# Prediction Of Student Achievement Using Artificial Neural Network And Support Vector Regression At SMK TELKOM Lampung

Desi Susanti[1,*], Joko Triloka[2],

[1,2]Master of Informatics Engineering, Institut Informatika Dan Bisnis Darmajaya, Lampung,

35141, Indonesia

*Corresponding email: desi.2321211030P@mail.darmajaya.ac.id

**Abstract:** The analysis of student performance is crucial in vocational schools because it helps identify the challenges students face in preparing themselves for the workforce. By integrating data mining techniques such as Artificial Neural Networks (ANN), educators can enhance their understanding of factors that improve student learning outcomes. An artificial neural network (ANN) is composed of interconnected artificial neurons that can learn from input data and make complex predictions, including academic achievements. The structure and function of the human brain inspire ANN. This study compares the effectiveness of the artificial neural network (ANN) method with other methodologies, such as support vector regression (SVR), to predict student achievement at SMK Telkom Lampung. Primary data collected from SMK Telkom Lampung includes 4939 examples with 550 cases, 26 features, and 4 meta-attributes. Performance evaluation involves metrics such as Mean Squared Error (MSE), Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), and coefficient of determination (R2). The coefficient of determination (R2) value of the Neural Network at 0.001 is higher than the R2 value of SVR, which reaches -0.036. Research findings indicate that the Artificial Neural Network model slightly outperforms the Support Vector Regression model, with lower prediction error rates and better ability to explain data variability.

**Keywords:** Student Performance, Vocational School, Coefficient Determination, Prediction

# 1   Introduction

Vocational High Schools (SMK) are institutes of vocational education designed to equip students with the necessary skills to work in the business/industry sector or launch their independent entrepreneurial ventures. Therefore, SMK students are expected to be ready to enter the workforce. Thus, students are required to have skills and professional attitudes in their respective fields [1]. However, in practice, many students still lack skills or expertise in their vocational fields, even though they are given opportunities to learn through projects and hands-on practice in laboratories, as well as receiving materials either directly or through digital platforms such as Learning Management Systems (LMS). Therefore, to evaluate the effectiveness of educational programs at SMK Telkom Lampung, it is necessary to predict student achievements by comparing predictions with actual results. Schools can assess how well their educational programs are working and identify areas that need improvement.

The educational organization is among the establishments that own a substantial quantity of data. [2]. This organization uses data to gather information, particularly regarding students. Numerous characteristics of student data enable us to anticipate things like the academic and extracurricular accomplishments of children during their time in school. Academic performance is frequently employed as a measure of educational success, even though in the autonomous curriculum, pupils are not assessed using a ranking or rating system in their academic achievements. It is seen as inaccurate to portray students' potential and talents. The ranking of pupils' academic performance serves as one measure of this accomplishment. The responsibilities that educators play, the drive and self-control of students, their socioeconomic backgrounds, and previous learning results all play a part in achieving high-quality education. Data mining is the process of identifying patterns in a large database to uncover hidden information using statistical approaches [3]. Data mining with the assistance of machine learning is very useful for solving various kinds of problems [4]. Sentiment analysis approaches are applied to data using machine learning (ML) techniques [5]. Schools can ensure students' academic and non-academic success by taking appropriate measures promptly after assessing their performance. The ultimate goal is for all students, regardless of their background, to reach their full potential both academically and in extracurricular activities.

In the field of education, it is essential to conduct research that predicts the success rates of students' academic studies to propose enhancements for future academic performance [6]. Numerous research studies have been conducted to forecast the academic success of children. For example, in the study 'Comparison of Data Mining Algorithm Performance for Student Graduation Prediction' by Sadimin and Handoyo Widi Nugroho, students' cumulative grade point average (GPA) was utilized as an indicator of their learning achievement [2]. Yahia Bashar et al. conducted a systematic literature review titled 'Towards Predicting Academic Performance of Students Using Artificial Neural Networks (ANN)', in which they investigated the topic. The comprehensive review concluded that artificial neural networks (ANN) have shown high accuracy in predicting academic achievement outcomes despite similar results being achieved with other data mining methodologies [7]. Ashraf Mohamed Hemeida conducted similar research in his journal titled "Nature-Inspired Algorithms for Feed-Forward Neural Network Classifiers: A Decade of Research Survey." The efficiency of training feed-forward neural networks (FFNN) is influenced by data quality, which includes accuracy, flexibility, and precision [8]. Another study

