistics for Data of Research Variables

Variables	Mean	Standard Deviation	Minimum	Maximum
Y	228.66	121.65	18.5	668
X₁	27.47	0.49	26.48	28.9
X₂	84.53	3.43	76	91.47
X₃	1,010.51	1.13	1,007.45	1,012.8
X₄	3.27	1.35	1.14	6.3

Based on Table 2, the mean monthly rainfall in Balikpapan is 228.66 mm, with a standard deviation of 121.65 mm. The highest rainfall, 668 mm, occurred in June 2019, and the lowest, 18.5 mm, in September 2018. Monthly rainfall is displayed in Figure 1. Average air temperature was 27.47 °C, with a standard deviation of 0.49 °C. The highest temperature, 28.90 °C, occurred in October 2015, while the lowest, 26.48 °C, was recorded in February 2020. Mean air humidity was 84.53 percent, with a standard deviation of 3.43 percent. It peaked at 91.47 percent in October 2020 and dipped to 76 percent in each of September and October 2023. The average air pressure was 1,010.51 mb, with a standard deviation of 1.13 mb. The highest pressure, 1,012.80 mb, was observed in October 2015, and the lowest, 1,007.45 mb, in December 2020. Finally, the mean wind speed was 3.27 m/s, with a standard deviation of 1.35 m/s. The highest wind speed, 6.30 m/s, occurred in August 2019, while the lowest, 1.14 m/s, was in January 2020.

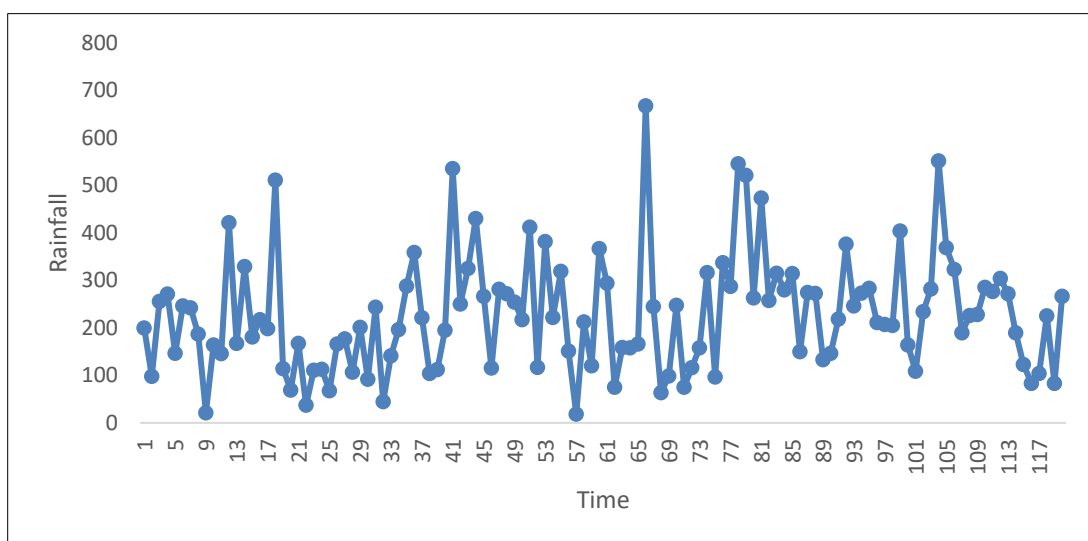


Figure 1. Scatter Plot of Monthly Rainfall in Balikpapan for the Period 2014–2023.

Terasvirta Test

Before analyzing the data with the SVR method, a Terasvirta test was performed to determine whether each model is linear or nonlinear. The hypotheses were:

H_0 : The model is linear

H_1 : The model is non-linear.

