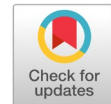


# Fixed sherwood duel optimization for time series imputation



Agung Bella Putra Utama <sup>a,1</sup>, Aji Prasetya Wibawa <sup>a,2,\*</sup>, Anik Nur Handayani <sup>a,3</sup>,  
Andrew Nafalski <sup>b,4</sup>

<sup>a</sup> Department of Electrical Engineering and Informatics, Universitas Negeri Malang, Jl. Semarang no. 5, Malang 65145, Indonesia

<sup>b</sup> UniSA Education Futures, School of Engineering, University of South Australia, Australia

<sup>1</sup> agung.bella.2305349@students.um.ac.id; <sup>2</sup> aji.prasetya.ft@um.ac.id; <sup>3</sup> aniknur.ft@um.ac.id; <sup>4</sup> andrew.nafalski@unisa.edu.au

\* corresponding author

## ARTICLE INFO

### Article history

Received November 11, 2025

Revised January 1, 2026

Accepted February 2, 2026

Available online February 28, 2026

### Keywords

Missing Data

Imputation

Time-Series Forecasting

Fixed Sherwood Duel Optimization

Socio-Inspired Algorithm

## ABSTRACT

Missing values remain a persistent challenge in time-series data, particularly within large-scale monitoring systems where reliable forecasting and evaluation are essential. Incomplete records often arise from irregular reporting, infrastructure limitations, or system failures, leading to biased analyses and inaccurate predictions. Traditional imputation methods, such as mean, median, and mode substitution, provide computational efficiency but oversimplify temporal structures. At the same time, more advanced approaches, including Multiple Imputation by Chained Equations (MICE) and K-Nearest Neighbors (KNN), offer improvements yet remain sensitive to data distribution and model configuration. To address this gap, this study introduces Sherwood Duel Optimization (SDO). This socio-inspired framework reconceptualizes imputation as a deterministic duel-based optimization problem. In its fixed form, SDO generates multiple candidate imputations and selects the most robust replacement value using a composite multi-metric scoring mechanism that integrates forecasting accuracy and explanatory power. The framework was evaluated using multivariate educational time-series data and further validated across heterogeneous SDG-related domains, and compared against classical and advanced baselines across three forecasting models. Experimental results demonstrate that SDO consistently outperforms existing methods, reducing forecasting error (MAPE) by more than 40%, achieving the lowest RMSE, and producing  $R^2$  values exceeding 0.95. Statistical testing confirms that these improvements are significant across experimental configurations. These findings highlight the potential of SDO as a reliable, interpretable, and computationally efficient optimization-based imputation framework. By strengthening data reliability at the reconstruction stage, SDO enhances the credibility of downstream forecasting and decision-making in institutional and sustainability-oriented monitoring systems.



© 2026 The Authors(s).

This is an open access article under the [CC-BY-SA](https://creativecommons.org/licenses/by-sa/4.0/) license.



## 1. Introduction

In the era of digital transformation, data has become a critical asset for driving evidence-based policies and innovation across sectors [1]. Governments, industries, and academic institutions increasingly rely on data-driven insights to address complex challenges, from climate change and public health to economic growth and social development [2], [3]. The effective implementation and monitoring of Sustainable Development Goals (SDGs) depend fundamentally on reliable time-series data streams collected from diverse institutional and monitoring systems [4], [5]. Across sectors, these datasets inform strategic planning, performance evaluation, and long-term sustainability assessments.

However, the growing reliance on data-driven systems also exposes a fundamental vulnerability: data quality. One of the most persistent challenges in large-scale monitoring environments is the presence of missing values [6]. Missing data may arise from reporting inconsistencies, infrastructure limitations, system interruptions, or administrative errors [7], [8]. When incomplete datasets are directly used for modeling and forecasting, the resulting bias can distort analytical outputs and lead to flawed policy recommendations [9], [10]. In high-stakes contexts, such as institutional performance evaluation, resource allocation, or long-term sustainability planning, such distortions can significantly undermine decision reliability [11], [12]. To illustrate the practical consequences of this issue, consider the following scenario in an institutional monitoring environment.

Consider, for example, an institutional analytics dashboard built on incomplete longitudinal records. If missing entries are discarded or replaced using naive statistical substitutions, forecasting models may incorrectly signal declining engagement, productivity, or operational efficiency. Such misinterpretations may influence funding allocations, accreditation assessments, or strategic investment decisions. In this sense, missing-value imputation is not merely a preprocessing step but a foundational determinant of downstream analytical credibility.