titled "Support Vector Regression" found that the SVR filtering approach reduces false positives in automated mass detection systems. SVR is also used to estimate the Target-to-Interferer Ratio (TIR), and the equivalence between TV regularization and SVR has been proven in the context of tube formulation [9].

Some identified gaps in this research include an incomplete understanding of predictive phenomena, constraints in modifying input variables, and a lack of focus on the coefficient of determination (R2) as an evaluation metric, which could be crucial in regression contexts to assess how well models predict dependent variables.

Educational data mining has emerged as one of the most popular areas of scientific inquiry in recent times [10]. One machine learning and data mining approach that has been employed in various research publications is artificial neural networks (ANN), which are claimed to produce better and more accurate results when predicting student performance [7]. This model disregards physical processes entirely and instead relies on a collection of linear and non-linear mathematical equations. Its ability to closely approximate real outcomes is its most crucial feature. Due to its training and pattern recognition capabilities, ANN is frequently utilized [11]. Support Vector Regression (SVR), along with ANN, is a well-known machine learning technique used in statistical methods for constructing regression functions [9]. SVR is an effective method for making predictions based on past data, and it works especially well with high-dimensional datasets that use Kernel functions to handle nonlinear scenarios [12].

This study utilizes an artificial neural network to predict student performance at SMK Telkom Lampung, thereby addressing the issues mentioned above. By contrasting the ANN and SVR approaches in predicting student achievement and utilizing the orange program as an analytical tool, the research aims to systematically and comprehensively compare these two methods, leveraging advancements in technology.

Therefore, this study compares the effectiveness of ANN and SVR approaches in forecasting student achievement at SMK Telkom Lampung and explores the potential of artificial intelligence technology in the educational setting. This analysis seeks to understand the efficacy of both approaches in predicting student performance in this unique educational environment.

## 2 Research Method

The phases of this study correspond to the actions shown in the flowchart in Figure 1. The figure is an explanation of the research stages mentioned.

### 2.1 Problem Scoping

The Vocational High School (SMK) Telkom Lampung is the primary research site for this study, where several core variables—such as current class, absenteeism, tardiness, ranking, extracurricular activities, and academic and non-academic student achievement data—are used in the problem scoping process. Orange software is employed to implement and analyze the focused analytic methods, which include Artificial Neural Network (ANN) and Support Vector Regression (SVR).
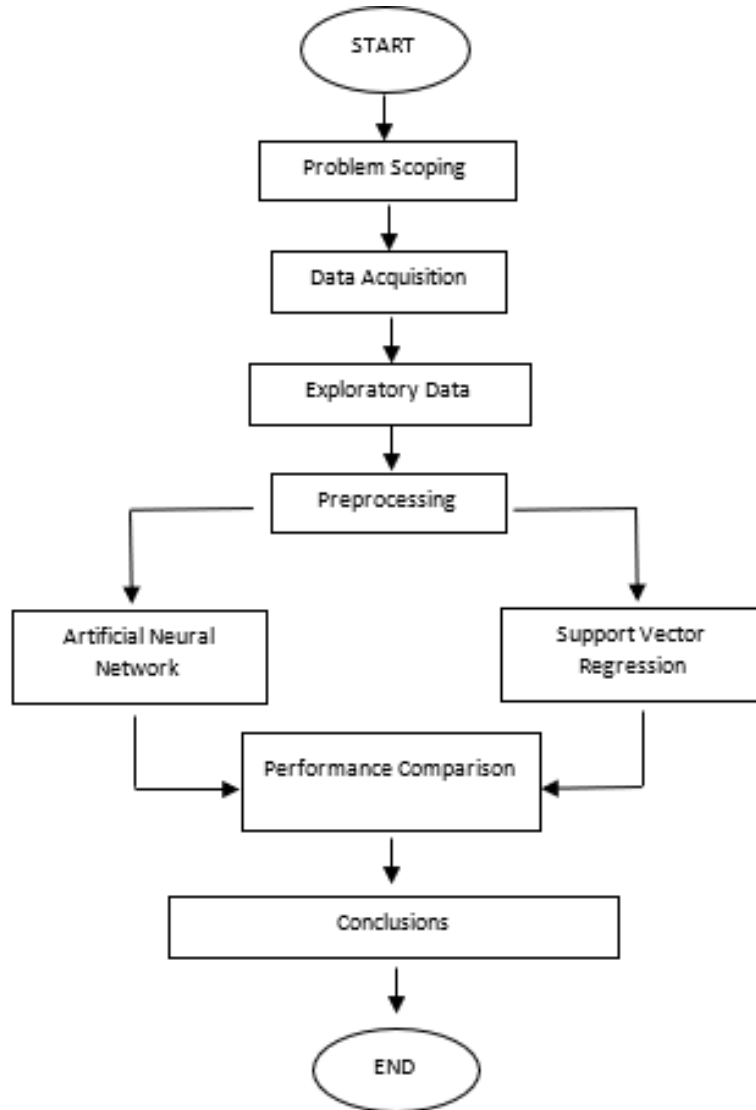
Figure 1: Research Stages

