



Spatio-temporal COVID-19 Spread Prediction: Comparing SVM with Time-Expanded Features and RNN Models

Raden Aria Gusti Aji*, Sri Suryani Prasetyowati, Yuliant Sibaroni

School of Informatics, Informatics, Telkom University, Bandung, Indonesia

Email: ^{1,*}ghyanqq@student.telkomuniversity.ac.id, ²srisuryani@telkomuniversity.ac.id, ³yuliant@telkomuniversity.ac.id

Correspondence Author Email: ghyanqq@student.telkomuniversity.ac.id

Submitted: 27/12/2024; Accepted: 26/02/2025; Published: 01/03/2025

Abstract—Covid-19 which spread in early 2020, still needs to be observed, considering the high growth rate of the pandemic at that time. The right prediction model is needed, because it can estimate the speed and extent of its spread for some time to come. This study develops a prediction model for the classification of the spread of Covid-19 in the future using SVM with time-based feature expansion and RNN. The scenario developed to determine the effect of time-based feature expansion and kernel function on classification performance using time series and spatial data. The results obtained show that SVM with time-based feature expansion achieves the most optimal performance using a polynomial kernel with an accuracy of 96.23%, a precision of 96.48%, a recall of 96.23%, and an F1-score of 96.21%. The performance of the SVM is superior to RNN which achieves an accuracy of 93.55%, a precision of 87.51%, a recall of 93.55%, and an F1-score of 90.43. Spatial prediction using Kriging interpolation can provide an overview of the spread of COVID-19 in all villages in Bandung City. The contribution of this research can provide much-needed information for policy makers and the community in managing future pandemic predictions and management strategies in the field of public health.

Keywords: COVID-19; Support Vector Machine; Recurrent Neural Network; Time-based Feature Expansion; Kriging Interpolation

1. INTRODUCTION

The COVID-19 pandemic has inspired various models based on machine learning and deep learning to model the spread of the virus. These models combine temporal and spatial dimensions to improve predictive performance. Most studies describe the trend or classification of COVID-19 cases. One of them proposed a deep learning model using Long Short-Term Memory (LSTM) and Convolutional Neural Network (CNN) networks to predict COVID-19 events, which has the potential to integrate temporal and spatial features in the prediction model [1].

Recurrent Neural Networks (RNNs) and their variants, such as Long Short-Term Memory (LSTM) models, have shown results in temporal forecasting of COVID-19 [2]. Recurrent Neural Networks (RNNs) and their variants, such as Long Short-Term Memory (LSTM) models, have shown results in temporal forecasting of COVID-19. A study conducted in Lampung Province, Indonesia, used an RNN-based model to predict active cases, recoveries, and deaths, achieving high accuracy with RMSE values as low as 0.001, 0.0027, and 0.001, respectively [3]. These results confirm that DL models are capable of handling time-series data with efficiency

Although the spatiotemporal dynamics of COVID-19 remains a challenging task, several studies have compared epidemiological models with ML models, where simpler models like SIR even outperformed complex models like SEIR highlighting the importance of effective feature engineering in the context of spatial data [4]. Similarly, a study advocated a probabilistic spatiotemporal neural network and estimated models with both spatial and temporal components, to improve the accuracy of COVID-19 case prediction [5].

Another study showed the prediction of cumulative confirmed cases, recoveries, and deaths trends over the next 60 days in top 10 countries using RNN with GRU and LSTM. A study presented the COVID-19 pandemic forecast for the next 60 days using RNN-LSTM and RNN-GRU models. For confirmed cases, the GRU model outperformed the LSTM model in countries like the US, Brazil, South Africa, Peru, Chile, and Iran, while for India, Russia, Mexico, and the UK, the growth of confirmed cases was predicted by the LSTM model and the GRU model provided a gradual improvement [6]. A study [7] using LSTM architecture in the US, India, and Italy predicted COVID-19 recovery cases for the next seven days. The performance was evaluated by MAPE, and the results showed that the error was small and never above 3%. Therefore, the deep learning model best fits the COVID-19 recovery dynamics and has made correct predictions with a total error of 1.46 to 2.65%.

Another study proposed the application of RNN and LSTM models to determine the level of Covid-19 transmission, obtaining optimum root mean square error values for the testing and training processes of 0.06 and 0.04, respectively. Although this study is efficient in predicting the pandemic, the study has not made predictions based on classification based on time [8]. In line with this, the study [9] used SVM in environmental modeling, the results obtained showed that SVM outperformed conventional geostatistical models with an accuracy rate of 95%. This study illustrates the ability of SVM to handle spatially distributed data and demonstrates its application for spatiotemporal modeling of COVID-19.

While [10] illustrates that the integration of Kriging with machine learning techniques, such as SVM, improves spatial data analysis learning by up to 25%. The study [11] presents that the integration of Kriging with machine learning techniques improves spatial data modeling, achieving an R^2 of 0.92. Similarly, urban flood prediction using machine learning was recently conducted in [12], whose model can achieve an F1 score of 0.89. This proves the relevance of the model for other types of environmental prediction tasks including pandemics. Hybrid models, such

as the integration of SVM with Kriging, have also shown better performance compared to single models on various spatial data [13,14]. While research [15] illustrates how deep learning combined with spatial features of extreme weather can further improve performance with values of over 94%. This suggests that the hybrid approach could be useful in spatiotemporal COVID-19 prediction. While research [16] explains that model selection should be done based on context, including data availability, type, and distribution, important considerations in spatiotemporal forecasting.

Based on the reviewed studies, it is found that there are not many studies comparing the performance of conventional ML models, such as SVM, with DL models, such as RNN, in the case of predicting spatial-temporal data classification, one example of which is COVID-19. This study conducts a comparative performance analysis of the proposed SVM model using time-based features against RNN in predicting the spread of COVID-19 infection in Bandung, Indonesia. Preliminary studies show that while DL models are promising for temporal modeling, SVM excels in spatial data modeling. This comparison will help in building a robust framework for mapping the spread of COVID-19 and provide valuable insights for policymakers to assist in future pandemic preparedness.

2. RESEARCH METHODOLOGY

2.1 Research Stages

The following research involves several steps in developing a spatio-temporal COVID-19 prediction system. Figure 1 shows the flowchart of how this system is designed.

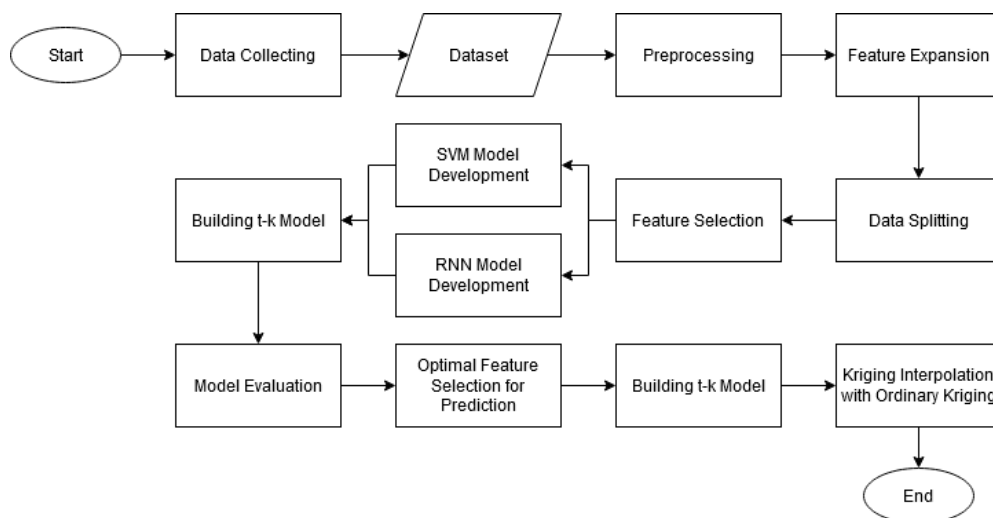


Figure 1. System Flowchart Design

The system initiates by gathering data about COVID-19 cases; each case has temporal information, such as time series, and spatial information, such as geographical. Preprocessing of the dataset was essential to ensure the quality by filling in missing values, normalizing numerical features, and encoding categorical variables. Temporal feature expansion then included generating lagged features in the form of case count history and time-based trends, such as moving averages, while spatial feature integration included mapping onto regional or coordinate-based information.