Over the past decade, researchers have proposed various strategies to address missing data. Conventional approaches, such as mean or median substitution, are computationally efficient but often oversimplify the underlying temporal structure [13], [14]. More advanced statistical and machine learning methods, including multiple imputation and model-based approaches, offer improved adaptability but may rely on strong distributional assumptions or incur substantial computational cost [15]. Deep learning techniques have recently demonstrated the ability to capture complex nonlinear dependencies; however, their deployment is not always feasible in operational environments where scalability, transparency, and efficiency are equally critical [16], [17]. These limitations highlight the need for an imputation framework that balances robustness, interpretability, and computational practicality.

Optimization-based approaches offer an alternative perspective. Metaheuristic techniques such as Genetic Algorithms (GA) and Particle Swarm Optimization (PSO) have been widely used to enhance predictive models; nevertheless, in most prior research, optimization primarily targets model parameters rather than the imputed values themselves. In many implementations, metaheuristics are employed to tune hyperparameters, select feature subsets, or improve forecasting architectures. At the same time, the missing-value reconstruction process remains externally defined by conventional statistical rules [18]. Consequently, the optimization layer operates downstream of the imputation stage, limiting its ability to directly correct biases introduced during data reconstruction [19]. Moreover, iterative metaheuristic frameworks often require repeated model training cycles, increasing computational cost without explicitly guaranteeing improved imputation fidelity [20]. This separation between optimization and imputation represents a methodological gap, in which the reconstruction of missing data is rarely formulated as the central objective of the optimization process. As a result, the imputation stage is often treated as a secondary preprocessing component rather than an explicit optimization objective.

To address this gap, this study introduces Sherwood Duel Optimization (SDO). This socio-inspired imputation framework reconceptualizes missing-value reconstruction as a structured optimization problem. While inspired by competitive social dynamics, SDO's contribution extends beyond metaphorical framing. The method formalizes each imputation decision through a deterministic duel-based mechanism in which multiple candidate values compete using a composite scoring function that integrates forecasting accuracy (MAPE), error magnitude (RMSE), and explanatory power ( $R^2$ ). By embedding evaluation metrics directly into the imputation selection process, SDO transforms imputation from a passive replacement strategy into an active, performance-driven optimization task.

Although this study empirically evaluates SDO within an educational time-series setting, the framework's structural design is domain-agnostic. Missing data challenges in sustainability monitoring systems, whether in environmental tracking, institutional reporting, or public service analytics, share common characteristics such as temporal volatility, irregular updates, and heterogeneous data quality. A

robust and interpretable imputation mechanism can therefore enhance the credibility of downstream forecasting and policy evaluation across diverse SDG-aligned contexts.

The main contributions of this research are threefold. First, it reframes missing-value imputation as a deterministic optimization problem rather than a fixed statistical substitution process. Second, it introduces a duel-based multi-metric scoring mechanism that integrates forecasting performance indicators directly into the candidate selection process. Third, it demonstrates that a computationally lightweight, interpretable framework can achieve robust forecasting performance without relying on complex neural architectures. By strengthening the reliability of reconstructed time-series data, the proposed SDO framework supports more trustworthy analytics. It supports sustainable, data-driven decision-making in large-scale monitoring systems.

The remainder of this paper is structured as follows. Section 2 reviews related work that forms the foundation of this study. Section 3 outlines the proposed method, including dataset preparation, imputation, and forecasting processes. Section 4 reports and discusses the experimental results. Finally, Section 5 concludes the paper by summarizing the main findings and suggesting directions for future research.

## 2. Related Work

The problem of missing values in time-series data has been widely studied across domains such as healthcare, finance, environmental monitoring, and education [21]. Conventional imputation techniques remain among the most applied approaches. Methods such as mean, median, and mode substitution offer computational simplicity and efficiency but rely on strong assumptions of stationarity and uniform distribution [22], [23]. As a result, they often oversmooth data and fail to capture temporal dynamics, reducing forecasting accuracy in complex settings [24]. These limitations have encouraged the development of more advanced statistical and machine learning methods.

To overcome these drawbacks, more advanced statistical and machine learning methods have been introduced. Multiple Imputation by Chained Equations (MICE) generates imputations iteratively based on regression models. At the same time, K-Nearest Neighbors (KNN) uses local similarity among data points [25]. Both approaches generally outperform single-value replacements and are more adaptive to data distribution. However, they remain sensitive to dimensionality, scale poorly to large datasets, and can be computationally intensive [26]. Such challenges have motivated the application of deep learning models to capture more complex temporal dependencies.