## 2.2  Data Acquisition

Data acquisition for this research involves several steps, as:

### 2.2.1  Data Collection

At this point, decisions are made about which data features to use, where to save the data from, and whether the goals of the data collection are met by the acquired data [13]. Therefore, in this study, the collected data is sourced from the school database from 2021 to 2024.

The quantity of records in the dataset is referred to as instances. The dataset contains a total of 550 samples.

### 2.2.2 Data Cleansing

The data cleansing process involves identifying and handling duplicate data, printing errors, and data inconsistencies. This process also includes verifying and correcting inaccurate or incomplete data. These steps are essential to ensure good data quality before further analysis.

### 2.2.3 Data Transformation

The format used for data mining processing is the Excel format. Data transformation is performed by importing the data into the Orange Data Mining software and then carrying out various preprocessing operations such as data cleansing, removal of missing values, normalization, or recoding of categorical features to prepare the data for further analysis.

### 2.2.4 Data Partitioning

In this study, the data is split into 80% for training the model and 20% for testing its performance. The training set is utilized to develop the prediction model, while the testing set assesses how well the trained model performs.

### 2.2.5 Data Validation

In this stage, consistency checks, research objective verification, outlier testing, cross-validation, and sensitivity testing are conducted.

## 2.3 Exploratory Data Analysis

To understand the characteristics of data in depth, it is necessary to examine descriptive statistics of the data and create graphs or plots to visualize data distributions, relationships between variables, and other patterns using Box plots and scatter plot widgets.

## 2.4 Preprocessing

Cleaning is typically a part of the data preprocessing steps [14]. Data cleaning involves checking for outliers, duplicate data, and missing values [15]. Clear and highly accurate information can be obtained through appropriate preprocessing procedures [16]. As part of the cleaning procedure, the researcher uses the Microsoft Excel tool to filter each column separately and identify empty data. There are numerous inaccuracies in data entries across columns, such as date of birth, class, number of siblings, and distance from home. These data points are assigned four scores in the prediction column based on factors like attendance, tardiness, academic achievement, and non-academic achievement. Students who still have records after expulsion or transfer are given a score of 0 to indicate erroneous data. Data with minimal influence is assigned a value of 1, moderate influence receives a score of 2, and significant influence receives a score of 3.

## 2.5 Prediction

Prophecy and estimation are synonymous with the term prediction [17]. Predictability refers to the ability to foresee future events or outcomes based on measurements, observations, gathered data, or study results that indicate specific patterns in phenomena [18].