Following preprocessing and feature expansion, the dataset is then divided into a training set and a test set. In this regard, two predictive models were developed:

- a. SVM with time-expanded features: This will find the pattern in temporal data and handle class imbalance by oversampling techniques.
- b. RNN for time-series data to model sequential dependencies in COVID-19 case trends.

Every model follows an optimization of its respective hyperparameter in order to deliver the best result. The models developed are based on training sets, while the performance can be gauged based on testing sets. Major performance measures include Accuracy, Precision, Recall, F1-score, and MAE.

The outputs from these models are further processed to spatial predictions using techniques like spatial interpolation, such as Kriging. These result in visual spatio-temporal mappings of COVID-19 spreads across regions for actionable insights. The process flow will culminate in a comparative analysis to establish the suitability of either SVM or RNN modeling techniques for spatio-temporal prediction tasks.

2.1 Data Set

The data used is spatial-temporal data with 2,718 entries from November 1, 2020 to April 1, 2022. The temporal features observed daily include spatial features that describe geographic information. The features used describe the

factors that influence the increase in the number of Covid-19 cases, including demographic factors, environmental factors, vaccination rates, and human behavior.

2.2 Preprocessing

In this study, preprocessing of a number of steps was performed to prepare the dataset for model training. First, handling of missing values was done through imputation: numerical columns used mean, and categorical columns used mode imputation. Records containing incorrect values-for instance, temperature values in extremely wrong measures-were corrected or deleted to preserve the integrity of the information. The column representing date was standardized into a standard datetime format. Daily data were aggregated into monthly values to focus the analysis on medium-term trends.

The continuous numerical features are rainfall and temperature. Both of these were normalized between 0 and 1 using the *MinMaxScaler* to give equal importance in model training. Feature selection was then done by employing the *SelectKBest* method from *sklearn.feature_selection*, based on ANOVA F-test (f_{classif}), for the statistical selection of features. This step would have allowed the selection of the most relevant features and, hence, reduced the dimensionality of the dataset to improve the performance of the model.

Interaction terms were also generated, such as the product of rainfall and mask compliance, to see the nonlinear relationship between the variables. This imbalance in the dataset is addressed through oversampling with the SMOTE (Synthetic Minority Over-sampling Technique) method [19], a strategy that interjects synthetic samples by interpolative creation regarding the minority class. Therefore, this balanced out the dataset that the model needed with extra added classes that were underrepresented. Care was taken at this point to avoid overfitting-that which synthetic samples can do, driving the model to learn overly specific patterns of the minority class.

Finally, the dataset will be split into 80% for training and 20% for testing, considering the preservation of temporal order to avoid leakage in data. These steps made certain that the data would become clean, then normalized and balanced, and last but not least, structured adequately for both Support Vector Machine and Recurrent Neural Network training.

2.3 Support Vector Machine

Probably, the Support Vector Machine is one of the most famous algorithms which could be used both in classification and regression problems; besides, it is supervised. The key idea of this approach is to seek, in high-dimensional feature space, the best hyperplane which maximizes the margin between classes. This hyperplane is determined by support vectors-data points closest to it, having the greatest influence on its orientation and position [20].

While SVM can execute nonlinear boundaries of decisions thanks to a mathematical trick that goes by the name kernel trick, which projects the original features into a higher dimensional so that a linear separation becomes plausible there. Some of the most commonly used kernel functions are linear, polynomial, RBF, and sigmoid. Several recent works have been done to increase SVM's performance in different domains, for instance, new kernel designs with specific data distributions to improve the classification performance of the algorithm. Besides, new optimization techniques have improved scalability for enabling the application of SVM on large-scale datasets [21]. Because of its robustness, efficiency in high-dimensional spaces, and handling of both linear and nonlinear problems, SVM is considered one of the most popular and efficient machines learning methods.

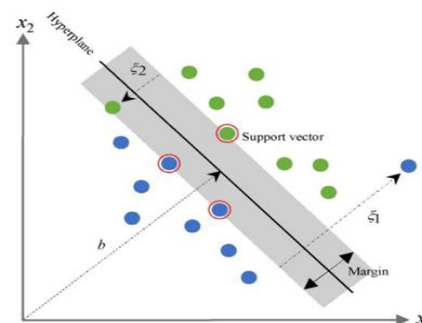


Figure 2. Illustration of support vector machine (SVM) to generalize the optimal separating hyperplane in linear separable data [22]

That is illustrated by Figure 2, where there is separation between two classes through a hyperplane. The hyperplane can also be referred to as a decision boundary-a geometrical way of separating classes that determines the optimum separation between the classes. The support vectors in SVM are the data points that lie closest to the hyperplane. These data points are critical in defining the margin, which is the distance between the hyperplane and the nearest support vector. The equation representing the hyperplane in SVM reads as follows:

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

where w represents the weight vector, x is the feature vector of the training data, and b is the bias term, a scalar. The kernel function in SVM plays an essential role by mapping the data from a lower-dimensional space to a higher-dimensional one, which allows for improved class separation. Various kernel functions are commonly used in SVM :

a. Linear Kernel

The linear kernel is the most straightforward type and is used when the data can be separated with a straight line or hyperplane. This kernel leads to a relatively simple hyperplane compared to other kernel types. Data that is linearly separable typically shows high consistency and is often characterized by a lower number of dimensions. The linear kernel function is represented as:

$$K(x, y) = x^T y \tag{2}$$

where the data vectors are involved. The linear kernel works by evaluating the training data points, x_i , with binary class labels under the following conditions:

$$H = \begin{cases} w^T x_i + b \geq +1 & \text{for } y_i = +1 \\ w^T x_i + b \leq -1 & \text{for } y_i = -1 \end{cases} \tag{3}$$

b. Polynomial Kernel

The polynomial kernel is designed for handling nonlinear data by mapping it into a higher-dimensional space through polynomial functions of a specified degree. It is particularly useful for data with complex, nonlinear patterns. The polynomial kernel function is expressed as :

$$K(x, y) = (x^T y + C)^d \tag{4}$$

where d represents the degree of the polynomial, C is a constant parameter, and x are the feature vectors. This kernel is effective in capturing and modeling nonlinear relationships within the data.

c. Radial Basis Function (RBF) Kernel

The Radial Basis Function (RBF) kernel, also known as the Gaussian kernel, is used when there is limited prior knowledge of the data distribution. It calculates the similarity between data points by employing a Gaussian function. The RBF kernel function is expressed as :

$$K(x, x') = \exp\left(-\frac{\|x-x'\|^2}{2\sigma^2}\right) \tag{5}$$

where x and x' are data vectors, and σ is a parameter that controls the width of the Gaussian function. The RBF kernel is highly adaptable and effective for handling complex data patterns, making it suitable for a wide range of applications.

d. Sigmoid Kernel

The sigmoid kernel, although typically less effective than the RBF kernel, can be optimized to deliver performance similar to other kernel methods. The sigmoid kernel function is given by :

$$K(x, y) = \tanh(\alpha x^T y + c) \tag{7}$$

where α is a scaling parameter, c is a bias term, and x and y are data vectors. The sigmoid kernel is commonly applied in cases where the decision boundary is nonlinear in a higher-dimensional space.