Table 3. Terasvirta Test Results

Model	df	χ^2	$\chi^2_{(\alpha, df)}$	p-value
$f(X_1) = \alpha_1 + \beta_1 X_1 + \varepsilon$	2	1.814	5.991	0.404
$f(X_2) = \alpha_2 + \beta_2 X_2 + \varepsilon$	2	0.332	5.991	0.847
$f(X_3) = \alpha_3 + \beta_3 X_3 + \varepsilon$	2	2.851	5.991	0.240
$f(X_4) = \alpha_4 + \beta_4 X_4 + \varepsilon$	2	0.029	5.991	0.986

As shown in Table 3, for all models where air temperature, air pressure, wind speed, and air humidity serve as predictor and rainfall as the response, the observed χ^2 values are below the critical value $\chi^2_{(\alpha, df)}$ and the p -values exceed the significance level ($\alpha = 0.05$). Hence, we fail to reject H_0 all models are linear, and accordingly the SVR model uses a linear kernel.

Rainfall Prediction in Balikpapan

1. SVR Model with 60:40 Data Split

Using a 60:40 split yields 72 training samples (60% of the 120 observations) and 48 testing samples. The resulting SVR model is:

$$f(x_i) = \sum_{n=1}^{72} (\alpha_n - \alpha_n^*) (x_n x_n^T) + 0.084.$$

The optimal tuning parameters are $\varepsilon = 0.7$ and $C = 16$. This model uses 27 support vectors and has a bias term $b = 0.084$.

2. SVR Model with 70:30 Data Split

With a 70:30 split, there are 84 training samples (70%) and 36 testing samples (30%). The fitted model is:

$$f(x_i) = \sum_{n=1}^{84} (\alpha_n - \alpha_n^*) (x_n x_n^T) + 0.039.$$

The best parameters from tuning are $\varepsilon = 0.9$ and $C = 16$. This configuration yields 22 support vectors and a bias $b = 0.039$.

3. SVR Model with 80:20 Data Split

An 80:20 split uses 96 training samples (80%) and 24 testing samples (20%). The model is:

$$f(x_i) = \sum_{n=1}^{96} (\alpha_n - \alpha_n^*) (x_n x_n^T) + 0.011.$$

Optimal parameters here are $\varepsilon = 0.9$ and $C = 512$. This version uses 23 support vectors and has a bias $b = 0.011$.

4. The SVR Model with 90:10 Data Proportion

With 108 training samples (90%) and 12 testing samples (10%), the model becomes:

$$f(x_i) = \sum_{n=1}^{108} (\alpha_n - \alpha_n^*) (x_n x_n^T) + 0.028.$$

The chosen tuning parameters are $\varepsilon = 0.9$ and $C = 128$. This final model uses 24 support vectors and has a bias $b = 0.028$.

Evaluation and Selection of the Best Model

The RMSE metric was used to select the best SVR model from the four train–test split configurations, as shown in Table 8.

Table 8. RMSE Comparison

Train-Test Split	RMSE
60:40	0.163
70:30	0.153
80:20	0.158
90:10	0.150

As Table 8 shows, the SVR model with a 90:10 split achieved the lowest RMSE and was therefore selected as the best model. The predictions on the test set for this model are presented in Table 9.

Table 9. Test Set Predictions of 90:10 Split

No	Prediction Value	No	Prediction Value
1	237.338	7	210.555
2	318.861	8	134.606
3	256.346	9	70.268
4	216.190	10	17.333
5	166.847	11	135.208
6	216.291	12	128.080

CONCLUSION

Support Vector Regression (SVR) is a robust machine learning method for rainfall prediction in Balikpapan, demonstrating high accuracy. Four different training-to-testing data splits were assessed (60:40, 70:30, 80:20, and 90:10) to determine optimal model performance. The selection of the SVR kernel function was guided by the Terasvirta test, which identified the linear kernel as the most suitable.

The best SVR model utilized a 90:10 data split with hyperparameters set at $\varepsilon = 0.9$ and $C = 128$, resulting in 24 support vectors and a bias term of 0.028. The model exhibited strong predictive performance, as evidenced by root mean square error (RMSE) values of 0.159 for the training set and 0.15 for the testing set. The close alignment between training and testing RMSE values indicates that the model generalizes well without overfitting.