## 2.6 Artificial Neural Network

One of the most popular methods in artificial intelligence is the neural network (ANN), which allows the algorithm to learn on its own using training data [19]. Artificial neural networks (ANNs) can broadly be classified into three categories: (1) topology, which refers to the layout of connections between neurons; (2) training or learning algorithms, which are methods used to adjust the weights on connections; and (3) activation functions [20]. In this study, the method employed is a feedforward neural network with parameter settings of two hidden layers, each containing 100 and 20 neurons respectively. The activation function utilized is ReLU (Rectified Linear Unit), and the optimization method applied is Adam. The Feedforward Neural Network model, which consists of many interconnected neurons arranged in complex layers, is a flexible model for constructing non-linear regression models, data reduction, and non-linear dynamic systems, capable of processing large volumes of data and producing accurate predictions [21].

## 2.7 Support Vector Regression

Support Vector Regression (SVR) is an extension of Support Vector Machine (SVM), originally used for classification problems, that applies SVM principles to solve regression problems [22]. This research focuses on tuning the parameters of the radial base function kernel. The parameter C (Complexity cost) is set to 1.00, while the parameter $\gamma$ (gamma) is set to 0.10. The performance of the model is measured using the R-square accuracy value, where a value approaching 1 indicates better model performance [22].

## 2.8 Performance Comparison

Absolute Error (MAE), and Coefficient of Determination (R2) are compared to assess the performance of Artificial Neural Networks (ANN) and Support Vector Regression (SVR). The following are the formulas for each evaluation metric frequently used in prediction: The mean squared error, or MSE, is the average of the squared discrepancies between the actual values ($Y_{\text{true}}$) and the predicted values ($Y_{\text{pred}}$).

$$MSE = \frac{1}{n} \sum_{i=1}^{n} (Y_{\text{true}_i} - Y_{\text{pred}_i})^2 \tag{1}$$

Where $n$ is the number of samples. $Y_{\text{true}_i}$ is the actual value of sample $i$. $Y_{\text{false}_i}$ is the predicted value of sample $i$.

The root mean square error or RMSE, is a metric that quantifies the average deviation between anticipated and observed values.

$$RMSE = sqrt\frac{1}{n} \sum_{i=1}^{n} (Y_{\text{true}_i} - Y_{\text{pred}_i})^2 \tag{2}$$

The mean absolute error or MAE is the average of the absolute discrepancies between the expected and actual values.

$$MAE = \frac{1}{n} \sum_{i=1}^{n} |Y_{\text{true}_i} - Y_{\text{pred}_i}| \tag{3}$$

Coefficient of Determination ($R^2$ score): This metric quantifies how much of the variability observed in the dependent variable can be explained by the model.

$$R^2 \text{score} = 1 - \frac{\sum_{i=1}^{n}(Y_{\text{true}_i} - Y_{\text{pred}_i})^2}{\sum_{i=1}^{n}(Y_{\text{true}_i} - \bar{Y}_{\text{true}})^2} \tag{4}$$

Where $\bar{Y}_{\text{true}}$ is the mean of the actual values true $Y_{\text{true}}$.

# 3 Results

The results of predicting student achievement at SMK Telkom Lampung using ANN and SVR methods are as follows:

## 3.1 Dataset

The data used in this study is primary data obtained from SMK TELKOM Lampung. The total number of data is 4939, with 550 instances, 26 features, and 4Meta-Attributes.

## 3.2 Data Mining Process

To select the best method with high accuracy, a comparison of several data mining methods is conducted by ANN and SVR as shown in Figure 2.
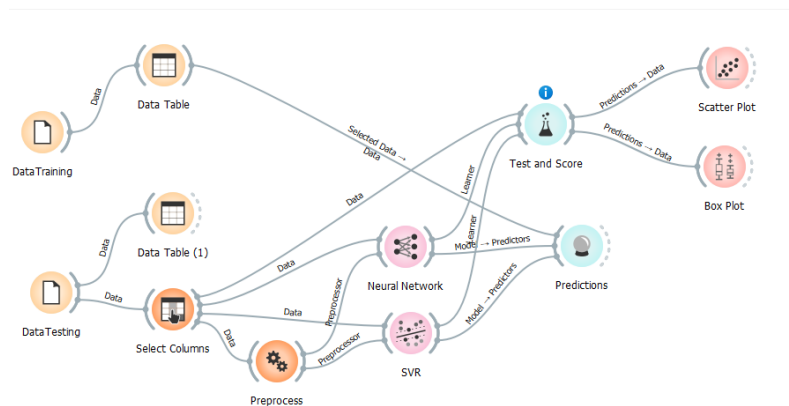