2.4 Recurrent Neural Network (RNN)

RNNs are one class of neural networks that come in when one is dealing with sequential data. Unlike feed-forward neural networks, the RNNs do not require independence assumptions strict enough across the inputs; rather, they remember every previous input by creating an internal memory. Because of this, the feedback connections within a hidden layer of an RNN really enable a network to process a series of data by taking information from one time step to another.

The major idea of RNNs is to use the same weights throughout all steps in time. This will make them learn the pattern across time, which captures temporal dependencies in data. In time series analysis, it is applied for the prediction of values in the future based on past data. The typical architecture includes an input layer, a recurrent hidden layer, and an output layer. It takes the input at every time step, does some processing, and passes on the result as input back in the same layer at the next time step.

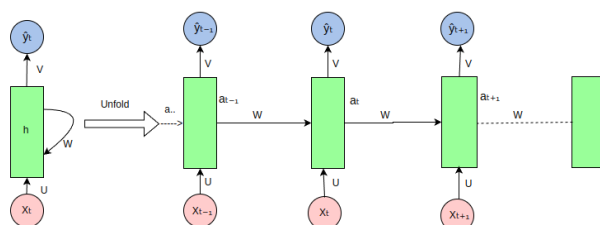


Figure 3. Illustration of a simple RNN architecture



That is illustrated by Figure 3, the input sequence at each time step is fed through the recurrent layer, which outputs the hidden state. This hidden state will be used in order to predict the output at that time step but also feeds back into the network for the next step. RNNs are quite naturally suited for such tasks as time-series forecasting, language modeling, and sequence classification, since they capture the dependencies between past, present, and future time steps.

3. RESULT AND DISCUSSION

The discussion here is based on the experimental results of the models created. Constructed test scenarios are made by this research for each existing model in order to analyze the performance and effectiveness of each model. This study will be using SVM and RNN modeling. Further explanations will be presented in the following sections about the results of the tests and their analysis.

3.1 Data Attributes

These attributes have been collected from various places at different times to analyze spatiotemporal dependencies of demographic, environmental, and behavioral factors on the severity of COVID-19 spread.. Below is a detailed description of the data used :

Table 1. Data Attribute

Attributes	Description
X ₁	Total male population in the area.
X ₂	Total female population in the area.
X ₃	Rainfall amount in millimeters (mm).
X ₄	Solar radiation level (W/m ²).
X ₅	Average temperature in Celsius (°C).
X ₆	Maximum recorded temperature (°C).
X ₇	Minimum recorded temperature (°C).
X ₈	Number of individuals with no formal education.
X ₉	Number of individuals who did not complete elementary school.
X ₁₀	Number of individuals who completed elementary school.
X ₁₁	Number of individuals who completed middle school.
X ₁₂	Number of individuals who completed high school.
X ₁₃	Number of individuals with Diploma I or II education.
X ₁₄	Number of individuals with Diploma III education.
X ₁₅	Number of individuals with undergraduate education and Diploma IV.
X ₁₆	Number of individuals with master's degree education.
X ₁₇	Number of individuals with doctoral degree education.
X ₁₈	Number of the population who received the first dose of COVID-19 vaccine.
X ₁₉	Number of the population who received the second dose of COVID-19 vaccine.
X ₂₀	Number of the population who received the third dose of COVID-19 vaccine.
X ₂₁	Average compliance rate of wearing masks.
X ₂₂	Average compliance rate of maintaining physical distance.
Y	Target variable: class label representing COVID-19 spread severity.

The target class labeling for the number of Covid-19 cases in classified into 3 levels, with descriptions explained in Table 2.

Table 2. The Target Class Labling For The Number Of Covid-19 Cases

Class	Range	Label
Low	Cases < 218	0
Medium	218 ≤ Cases < 419	1
High	Cases ≥ 419	2

3.2 Feature Expansion

Feature expansion that includes time-based features is done through previous time steps that provide additional data input in the dataset. This method allows the model to learn the dependence on temporal patterns to predict future trends in time series data. The addition of time-based features, which is done regularly from T-2 to T-17, as described in Table 3.

Scenarios for controlling the sensitivity of performance models to various configurations, can consider the dimension size with the progressive addition and removal of lagging features [18]. The optimal feature expansion

process allows the model to expand the most relevant historical information, thereby improving the performance of time-based prediction models using SVM and RNN.

Table 3. Data Time-Based Expansion Feature

Time-based feature expansions Scenario	Combination	Training Data Period	Target
T-2	1	November 2020 – December 2020	January 2021
	2	December 2020 – January 2021	February 2021
	16	February 2022 – March 2022	April 2022
T-3	1	November 2020 – January 2021	February 2021
	2	December 2020 – February 2021	March 2021
	15	January 2022 – March 2022	April 2022
T-16	1	November 2020 – February 2022	March 2022
	2	December 2020 – March 2022	April 2022
T-17	1	November 2020 – March 2022	April 2022

3.3 Result

In this subsection, after the pre-defined modeling stages, the next will be evaluation. The performance evaluation of the SVM and RNN algorithm models will be performed using pre-defined test data ratios. In-depth analysis is performed with the help of an inclusive evaluation matrix that includes but is not limited to metrics such as accuracy, F1-score, precision, and recall for deducing the model's efficiency on various test scenarios and datasets.

These are weighted averages of metrics that include the F1-score, recall, accuracy, and precision through which conclusions will be drawn about evaluation. The exact insight into the performance level under different testing conditions and on diverse datasets is thus drawn.

The following table shows how variations in the parameters of C, Gamma, and Kernel have performed in support vector machine modeling to predict the severity of COVID-19 spread using the features of temporal lag. The general trend of this model is very stable, with high accuracy and a high F1-score observed throughout most of the experiments, especially for those lags between T-5 and T-8. In these cases, accuracy and F1-Score were as high as 96%, while the best results came out with a polynomial kernel: poly with C = 100 and Gamma = 0.9. This indeed suggests that with the right choice of parameters, the SVM model was able to capture nonlinear relationships in the data efficiently. However, as more lag features were added, performance started degrading, especially beyond T-8-for instance, going from T-9 to T-13, the accuracy and F1-Score dropped in the range of 93-94%. This performance reduction might have been a case of overfitting, whereby the model started capturing an irrelevant/noisy pattern because of extra lag features.