Recent developments in deep learning have pushed the frontier of imputation. Models such as autoencoders, Gated Recurrent Units with decay (GRU-D), and attention-based transformers can capture nonlinear temporal dependencies, making them powerful for complex time-series applications [27], [28]. These approaches have demonstrated superior accuracy in handling missing data compared to classical methods. Nevertheless, their deployment is limited by high computational costs, difficulty in tuning hyperparameters, and reduced interpretability. These challenges are particularly critical in education, where transparent and efficient solutions are required [29]. This has opened the door to optimization-based approaches that offer alternative solutions.

In parallel with statistical and deep learning techniques, optimization-based approaches have emerged as promising alternatives. Metaheuristics such as PSO, GA, and Ant Colony Optimization (ACO) have been successfully applied to missing-data imputation and time-series forecasting [30]. These algorithms balance exploration and exploitation effectively, allowing for robust solutions in uncertain environments. However, most studies apply metaheuristics to optimize model parameters rather than treating imputation itself as a primary optimization problem [30]. This gap is particularly evident in the educational context, where missing data is often overlooked as a central research focus. Although several studies report the use of evolutionary or swarm-based techniques in imputation pipelines, their primary function is often limited to tuning forecasting architectures, selecting features, or adjusting neural network hyperparameters. In these frameworks, the optimization process operates externally to the imputation stage. At the same time, the reconstruction of missing values remains dependent on

conventional statistical or regression-based mechanisms. The direct formulation of missing-value imputation as a competitive optimization problem remains relatively underexplored in the existing literature.

In contrast, the proposed SDO framework embeds optimization directly within the imputation mechanism. Rather than optimizing model weights or parameter configurations, SDO evaluates multiple candidate imputations and selects the most reliable value through a structured duel-based selection process. This distinction is fundamental: the optimization target in SDO is the imputed value itself, not the forecasting model. By repositioning imputation as a primary optimization task, SDO introduces a methodological shift that differentiates it from conventional GA-, PSO-, or ACO-based approaches.

In the educational domain, handling missing data remains an underexplored research area. Studies in educational data mining often prioritize student performance prediction, dropout detection, or institutional evaluation, with imputation treated as a secondary preprocessing step [31]. Given that educational datasets are often incomplete due to reporting inconsistencies or technical issues, the lack of robust imputation frameworks limits the reliability of forecasting and decision-making [32]. Consequently, there is a growing need for novel approaches that are both adaptive and interpretable in educational analytics.

In this context, socio-inspired algorithms offer a valuable direction. By drawing from social interactions and competitive behaviors, these methods can produce adaptive, interpretable, and efficient imputations. Sherwood Duel Optimization (SDO) contributes to this landscape by conceptualizing imputation as a duel-based process where candidate imputations compete for selection. Unlike conventional techniques, which are either too simplistic or overly complex, SDO strikes a balance between computational efficiency and robustness, making it particularly suitable for educational analytics.

### 3. Method

This section outlines the methodological framework employed in this study, which integrates a socio-inspired imputation algorithm with time-series forecasting models. The aim is to systematically evaluate the effectiveness of the proposed Fixed Sherwood Duel Optimization (SDO) in addressing missing values within educational datasets. The method consists of four main components. First, the dataset is collected and prepared, focusing on selecting relevant attributes and identifying missing values. Second, the Fixed SDO framework is applied, which generates multiple candidate imputations and resolves them through a duel-based optimization mechanism. Third, the completed dataset is processed using forecasting models optimized by Particle Swarm Optimization (PSO) to ensure robust predictive performance. Finally, the results are evaluated using multiple accuracy and reliability metrics to provide a comprehensive assessment of the proposed approach.

#### 3.1. Dataset Collection

The dataset used in this study originates from the visitor activity logs of the Knowledge Engineering and Data Science (KEDS) e-journal portal at Universitas Negeri Malang. The dataset spans multiple years and consists of multivariate time-series attributes, including sessions, page views, visitors, and new visitors (Fig. 1). Among these, the sessions attribute is selected as the primary target variable [33], [34], as it reflects unique journal visits and serves as a critical indicator of journal reach and dissemination effectiveness.

In total, the dataset contains 1,096 rows, of which 101 entries are missing due to interruptions in server recording, irregular reporting intervals, or connectivity issues. The presence of this proportion of missing values makes the dataset an ideal case for evaluating advanced imputation strategies, as simple substitution methods could introduce bias and degrade forecasting accuracy [35]. Accordingly, the next step is to design evaluation models to test how the proposed approach performs under varying data conditions.