Figure 2: Orange Tools Widget Design

The data mining process, from raw data processing to evaluation and result visualization. It includes several important steps to build and evaluate prediction models using two

Table 1: Data Exploration

| No | Name Parameter | Status | Data Type |
|----|----------------|--------|-----------|
| 1 | Name | Input | Nominal |
| 2 | NISN | Input | Numerik |
| 3 | Place of Birth | Input | Nominal |
| 4 | Date of Birth | Input | Numeric |
| 5 | Gender | Input | Nominal |
| 6 | Religion | Input | Nominal |
| 7 | Address | Input | Nominal |
| 8 | Subdistrict | Input | Nominal |
| 9 | Child Number | Input | Numeric |
| 10 | Number of Siblings | Input | Numeric |
| 11 | Type of Residence | Input | Nominal |
| 12 | Transportation | Input | Nominal |
| 13 | Current Class | Input | Nominal |
| 14 | Distance from Home to School (KM) | Input | Numeric |
| 15 | KPS Recipient Input | Input | Nominal |
| 16 | KIP Recipient | Input | Nominal |
| 17 | Eligible for PIP (proposed by the school) | Input | Nominal |
| 18 | Reason for Eligibility for PIP | Input | Nominal |
| 19 | Special Needs | Input | Nominal |
| 20 | Sum of Value Scores | Input | Numeric |
| 21 | Mean | Input | Numeric |
| 22 | Rank | Input | Numeric |
| 23 | Nonattendance | Input | Numeric |
| 24 | Delay Input | Input | Numerik |
| 25 | Extracurricular Activities | Input | Nominal |
| 26 | Academic Achievement | Input | Nominal |
| 27 | Non-Academic Achievement | Input | Nominal |
| 28 | Prediction | Output | Numeric |

different techniques: Neural Network and Support Vector Regression (SVR). The dataset is divided into two parts, training data, and testing data, which are then loaded into their respective Data Table widgets. The next step is featuring selection using the Select Columns widget and data preprocessing with the Preprocess widget to ensure the data is clean and ready for analysis. Two prediction models are built using Neural Networks and SVR, with each model trained using processed data. The trained models are evaluated using the Test and Score widget, which produces performance metrics such as Mean Squared Error (MSE), Root Mean Squared Error (RMSE), and R-squared ($R^2$). Predictions are generated using the Predictions widget, and the results are visualized using scatter plots and box plots to understand the distribution and performance of the models in more detail.

## 3.3   Performance Comparison of ANN and SVR

Figure 3 displays the results of each model's calculations based on the tested data. Metric-based evaluation indicates that the Neural Network model outperforms the Support Vector Regression (SVR) model. Specifically, the Neural Network model produces the following results: $R^2$: 0.001, MAE: 0.392, RMSE: 0.594, MSE: 0.353.

| Model | MSE | RMSE | MAE | R2 |
|---|---|---|---|---|
| SVR | 0.366 | 0.605 | 0.303 | -0.036 |
| Neural Network | 0.353 | 0.594 | 0.392 | 0.001 |

Compare models by: Coefficient of variation of the RMSI ⌄   ☐ Negligible diff.:   0.1

| | SVR | Neural N... |
|---|---|---|
| SVR | | 0.923 |
| Neural Network | 0.077 | |