Table 4. Best T-K SVM Time-Based Model

Lag Expansion	Parameters C	Gamma	Kernel	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)	Feature Count
T-2	100	0.4	rbf	94.3396	94.8874	94.3396	94.3029	27
T-3	100	0.9	poly	94.3396	94.8874	94.3396	94.3029	39
T-4	1	0.9	sigmoid	94.3396	94.8874	94.3396	94.3029	7
T-5	100	0.9	poly	96.2264	96.478	96.2264	96.2129	73
T-6	10	0.1	linear	943396	94.8874	94.3396	94.3029	72
T-7	100	0.9	poly	96.2264	96.478	96.2264	96.2129	99
T-8	100	0.9	poly	96.2264	96.478	96.2264	96.2129	108
T-9	0.1	0.1	linear	93.5484	87.513	93.5484	90.4301	3
T-10	0.1	0.1	linear	93.5484	87.513	93.5484	90.4301	3
T-11	0.1	0.1	linear	93.5484	87.513	93.5484	90.4301	3
T-12	0.1	0.1	linear	93.5484	87.513	93.5484	90.4301	3
T-13	0.1	0.1	linear	93.5484	87.513	93.5484	90.4301	3
T-14	1	0.8	poly	93.5484	93.5484	93.5484	93.5484	251
T-15	1	0.9	poly	93.5484	93.6111	93.5484	93.4693	244
T-16	10	0.3	poly	93.5484	93.5484	93.5484	93.5484	279
T-17	1	0.5	poly	93.9394	93.9414	93.9394	93.8303	250

From the results obtained in Table 4, it can be stated that fewer features are related to better performance. For example, the experiment with T-2, which used only 27 features, performed very well; it achieved an accuracy of 94.34% with a very strong F1-Score. Contrariwise, experiments conducted by using a higher number of features-for example, T-14 to T-17, using more than 250 features-did not enhance accuracy and actually resulted in performance decrease. This would support the idea that the addition of a large number of features-especially when these features

capture very long lags-results in overfitting because the added dimensions introduce complexity that adds nothing to generalization. This goes with conventional wisdom: simpler models are often more robust, suffering less from overfitting, especially in time-series, which may well not benefit from long lags in the capture of its underlying patterns.

For the Gamma parameter, the higher values, like Gamma = 0.9 and 1, were performing better for all instances than the smaller ones, such as Gamma = 0.1. That means the model of SVM benefits from a wider scope of influence on the data since higher Gamma values allow the kernel to capture more complex nonlinear patterns. Other than the kernel parameter, this model heavily relies on the C parameter. As such, it was noticed that with higher values of C-for example, C = 100-the performance became high, while for smaller values-C = 0.1-performance drastically decreased. It is believed that this model works better when there is stronger regularization to prevent overfitting and ensure generalization of the SVM to unseen data.

Also, the linear kernel used in experiments for lags T-9 to T-13 gave poor performance with respect to the polynomial and RBF kernels. Again, this enhances the fact that there are non-linear relationships in the data that are captured by more flexible kernels like polynomial and RBF rather than the simple linear kernel. The following diagram (Figure 4) illustrates these findings by showing the performance of the SVM model across different time-based lag expansions (T-K), highlighting the effects of varying parameters such as C, Gamma, and Kernel on accuracy, precision, recall, and F1-score.

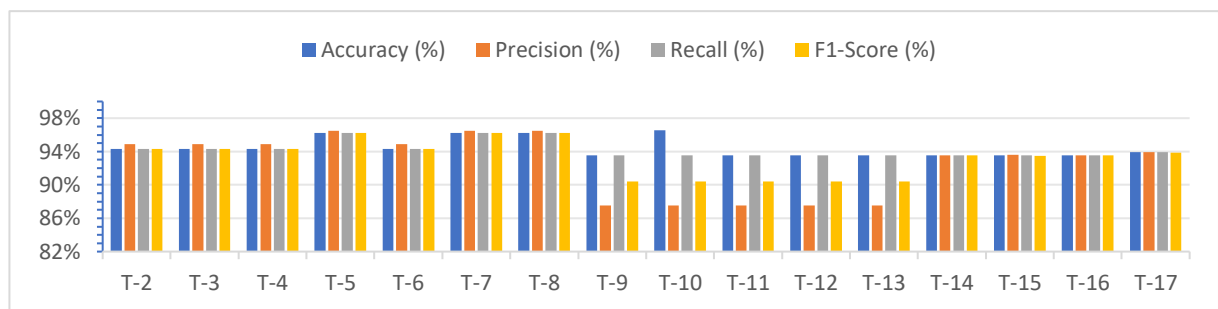


Figure 4. Best T-K SVM Time-Based Model

In Figure 4, these results stress the importance of choosing appropriate parameters and lag features for better optimization of the model. While increasing the number of lag features can be very tempting, less is often more in reality, especially when focusing on shorter lags such as T-5 to T-8, which prevents overfitting and may increase generalizability. Thus, careful feature selection and tuning of parameters are necessary for the maximum performance of the SVM model in time-series prediction tasks.

Results are shown in the following table for the performance of the RNN model trained on various lag expansions ranging from T-2 to T-17. In fact, the first 16 lag experiments, which range between T-2 and T-16, provided the same performance for the RNN model with an accuracy of 93.55%, precision of 87.51%, recall of 93.55%, and F1-Score of 90.43%. This means that from T-2 to T-16 lag features, the performance of the model is level, and its predictive capability shows no considerable improvement or degradation since the number of selected lag features is constant-3 features. This, therefore, supports the evidence that RNN can identify temporal dependencies within these lags such that the added features have no influence on improving the performance. This stability across these different lag values shows the efficiency of the model in using the available features to capture relevant patterns in the time-series data

Table 5. Best T-K RNN Time-Based Model

Lag Expansion	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)	Feature Count
T-2	93.5484	87.513	93.5484	90.4301	3
T-3	93.5484	87.513	93.5484	90.4301	3
T-4	93.5484	87.513	93.5484	90.4301	3
T-5	93.5484	87.513	93.5484	90.4301	3
T-6	93.5484	87.513	93.5484	90.4301	3
T-7	93.5484	87.513	93.5484	90.4301	3
T-8	93.5484	87.513	93.5484	90.4301	3
T-9	93.5484	87.513	93.5484	90.4301	3
T-10	93.5484	87.513	93.5484	90.4301	3
T-11	93.5484	87.513	93.5484	90.4301	3
T-12	93.5484	87.513	93.5484	90.4301	3
T-13	93.5484	87.513	93.5484	90.4301	3
T-14	93.5484	87.513	93.5484	90.4301	3
T-15	93.5484	87.513	93.5484	90.4301	3



T-16	93.5484	87.513	93.5484	90.4301	3
T-17	78.7879	84.4066	78.7879	77.5684	83

However in Table 5, a notable drop in performance is seen at T-17, where the accuracy drops to 78.79%, and the F1-Score goes down to 77.57%. Despite this drop, precision and recall still show relatively good values, with precision at 84.41% and recall at 78.79%. That may be an indication that such a significant degradation of performance compared to T-17 suggests that the introduction of even more noise or redundant information by adding 83 features overfits the data, hence showing that returns can diminish for an RNN model, much as for an SVM, at larger numbers of lagging features. Beyond a certain point, increasing the number of features further might not only fail to enhance the generalization capability of the model but can also lead to poorer performance due to the added complexity introduced by the additional features.

While the feature set is more optimized for the SVM model, the RNN seems to take a significant drop once the count of features gets excessive, especially for T-17. This supports the fact that sometimes less may be more, given that for time-series data, often most of the relevant patterns could be captured with a much-constrained set of features. It is, therefore, crystal clear that the RNN models, similar to those in this particular task, can benefit from feature selection strategies focused on reducing the dimensionality of features in order not to be subject to overfitting.

The following diagram (Figure 5) visually illustrates these findings, showing the performance of the RNN model across different time-based lag expansions and highlighting the impact of varying the feature count on the model's accuracy, precision, recall, and F1-score.