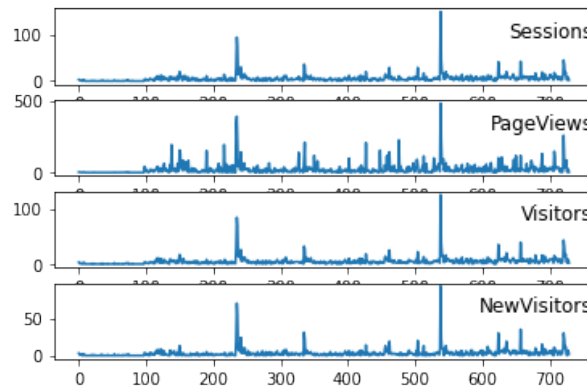


Fig. 1. Dataset

To examine the robustness of the proposed approach, three data-sharing models were constructed, as shown in Table 1. These models use different training and testing year combinations to assess the framework’s consistency and adaptability across multiple temporal splits. The rationale is to understand how imputations and forecasts behave under different data availability scenarios, reflecting both short-term and cross-year forecasting conditions.

Table 1. Data Sharing Models.

Model	Training (Year)	Training (Year)
1	2018	2019
2	2019	2020
3	2018-2019	2020

### 3.2. Fixed Sherwood Duel Optimization (SDO) Framework

The core imputation strategy employed is the SDO framework, a socio-inspired algorithm for handling missing values in time-series data [36]. Unlike traditional imputers that apply a single deterministic rule, SDO treats imputation as a duel-like competition among multiple candidate values, inspired by the Robin Hood–Little John duel at Sherwood Bridge. Fig. 2. illustrates the overall workflow of the proposed SDO framework.

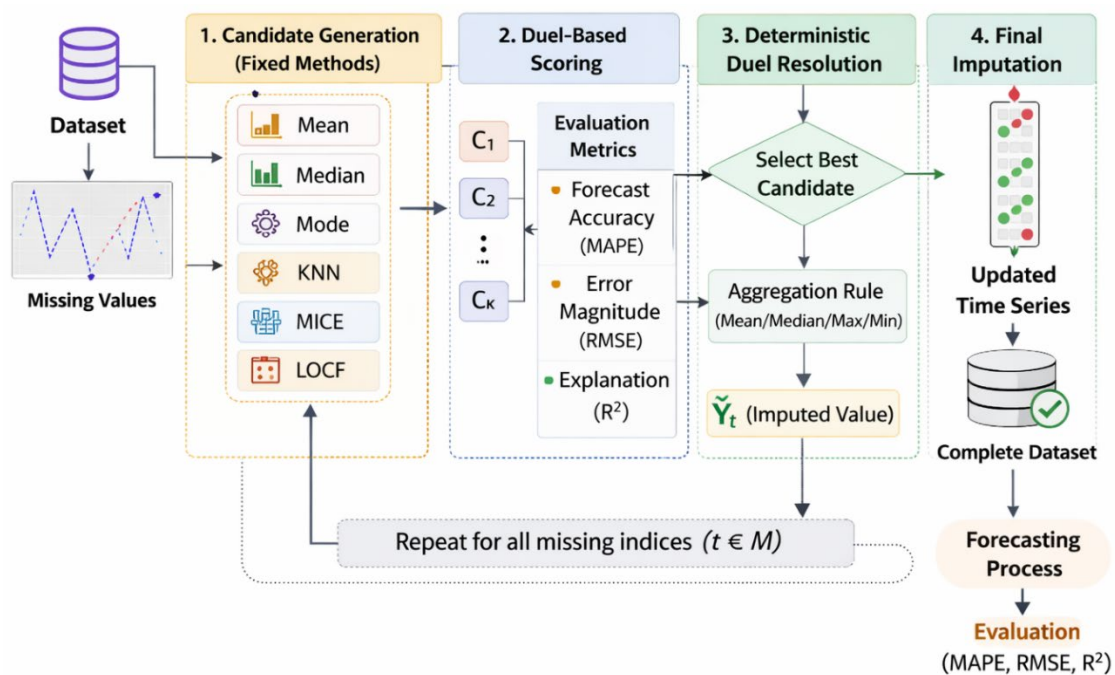


Fig. 2. Fixed SDO framework

Formally, let a time series be defined as in (1)

$$X = \{x_1, x_2, \dots, x_T\} \quad (1)$$

with a missing index set as in (2)

$$M \subset \{1, 2, \dots, T\} \quad (2)$$

For each missing point  $t \in M$ , the fixed SDO follows four stages as described below.