Figure 3: Test and Score Widget Results

However, despite the SVR model having stronger MSE, RMSE, and MAE values (0.366, 0.605, and 0.303, respectively), the $R^2$ value of -0.036 suggests that the model is not suitable for the data. In contrast, the neural network model can produce predictions that are more accurate than the actual values, even though its $R^2$ value is very close to zero. Therefore, the evaluation concludes that the neural network model performs better overall.
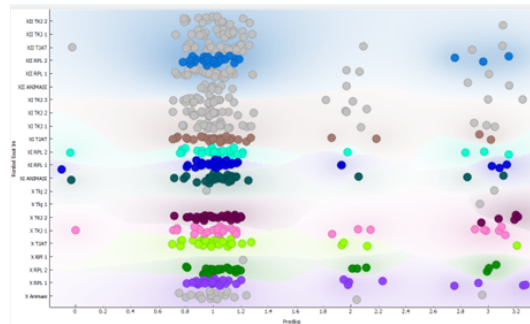


Figure 4: Scatter Plot of Prediction Data

Only the visualization between predictions and classes is displayed in the prediction scatter plot in Figure 4.

With a Pearson correlation coefficient (r) of 0.08, the scatter plot of predictions made using the Neural Network (NN) demonstrates a weak correlation between the model's predictions and the actual values. This indicates little association between the Neural Network model's predictions and the actual values, resulting in a nearly flat regression line.

With a Pearson correlation coefficient (r) of zero, or -0.00, the scatter plot of predictions using Support Vector Regression (SVR) indicates that there is no meaningful correlation
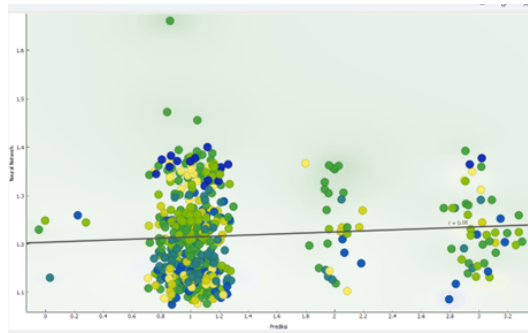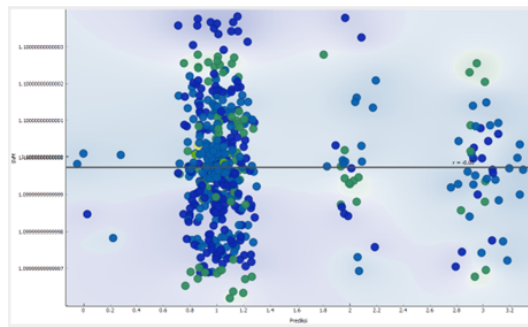
Figure 5: ANN Scatter Plot



Figure 6: SVR Scatter Plot

between the model predictions and the actual values. This suggests that there is no discernible regression line because the SVR model's predictions and the actual values do not follow a linear relationship.
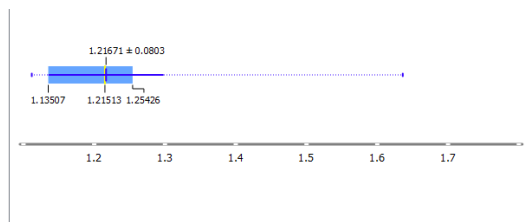


Figure 7: ANN Box Plot

The Neural Network (NN) box plot findings indicate that the model's prediction distribution has a mean of 1.21671 and a standard deviation of 0.0803. This suggests that the NN model's predictions frequently center around an average value slightly above 1.21 with low variability of about 0.0803.

The Support Vector Regression (SVR) box plot findings, as shown in Figure 8, indicate that the model's prediction distribution has a mean of 1.1 with negligible variance and a
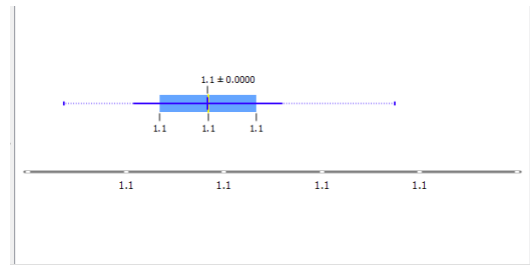
Figure 8: SVR Box Plot

standard deviation of 0.0000, which is very close to zero. This suggests that the predictions of the SVR model are generally stable and do not vary significantly around the value of 1.1.