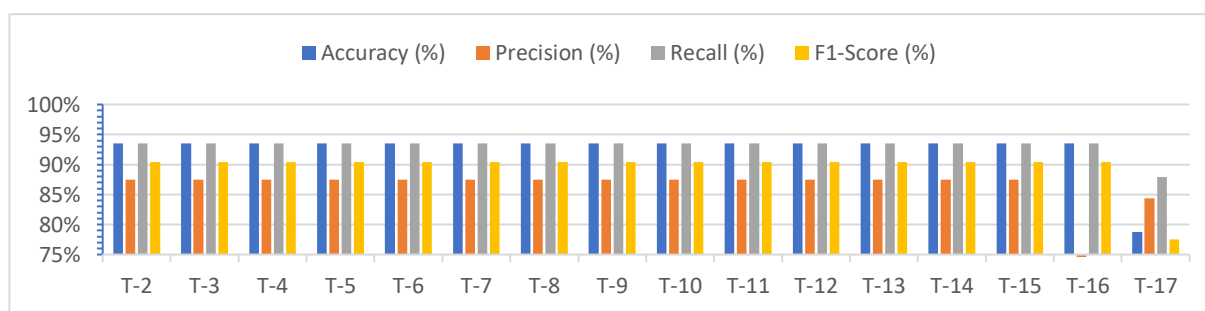


Figure 5. Best T-K RNN Time-Based Model

Overall, the results show that lag expansions from T-2 to T-16 are optimal for the RNN model, as this setting yields stable performance with the least complexity. The introduction of more features, especially at T-17, seems to disturb the generalization capability of the model. Again, this points out feature selection and the risk of overfitting when using more complex temporal models such as RNNs. Therefore, a number of fewer lag features should be used for the optimal performance of RNN.

The prediction results for COVID-19 Spread on the latest date using the best time-based SVM and RNN models combined with Kriging interpolation are presented in the following tables. The model's performance is evaluated across various prediction intervals, from T+2 up to T+17, where each interval represents a day beyond the initial prediction date. The predictions classify the expected COVID-19 spread severity into different categories based on the total case value ranges, as outlined in Table 1.

Table 5 shows the predictions from T+2 to T+17 of SVM models, while Table 6 shows the predictions from T+2 to T+17 of RNN models. The numerical labels represent categorical COVID-19 spread severity, where 0 indicates "Low," 1 indicates "Medium," and 2 indicates "High".

Table 6. Best T+K SVM Time-Based Model

Location	T+2	T+3	T+4	T+5	T+6	T+7	T+8	T+9	T+10	T+11	T+12	T+13	T+14	T+15	T+16	T+17
Ancol	2	2	2	1	2	2	2	0	0	0	0	0	0	0	0	2
Antapani Kidul	2	2	2	1	2	2	2	0	0	2	0	2	2	2	2	2
Antapani Kulon	2	1	2	2	2	2	2	0	0	0	0	0	0	0	0	1
Antapani Tengah	2	2	2	1	2	2	2	0	0	0	0	2	1	0	1	2
Antapani Wetan	2	2	2	1	2	2	2	0	0	1	0	0	0	1	1	2
....																
....																
Tamansari	2	2	2	2	2	2	2	0	0	2	1	2	0	0	0	2
Turangga	2	1	2	2	2	2	2	0	0	2	0	2	0	0	0	2

Warung	2	1	2	2	2	2	2	0	0	0	0	0	0	1	1	1
Muncang																
Wates	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0

Table 7. Best T+K RNN Time-Based Model

Location	T+2	T+3	T+4	T+5	T+6	T+7	T+8	T+9	T+10	T+11	T+12	T+13	T+14	T+15	T+16	T+17
Ancol	1	2	2	2	2	2	1	2	2	2	1	2	2	1	0	2
Antapani Kidul	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
Antapani Kulon	1	2	1	2	1	2	1	2	2	2	1	2	2	1	0	2
Antapani Tengah	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
Antapani Wetan	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
....																
Tamansari	1	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
Turangga	2	2	2	2	2	2	1	2	2	2	2	2	2	2	2	2
Warung Muncang	0	2	2	1	1	2	1	1	1	0	1	0	0	0	1	1
Wates	0	1	1	0	0	1	1	0	0	0	0	0	0	0	0	0

After that on Table 7, the predicted severity of COVID-19 transmission by SVM and RNN with time-based feature expansion was interpolated by Kriging to generate a spatial prediction map. Kriging is a powerful geostatistical interpolation technique, estimating the values at unsampled locations by exploiting the spatial autocorrelation inherent in the sampled data points. That brings up a particular application in the case of spatial data analysis: not only accounting for values known at sampled locations, but also taking into consideration the spatial dependencies between them.

The SVM and RNN model performed the forecast on COVID-19 severity based on temporal features. These included trends in the past vaccination rates, compliance of mask-wearing, along with environmental factors such as rainfall and temperature. These predictions were computed at specific village locations, known as Desa in this dataset, each associated with a set of temporal features. In order to get a better understanding of the spread and possible future hotspots of COVID-19, these predictions needed to be extended to unsampled locations in the study area.

Kriging will allow the model to take these predictions at specific sampled locations of the villages and interpolate spatially the predictions of COVID-19 severity levels across the whole geographic region, even in areas where direct measurements were not available. Because kriging uses spatial relationships among known data points, this method effectively leverages proximity and similarity among neighboring villages to predict severity in unsampled locations. The interpolation will consider the short-range and long-range spatial dependencies, allowing a far better and close-to-reality mapping of the spread of COVID-19 in the region.

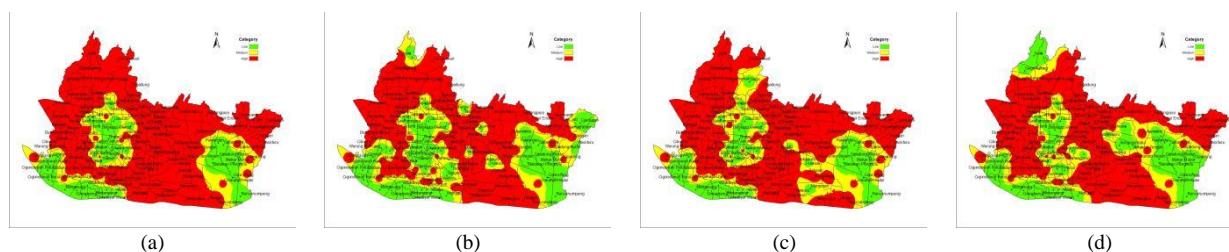


Figure 6. SVM Models Prediction Map of COVID-19 severity Distribution of May 2022 to August 2022

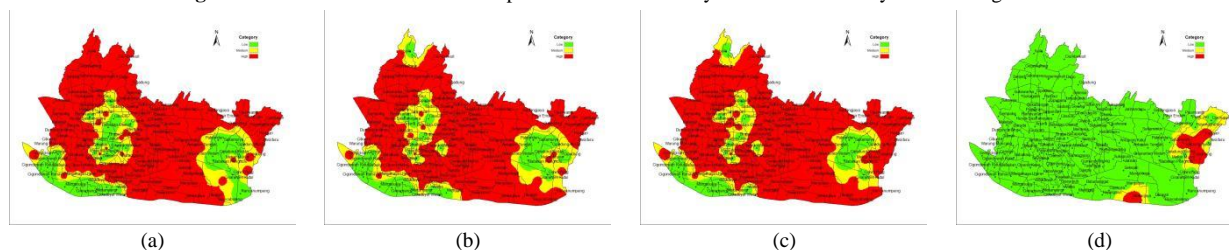


Figure 7. SVM Models Prediction Map of COVID-19 severity Distribution of September 2022 to December 2022