REFERENCES

- Akbar, F., Saputra, H. W., Maulaya, A. K., Hidayat, M. F., & Rahmaddeni, R. (2022). Implementation of Decision Tree C4.5 Algorithm and Support Vector Regression for Stroke Disease Prediction. *MALCOM: Indonesian Journal of Machine Learning and Computer Science*, 2(2), 61-67.
- Alaudin, W., Maslina, Melawardani, S., & Ryka, H. (2022). Analisis Sistem Drainase pada Wilayah Rawan Banjir Simpang Jalan Manunggal-MT Haryono Balikpapan. *Jurnal TRANSUKMA*, 4(2), 76-82.
- Badan Penanggulangan Bencana Daerah Provinsi Kalimantan Timur. (2023). *Rencana Strategis Tahun 2024 - 2026*. Samarinda: Badan Penanggulangan Bencana Daerah Provinsi Kalimantan Timur.
- Badan Pusat Statistik. (2024). *Kota Balikpapan Dalam Angka 2024*. Balikpapan: Badan Pusat Statistik Kota Balikpapan.
- Fajri, C. H., Siregar, H., & Sahara. (2019). Impact of Climate Change on Food Price in the Affected Provinces of EL NINO and LA NINA Phenomenon: Case of Indonesia. *International Journal of Food and Agricultural Economics*, 7, 329-339.
- Lagat, A. K., Waititu, A. G., & Wanjoya, A. K. (2018). Support Vector Regression and Artificial Neural Network Approaches: Case of Economic Growth in East Africa Community. *American Journal of Theoretical and Applied Statistics*, 7, 67-79.

- Larose, D. T., & Larose, C. D. (2015). *Data Mining and Predictive Analytics* (2nd ed.). New Jersey: John Wiley & Sons.
- Permana, I., & Salisah, F. N. (2022). The Effect of Data Normalization on the Performance of the Classification Results of the Backpropagation Algorithm. *IJIRSE: Indonesian Journal of Informatic Research and Software Engineering*, 2(1), 67-72.
- Prahotama, A., & Yasin, H. (2015). Prediction of Weekly Rainfall in Semarang City Using Support Vector Regression (SVR) with Quadratic Loss Function. *International Journal of Science and Engineering (IJSE)*, 9(1), 13-16.
- Purnama, D. I., & Setianingsih, S. (2020). Support Vector Regression (SVR) Model for Forecasting Number of Passengers on Domestic Flights at Sultan Hasanuddin Airport Makassar. *Jurnal Matematika, Statistika, dan Komputasi*, 16(3), 391-403.
- Refonaa, J., Lakshmi, M., Abbas, R., & Raziullha, M. (2019). Rainfall Prediction Using Regression Model. *International Journal of Recent Technology and Engineering (IJRTE)*, 8(2), 2277-3878.
- Siregar, N. A. (2022). Peramalan Curah Hujan di Kota Medan Menggunakan Metode Support Vector Regression. *Journal of Informatics and Data Science*, 1(1), 1-5.
- Smola, A. J., & Scholkopf, B. (2004). A Tutorial on Support Vector Regression. *Statistics and Computing*, 14(3), 199-222.
- Terasvirta, T., Linc, F., & Granger, C. W. (1993). Power of the Neural Networks Linearity Test. *Journal of Time Series Analysis*, 14, 159-171.
- Wati, P., Adriyansyah, & Sulistiana, I. (2024). Comparison of Rainfall Prediction Results in South Bangka Regency Using Support Vector Regression and SARIMA. *CoreID Journal*, 2(3), 86-92.
- Yulianto, F., Mahmudy, W. F., & Soebroto, A. A. (2020). Comparison of Regression, Support Vector Regression (SVR), and SVR-Particle Swarm Optimization (*Journal of Information Technology and Computer Science*, 5(3), 235-246.