# 4   Discussion

In the measurements, the Neural Network showed a lower error rate with an MSE of 0.353 and an RMSE of 0.594, while the SVR had an MSE of 0.366 and an RMSE of 0.605. This indicates that the predictions from the Neural Network are generally more accurate than those from the SVR. Although the coefficient of determination ($R^2$) of the Neural Network is still low at 0.001, this value is higher compared to the $R^2$ of the SVR, which is -0.036. This shows that the Neural Network is better at explaining student performance predictions at SMK Telkom Lampung.

This result is also supported by the scatter plot, where the Neural Network shows a low correlation between predictions and actual values with a Pearson correlation coefficient (r) of 0.08. In contrast, the SVR shows no significant correlation with a Pearson correlation coefficient close to zero (-0.00). The box plot shows that the predictions from the Neural Network have greater variability compared to the SVR, with a mean value of 1.21671 and a standard deviation of 0.0803. In contrast, the SVR has a mean prediction value of 1.1 with no significant variation.

This indicates that the Neural Network, with its potential to be more sensitive to variations in the data, could be more adaptable. However, this also suggests a potential for overfitting. In another study titled "Gold Price Prediction Using Algorithms Support Vector Regression (SVR) and Linear Regression (LR)", the research focused merely on one evaluation metric, explicitly MSE, without considering factors such as model complexity, prediction stability over time, or a more holistic interpretation of the results [23]. A more comprehensive evaluation considering various metrics and other factors would provide a better understanding of the overall performance of the model.

# 5   Conclusion

The results of the study show that the Neural Network (NN) outperforms Support Vector Regression (SVR) in predicting the academic performance of students at SMK Telkom Lampung. While SVR produces a negative $R^2$ value, the Neural Network (NN) approaches zero and makes more accurate predictions. Additionally, the data reveals that although

NN's predictions are not as well correlated with the actual values as SVR's, they are still better than SVR's weak correlation. Furthermore, NN's prediction distribution is slightly more variable than SVR's stable distribution. Therefore, to enhance forecast accuracy and deepen our understanding of this phenomenon, future studies should consider incorporating additional data, modifying the model, and further examining student performance.

# References

[1] S. M. Susilo and I. Ismiyati, "Pengaruh praktik kerja industri, informasi dunia kerja dan motivasi memasuki dunia kerja terhadap kesiapan kerja siswa," *Business and Accounting Education Journal*, vol. 1, pp. 290–296, Dec. 2020.

[2] S. Sadimin and H. W. Nugroho, "Perbandingan kinerja algoritma datamining untuk prediksi kelulusan mahasiwa," *Jurnal Teknoinfo*, vol. 17, pp. 512–520, July 2023.

[3] S. Mukodimah and C. Fauzi, "Comparison of tree implementation, regression logistics, and random forest to detect iris types," *Jurnal TAM (Technology Acceptance Model)*, vol. 12, p. 149, Nov. 2021.

[4] A. Ishaq, S. Sadiq, M. Umer, S. Ullah, S. Mirjalili, V. Rupapara, and M. Nappi, "Improving the prediction of heart failure patients' survival using smote and effective data mining techniques," *IEEE Access*, vol. 9, pp. 39707–39716, 2021.

[5] M. Baker, K. Jihad, and Y. Taher, "Prediction of people sentiments on twitter using machine learning classifiers during russian aggression in ukraine," *Jordanian Journal of Computers and Information Technology*, no. 0, p. 1, 2023.

[6] D. Kurniadi, E. Abdurachman, H. L. H. S. Warnars, and W. Suparta, "Predicting student performance with multi-level representation in an intelligent academic recommender system using backpropagation neural network," 2021.

[7] Y. Baashar, G. Alkawsi, A. Mustafa, A. A. Alkahtani, Y. A. Alsariera, A. Q. Ali, W. Hashim, and S. K. Tiong, "Toward predicting student's academic performance using artificial neural networks (anns)," *Applied Sciences*, vol. 12, p. 1289, Jan. 2022.