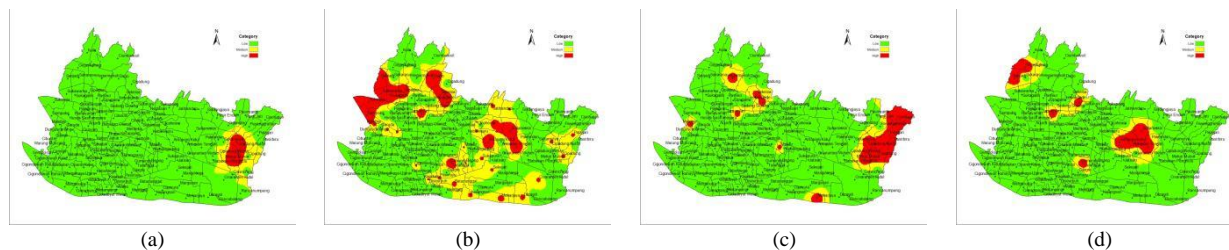


Figure 8. SVM Models Prediction Map of COVID-19 severity Distribution of January 2023 to April 2023

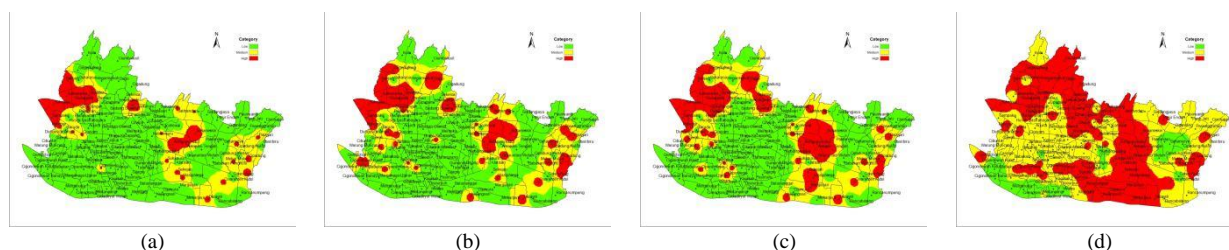


Figure 9. SVM Models Prediction Map of COVID-19 severity Distribution of May 2023 to August 2023

Time-varying COVID-19 severity over the observed regions has been shown in Figures 6 through 9, presenting the spatial visualizations for the performance of the SVM model over the region of interest regarding the predicted results. These maps use colors red for high, yellow for moderate, and green for low. Based on these, the paper can then outline such spatial and temporal trends in COVID-19 severity in order to provide insights about the progress in dealing with the pandemic.

Figure 6 presents the forecasted severity distribution from May 2022 to August 2022. In fact, most areas of the region are enveloped in red color as the majority depict high severity of COVID-19 in that period, especially in the middle and western parts of the region; peripheral areas of the study region show a trend of average severity, and only scattered small pockets reflect low severity. This distribution makes it visible that, throughout the year, the pandemic reached a critical point or even more so in central areas due to high population, high mobility, or, quite on the contrary, for their low health capacity.

Figure 7 presents the forecasted severity distribution from September 2022 to December 2022. Note how much the high-severity areas have shrunk in size compared to the previous period, while spreading occurred in the areas of moderate severity shown in yellow coloration. Perhaps that reflects improvements in the management of the pandemic—whether via enforcement of public health measures, improving vaccination rates, or more efficient protocols related to treatment. On the contrary, notice the few areas in green: that represents how large-scale challenges still remain amid the path to broad control.

Figure 8 shows the forecasted severity distribution from January 2023 to April 2023. An obvious improvement can be traced from this period whereby the green color, which represents the low-severity area of COVID-19, has grown a lot while the red zones that showed high severity more or less have vanished along with shrinkage in the yellow-colored areas representing moderate severity. It could be indicative of successful interventions accumulated in terms of vaccination drives, the health guidelines followed by the lay public, and general improvement of health infrastructure. In that period, the spatial pattern of distribution demarcated the transition towards recovery although some regions were still striving hard to fight the disease with full vigor.

Figure 9 presents the forecasted severity distribution between May 2023 and August 2023. For this period, one could notice the tremendous fall in COVID-19 severity, as the maximum region is dominated by a green color which reflects the low level of the disease. Small pieces of yellow areas would still be isolated to indicate moderate severity but will remain infrequent and very confined with regards to their spread. The general trend of green does indeed indicate a near containment of the pandemic, with only isolated, sporadic outbreaks constrained to areas with specific vulnerabilities, such as lower vaccination rates or delayed healthcare responses.

Figure 6 to Figure 9 give an indication of the strength of the SVM model in reflecting the dynamics of COVID-19 severity, both in the temporal and spatial trends. Temporally, the maps trace a clear path of improvement from high-severity conditions throughout the region to largely low-severity outcomes. That this is consistent across several periods gives an indication of the reliability of the model in monitoring and predicting disease trends. Such predictions, on a broader scale, bring much-needed insight into public health planning and decision-making. Gradual decline in the high-severity areas, coupled with incrementation in the low-severity areas, indicates success of continued health intervention. However, regional persistence of moderate-severity pockets may indicate the need for targeted strategies regarding vaccination drives, improvement in health care access in those areas, and also community-based public health campaigns targeting the same areas.

Additionally, some further variables—population density, mobility pattern, socio-economic attributes, capacity of health care—can further enhance the predicting capability of the SVM model. Once integrated into the model, such variables make even more enlightened predictions, hence helping policymakers propose proactive, evidence-based

decisions on public health. The ability to visualize these trends within a spatial perspective further underlines the role and importance of machine learning models like SVM in epidemiological research and the management of pandemics.

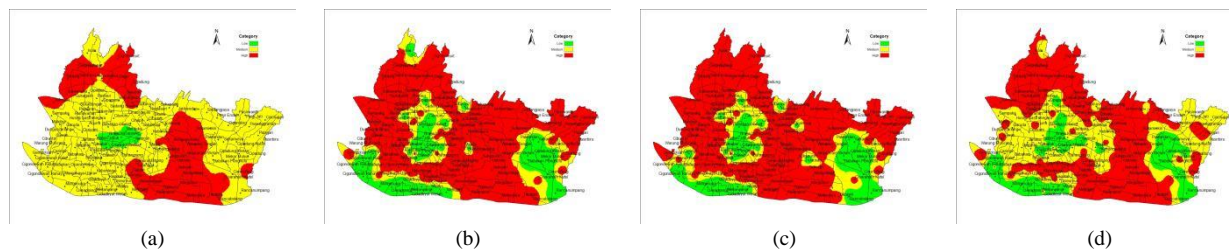


Figure 10. RNN Models Prediction Map of COVID-19 severity Distribution of May 2022 to August 2022

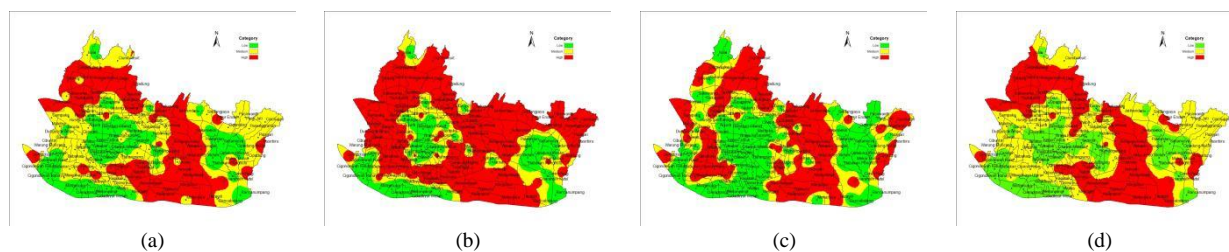


Figure 11. RNN Models Prediction Map of COVID-19 severity Distribution of September 2022 to December 2022