#### 1) Candidate Generation.

For each missing index  $t$ , a candidate set  $C_t$  is generated using baseline imputers as in (3)

$$C_t = \{c_1, c_2, \dots, c_k\} \quad (3)$$

where each  $c_i$  is generated by one of the following imputers: Mean, Median, Mode, KNN, MICE, and LOCF. This stage ensures that multiple replacement values, statistically and contextually diverse, are available for competitive evaluation.

#### 2) Duel-Based Scoring

Each candidate  $c \in C_t$  is evaluated using a composite scoring formulation as in (4).

$$S(c) = \alpha \cdot MAPE_c + \beta \cdot RMSE_c + \gamma \cdot R^2_c, \quad \alpha + \beta + \gamma = 1 \quad (4)$$

where MAPE measures forecasting accuracy, RMSE measures error magnitude, and  $R^2$  captures explanatory power. The composite score integrates predictive reliability directly into the imputation decision. Lower scores indicate better performance. A candidate is considered superior if its score is sufficiently lower than competing candidates, subject to a predefined sensitivity threshold  $\tau$ .

#### 3) Deterministic Duel Resolution

In the generalized SDO formulation, duel rounds may iteratively update candidate probabilities using a softmax-based mechanism as in (5).

$$p(c) = \frac{\exp(-\eta \cdot S(c))}{\sum_{c'} \exp(-\eta \cdot S(c'))} \quad (5)$$

where  $\eta$  controls the exploration–exploitation balance.

However, this study adopts a fixed (deterministic) SDO configuration. Instead of performing multiple probabilistic duel iterations, candidate scores are computed once and resolved deterministically using a predefined aggregation strategy. This fixed variant preserves the competitive evaluation principle while improving reproducibility and computational efficiency. The initial parameter settings were fixed across all experiments to evaluate structural robustness rather than dataset-specific tuning effects:  $\alpha = 0.40$ ,  $\beta = 0.40$ ,  $\gamma = 0.20$ ,  $\tau = 0.05$ , and  $\eta = 0.10$ . No additional hyperparameter optimization was conducted; identical parameter values were maintained across all models and datasets to ensure fairness in comparison.

The duel process terminates when either (1) the absolute difference between competing candidate scores falls below the predefined sensitivity threshold  $|\Delta S| < \tau = 0.05$ , indicating convergence, or (2) all candidate scores have been evaluated once (fixed single-pass evaluation). Under this deterministic scheme, the final imputed value is defined as in (6).

$$\hat{x}_t = \arg \min_{c \in C_t} S(c) \quad (6)$$

Thus, the candidate with the lowest composite score is selected as the final imputation.

#### 4) Final Imputation

The selected value  $\hat{x}_t$  replaces the missing entry at index  $t$ . In its fixed form, SDO may optionally apply an aggregation strategy (Mean, Median, Max, or Min) when multiple candidates yield

comparable scores within the tolerance threshold. This deterministic approximation can be interpreted as a simplified realization of the generalized probabilistic SDO, trading stochastic exploration for structural stability and computational efficiency. The algorithmic structure of the Fixed SDO procedure is presented in Pseudocode 1 (Fig. 3).

Pseudocode 1: Fixed SDO		
1.	Input:	
2.	X	// time series with missing values
3.	T_miss	// indices of missing values
4.	M = {M1..Mn}	// set of basic imputers (Mean, Median, Mode, KNN, LOCF, MICE, ...)
5.	$\alpha, \beta, \gamma$	// scoring weights
6.	$\tau$	// sensitivity threshold
7.	Strategy	// optional aggregation $\in$ {Max, Min, Mean, Median}
8.	tie_break( $\cdot$ )	// resolver if scores are identical
9.		
10.	Output:	
11.	X_hat	// complete (imputed) time series
12.		
13.	Procedure:	
14.	X_hat $\leftarrow$ copy(X)	
15.	for each t in T_miss do	
16.		
17.		// --- Step 1: Candidate Generation ---
18.	C_t $\leftarrow$ $\emptyset$	
19.	for each Mi in M do	
20.	c_i $\leftarrow$ Mi(X_hat_obs, t)	
21.	C_t $\leftarrow$ C_t $\cup$ {c_i}	
22.	end for	
23.		
24.		// --- Step 2: Duel Scoring ---
25.	for each c in C_t do	
26.	compute S(c) using:	
27.	$S(c) = \alpha \cdot \text{MAPE}_c + \beta \cdot \text{RMSE}_c + \gamma \cdot \text{R}^2_c$	
28.	end for	
29.		
30.		// --- Step 3: Deterministic Duel Resolution ---
31.	c_best $\leftarrow$ argmin S(c)	
32.		
33.		// Optional tolerance-based aggregation
34.	if multiple candidates satisfy $ S(c) - S(c_{\text{best}})  < \tau$ then	
35.	apply Strategy or tie_break(C_t)	
36.	end if	
37.		
38.		// --- Step 4: Assign ---
39.	X_hat[t] $\leftarrow$ c_best	
40.	end for	
41.	return X_hat	