[8] A. M. Hemeida, S. A. Hassan, A.-A. A. Mohamed, S. Alkhalaf, M. M. Mahmoud, T. Senjyu, and A. B. El-Din, "Nature-inspired algorithms for feed-forward neural network classifiers: A survey of one decade of research," *Ain Shams Engineering Journal*, vol. 11, pp. 659–675, Sept. 2020.

[9] F. Zhang and L. J. O'Donnell, "Support vector regression," in *Machine Learning*, pp. 123–140, Elsevier, 2020.

[10] D. Wang, D. Lian, Y. Xing, S. Dong, X. Sun, and J. Yu, "Analysis and prediction of influencing Factors of College Student Achievement Based on Machine Learning," *Frontiers in Psychology*, vol. 13, p. 881859, Apr. 2022.

[11] A. Cetinkaya and O. K. Baykan, "Prediction of middle school students' programming talent using artificial neural networks," *Engineering Science and Technology, an International Journal*, vol. 23, pp. 1301–1307, Dec. 2020.

[12] R. E. Cahyono, J. P. Sugiono, and S. Tjandra, "Analisis kinerja metode support vector regression (svr) dalam memprediksi indeks harga konsumen," *JTIM : Jurnal Teknologi Informasi dan Multimedia*, vol. 1, pp. 106–116, Aug. 2019.

[13] M. Yağcı, "Educational data mining: prediction of students' academic performance using machine learning algorithms," *Smart Learning Environments*, vol. 9, p. 11, Dec. 2022.

[14] S. S. Berutu, H. Budiati, J. Jatmika, and F. Gulo, "Data preprocessing approach for machine learning-based sentiment classification," *JURNAL INFOTEL*, vol. 15, pp. 317–325, Nov. 2023.

[15] A. Purwanto and H. W. Nugroho, "Analisa perbandingan kinerja algoritma c4.5 dan algoritma k-nearest neighbors untuk klasifikasi penerima beasiswa," *Jurnal Teknoinfo*, vol. 17, pp. 236–243, Jan. 2023.

[16] R. Toro and S. Lestari, "Perbandingan algoritma data mining untuk penentuan lokasi promosi penerimaan mahasiswa baru pada iib darmajaya lampung," *Techno.Com*, vol. 22, no. 1, pp. 223–234, 2023.

[17] F. Maylani, S. , and N. , "Implementasi metode data mining untuk memprediksi warna anak kucing pada proses pengembangbiakan kucing ras menggunakan algoritma support vector machine (svm)," pp. 114–125, 2021.

[18] R. Ginting and C. Humaira, "Penerapan data mining: Prediksi penjualan mobil toyota menggunakan artificial neural network pada software orange," *Talenta Conference Series: Energy and Engineering (EE)*, vol. 4, Oct. 2021.

[19] M. Uzair and N. Jamil, "Effects of hidden layers on the efficiency of neural networks," in *2020 IEEE 23rd International Multitopic Conference (INMIC)*, (Bahawalpur, Pakistan), pp. 1–6, IEEE, Nov. 2020.

[20] C. F. Rodríguez-Hernández, M. Musso, E. Kyndt, and E. Cascallar, "Artificial neural networks in academic performance prediction: Systematic implementation and predictor evaluation," *Computers and Education: Artificial Intelligence*, vol. 2, p. 100018, 2021.

[21] P. D. Saputri and P. P. Oktaviana, "Comparison of feedforward neural network and classical statistics methods: Application in finance," *Jurnal Matematika, Statistika dan Komputasi*, vol. 19, pp. 537–548, May 2023.

[22] D. Haryadi, A. R. Hakim, D. M. U. Atmaja, and S. N. Yutia, "Implementation of support vector regression for polkadot cryptocurrency price prediction," *JOIV : International Journal on Informatics Visualization*, vol. 6, p. 201, May 2022.

[23] B. A. Anisa Aulia, "Prediksi harga emas dengan menggunakan algoritma support vector regression (svr) dan linear regression (lr)," Apr. 2022.