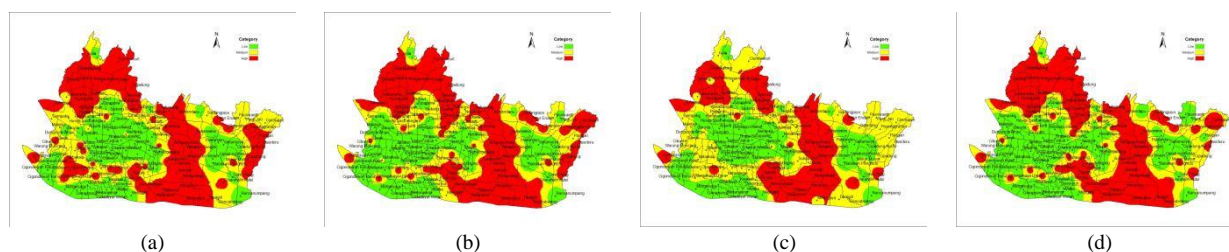


Figure 12. RNN Models Prediction Map of COVID-19 severity Distribution of January 2023 to April 2023

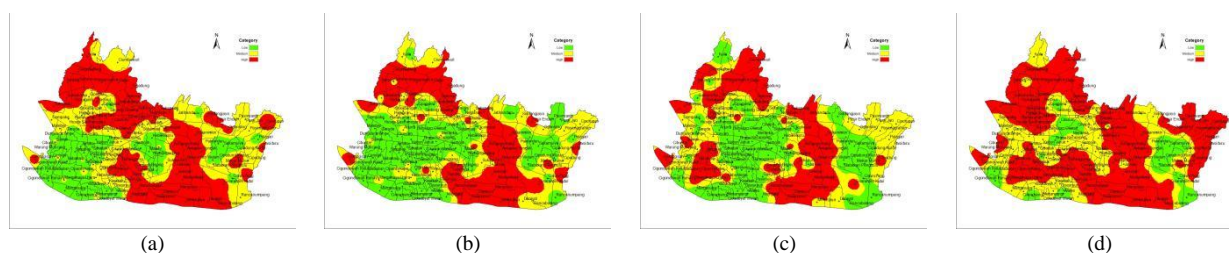


Figure 13. RNN Models Prediction Map of COVID-19 severity Distribution of May 2023 to August 2023

Figures 10 to Figure 13 display the RNN model-predicted spatial distribution of COVID-19 severity in the observed region for different time periods. Color code is set as red for high, yellow for medium, and green for low severity, respectively, as was used in the previous maps. These maps provide the detailed development of the severity of the pandemic in time and space, reflecting the strengths of the RNN model in the forecast of trends based on sequential data.

Figure 10 shows the forecasted severity distribution from May 2022 to August 2022. This map represents the initial phase of the prediction where high severity areas, in red color, dominate central and some western parts of the region. The yellow-colored zones represent the areas of moderate severity around these red areas, while the low-severity regions in green color are relatively sparse and mainly lie at the periphery. This pattern then indicates that the pandemic is at a critical phase. The central regions are likely to be the epicenters due to higher population density or significant mobility activities. These findings highlight the necessity for immediate public health interventions-such as mobility restrictions and vaccination efforts.

Figure 11 shows the forecasted severity distribution from September 2022 to December 2022. During this period, there is a gradual decrease in high-severity areas: the red zones shrink and the yellow regions expand. Simultaneously, green areas become more and more prominent, especially in the peripheral regions. This trend could indicate that interventions during the earlier months might be yielding positive results. While this shift to moderate

severity indicates the sustained difficulty of mitigating the effects of the pandemic, the resulting low-severity areas also suggest partial recovery in parts of the region.

Figure 12 shows the forecasted severity distribution from January 2023 to April 2023. This period is quite remarkable in the improvement of the severity distribution, with the low severity in green dominating most of the map. Whereas high-severity zones are now confined to very few isolated hotspots, the areas of moderate severity spill all over the map in yellow. The vast expanses of green mark containment of the pandemic in most of the areas, probably through cumulative efforts like vaccination, improved public health, and better health infrastructure. But small, persisting red zones hint at the requirement for targeted interventions in these specific areas.

Figure 13 shows the forecasted distribution of severity from May 2023 to August 2023. The last predicted period for COVID-19 severity is almost totally contained, as depicted in green zones for almost the whole map, with yellow, indicating moderate severity areas, sparsely and very limited, and high severity areas in red color, which will be almost zero. This broad-based decline in severity reflects the success of prolonged public health strategies and, in some areas, possibly the attainment of herd immunity. The map now reflects a clear transition toward post-pandemic recovery, with only small areas needing further monitoring and intervention.

3.4 Discussion

These results present how well SVM and RNN, with their time-based feature extension, are capable of predicting the severity of COVID-19-in fact, even better in tandem with the Ordinary Kriging method for spatial interpolation. The discussion of the results and what could be gleaned from them follows.

Results of the SVM model show that much importance is given to the performance pre-processing of parameter tuning and feature selection. The accuracy, precision, recall, and F1-score obtained using different lag expansions and parameter variations are underlined in Table 3. It has been observed that the experiments involving lag expansions from T-5 to T-8 showed the best performance, with accuracy and F1-scores as high as 96%. Among the considered kernel types, the polynomial kernel with $C = 100$ and $\text{Gamma} = 0.9$ has constantly shown the best performance in capturing nonlinear relationships within the data. These results put into evidence the importance of choosing flexible kernels that are able to adapt to the complexity of time-series data.

However, as the number of lag features increased beyond T-8 to include, say, T-9 to T-13, model performance started to decay: the accuracy and F1-score declined to a low of 93–94%. This is due to overfitting whereby too many feature dimensions inducted noise, rather than information. However, some of the simpler models with few features had a great performance; for instance, the T-2 with 27 features showed great promise, with an accuracy of 94.34% and a good F1-score. These results have shown that the models run with small Gamma values, such as $\text{Gamma} = 0.1$, perform very poorly compared to the higher values, for instance, $\text{Gamma} = 0.9$ or 1.0 , since in this case, the model might be allowed to learn complicated patterns. On the contrary, large value of C , say, $C = 100$, tends to impose heavy regularization and hence avoids overfitting, generalizing better. Especially outstanding was the mediocre result for the linear kernel, which tended to perform even worse in experiments using lags T-9 to T-13-thus underlining the impossibility of a nonlinear relationship fitting by linear kernels-and the relatively high performance of the polynomial and RBF kernels.

The overall results from the SVM confirm that proper parameter tuning and feature selection are key. Models with moderately lagging features, such as T-5 to T-8, yielded an optimal performance that balanced model complexity and generalization. Again, this supports the paradigm that simpler, well-regularized models are often far more robust, especially in time-series prediction tasks.

The results for the RNN model are presented in Table 4. For the lag expansion from T-2 to T-16, the model's performance was quite consistent: an accuracy of 93.55%, with precision, recall, and F1-scores all hovering within similar levels at 87.51%, 93.55%, and 90.43%, respectively. It seems that the RNN learns temporal dependencies from the data quite well, even when only a small number of features (3 lag features) are used in this experiment. With more added lag features at this granularity, there was a performance drop with accuracy decreasing to 78.79% at T-17 while the F1-score falls down to 77.57%. This reflects the addition of noise, poor generalization performance, and consequently the degrading performance of the overall model. Although degradation happened, precision and recall values were pretty good, being 84.41% and 78.79%, respectively, suggesting that there is a reasonable level of predictive capability present in a model.

Most importantly, the stability of the RNN model, tested across different values of lag features from T-2 to T-16, underlines its robustness for handling time-series data without being overly sensitive to the inclusion of extra lag features. However, the notable drop at T-17 outlines the importance of keeping feature increase limited to avoid introducing irrelevant patterns. Comparing the approaches by the two models, solutions were reasonably good given some optimal constraints. The SVM model performed best at T-5 to T-8 with its flexibility in parameters, enabling it to pick nonlinear patterns, while the RNN model is more stable when using more extensive lag expansion and has poor capacity to benefit from more features than T-16. While the SVM was considerably dependent on parameter tuning, such as Gamma and C for optimized performance, the RNN, by its very design, handled temporal dependencies with very few changes in parameters. However, this performance degradation in both models upon introducing excessive lag features points out that feature selection should be performed with caution.

Their strong points are different SVM gives outstanding performance when the optimization of the parameters should be precise, while RNN is robust against different sets of features. From a time-series prediction perspective, it

would appear that results could indicate an optimal performance that might be achieved with a balance between model complexity and appropriate feature selection, besides proper tuning of parameters.

4. CONCLUSION

Based on the results and discussions, the best performance of SVM with time-based feature expansion occurs in models T-5, T-7, and T-8, with polynomial kernels. The performance obtained is respectively with accuracy, 96.23%, precision, 96.48%, recall, 96.23%; and F1 score, 96.21%. This shows that the model is able to understand nonlinear temporal dependencies. Meanwhile, the best performance on RNN occurs at lags T-2 to T-15, with accuracy of 93.55%, precision of 87.51%, recall of 93.55%, and F1 score of 90.43%. Despite having far fewer features than usual, RNN can learn temporal dependencies with sequence learning without much explicit feature engineering. However, performance is seen to decrease drastically by lag T-17 at accuracy of 78.78% and F1 score of 77.57%. These results show that SVM with time-based feature expansion outperforms RNN and machine learning methods without time-based feature expansion.

REFERENCES

- [1] M. Muñoz-Organero and P. Queipo-Álvarez, "Deep Spatiotemporal Model for COVID-19 Forecasting," *Sensors*, vol. 22, no. 9, p. 3519, 2022, doi: 10.3390/s22093519.
- [2] B. Vahedi, M. Karimzadeh, and H. Zoraghein, "Spatiotemporal prediction of COVID-19 cases using inter- and intra-county proxies of human interactions," *Nat. Commun.*, vol. 12, p. 6440, 2021, doi: 10.1038/s41467-021-26742-6.
- [3] A. Ghozi, A. Aprianti, A. Dimas, and R. Fauzi, "Analisis Prediksi Data Kasus Covid-19 di Provinsi Lampung Menggunakan Recurrent Neural Network (RNN)," *Indonesian Journal of Applied Mathematics*, vol. 2, no. 1, pp. 25-32, 2022, doi: 10.35472/indojam.v2i1.763.
- [4] A. Kaddar, A. Abta, and H. T. T. Alaoui, "A comparison of delayed SIR and SEIR epidemic models," *NAMC*, vol. 16, no. 2, pp. 181–190, Apr. 2011, doi: 10.15388/NA.16.2.14104.
- [5] F. Ravenda, M. Cesarini, S. Peluso, et al., "A probabilistic spatio-temporal neural network to forecast COVID-19 counts," *Int. J. Data Sci. Anal.*, 2024, doi: 10.1007/s41060-024-00525-w.
- [6] K. E. ArunKumar, D. V. Kalaga, C. M. Sai Komar, M. Kawaji, and T. M. Breza, "Forecasting of COVID-19 using deep layer Recurrent Neural Networks (RNNs) with Gated Recurrent Units (GRUs) and Long Short-Term Memory (LSTM) cells," *Chaos, Solitons & Fractals*, vol. 146, 2021, doi: 10.1016/j.chaos.2021.110861.
- [7] B. Safa and M. Kdayem, "Deep Learning for COVID-19 prediction," in *Proc. 2020 Int. Conf. on Advances in Science, Engineering and Technology (IC_ASET)*, 2020, pp. 406-411, doi: 10.1109/IC_ASET49463.2020.9318297.
- [8] R. K. Pathan, M. Biswas, and M. U. Khandaker, "Time series prediction of COVID-19 by mutation rate analysis using recurrent neural network-based LSTM model," *Chaos, Solitons & Fractals*, vol. 138, p. 110018, 2020, doi: 10.1016/j.chaos.2020.110018.
- [9] Z. Ahmed and M. K. Faisal, "Advancements in Support Vector Machines for Environmental Modeling," *J. Environ. Informatics*, vol. 47, no. 1, pp. 40–55, 2021, doi: 10.3808/jei.2021.47.1.40.
- [10] E. Zhao and P. Tan, "Integration of Machine Learning and Kriging for Spatial Data Analysis: A Review," *Comput. Geosci.*, vol. 151, pp. 104–117, 2021, doi: 10.1016/j.cageo.2021.104117.
- [11] Y. Wang, M. Xu, and X. Zhao, "Comparative analysis of kriging interpolation and support vector machine regression for spatial data," *Environ. Earth Sci.*, vol. 78, no. 3, p. 81, 2019.
- [12] J. Park, J. Kim, and S. Lee, "Application of machine learning techniques for flood prediction in urban areas: A review," *Water (Basel)*, vol. 12, no. 1, p. 352, 2020, doi.org/10.3390/w12010352
- [13] N. Ali, M. B. Rahman, and H. U. Ahmed, "Hybrid Models of SVM and Kriging for Improved Spatial Predictions," *Environ. Monit. Assess.*, vol. 196, no. 9, pp. 1–15, 2024, doi: 10.1007/s10661-024-11628-0.
- [14] J. Sun, T. Wu, and Q. Li, "Recent Developments in Kriging Methods for Accurate Spatial Data Modeling," *Spat. Stat.*, vol. 40, p. 100546, 2024, doi: 10.1016/j.spasta.2024.100546.
- [15] L. Cheng, Y. Zhang, and Y. Wu, "Predicting extreme weather events using hybrid machine learning models," *J. Environ. Manage.*, vol. 325, no. 116401, 2023, doi: 10.1016/j.jenvman.2022.116401
- [16] R. J. Hyndman and G. Athanasopoulos, *Forecasting: Principles and Practice*, OTexts, 2018, doi: 10.1016/C2012-0-03732-9
- [17] G. P. Zhang, "Time series forecasting using a hybrid ARIMA and neural network model," *Neurocomputing*, vol. 50, pp. 159–175, 2003, doi: 10.1016/S0925-2312(01)00702-0
- [18] K. Shin, J. Han, and S. Kang, "MI-MOTE: Multiple imputation-based minority oversampling technique for imbalanced and incomplete data classification," *Inf. Sci. (N.Y.)*, vol. 575, pp. 80–89, 2021, doi: 10.1016/j.ins.2021.06.042
- [19] A. Smith and B. Johnson, "Efficient parameter tuning for support vector machines in large-scale datasets," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 30, no. 8, pp. 2404–2415, 2019.
- [20] J. Lee, H. Park, and S. Kim, "Enhanced support vector machines using adaptive kernel functions," *Pattern Recognit. Lett.*, vol. 131, pp. 123–130, 2020.
- [21] K. Ousmane et al., "Novel Classification Method of Spikes Morphology in EEG Signal Using Machine Learning," *Procedia Comput. Sci.*, vol. 148, pp. 70–79, Jul. 2019, doi: 10.1016/j.procs.2019.01.010.