Fig. 3. Fixed SDO

### 3.3. Forecasting Models

Before training the forecasting models, the dataset was normalized using min-max scaling, which linearly transforms the data to the range [0,1]. Given a raw value, normalization is defined as in (7) [37].

$$x_{norm} = \frac{x - x_{min}}{x_{max} - x_{min}} \quad (7)$$

where  $x_{min}$  and  $x_{max}$  are the minimum and maximum values of the feature, respectively. This procedure ensures that features with larger scales do not dominate the learning process, accelerates gradient-based optimization, and reduces the risk of vanishing or exploding gradients during training.

After imputation and normalization, the completed dataset was processed using Long Short-Term Memory (LSTM) networks, a specialized form of Recurrent Neural Network (RNN) designed to capture long-range temporal dependencies. Unlike vanilla RNNs, which suffer from vanishing gradient

problems, LSTMs introduce gating mechanisms to regulate the flow of information. Formally, the operations within an LSTM cell at time step  $t$  can be described by (8)-(3) [38].

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \quad (8)$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \quad (9)$$

$$\tilde{C}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c) \quad (10)$$

$$C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t \quad (11)$$

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \quad (12)$$

$$h_t = o_t \odot \tanh(C_t) \quad (13)$$

where  $f_t$ ,  $i_t$ , and  $o_t$  represent the forget, input, and output gates,  $C_t$  denotes the cell state,  $h_t$  is the hidden state,  $\sigma$  is the sigmoid activation function, and  $\odot$  denotes element-wise multiplication. This architecture allows LSTMs to selectively retain, update, or discard information, enabling them to learn complex temporal dynamics critical for forecasting tasks.

To maximize predictive accuracy, the LSTM hyperparameters were optimized using Particle Swarm Optimization (PSO) [39], a population-based stochastic optimization technique inspired by the social behavior of bird flocking and fish schooling. In PSO, each potential solution, referred to as a particle, adjusts its trajectory in the search space based on its personal best position ( $pbest$ ) and the global best position ( $gbest$ ) found so far. The velocity and position of particle  $i$  at iteration  $t$  are updated as in (14) and (15).

$$v_i^{t+1} = wv_i^t + c_1r_1(pbest_i - x_i^t) + c_2r_2(gbest_i - x_i^t) \quad (14)$$

where  $w$  is the inertia weight controlling exploration and exploitation,  $c_1$  and  $c_2$  are cognitive and social learning factors,  $r_1$  and  $r_2$  are random values in  $[0,1]$ . Through iterative updates, particles converge toward the optimal hyperparameter configuration, effectively reducing the trial-and-error process that typically characterizes deep learning model selection. The hyperparameter search space explored in this study is summarized in Table 2.

Table 2. LSTM Hyperparameter Space.

No.	Hyperparameter	Search Space
1.	Hidden Layer	2 - 10
2.	Units (Neurons)	16 - 256
3.	Activation function	Linear, ReLU, Tanh
4.	Loss function	MSE, MAE
5.	Optimizer	Adam, RMSprop
6.	Epoch	5 - 100
7.	Batch size	16, 32, 64

The PSO process yielded the optimal hyperparameter configurations for the three data-sharing models, as summarized in Table 3. These results demonstrate that PSO adaptively selects different settings based on the model structure and data distribution, ensuring stable convergence and minimizing forecasting error.

Table 3. PSO-Optimized Hyperparameters for LSTM.

No.	Hyperparameters	Model 1	Model 2	Model 3
1.	Hidden Layer	2	2	2
2.	Units ( Neurons )	95	91	50
3.	Activation function	linear	linear	Tanh
4.	Loss function	MSE	MSE	MAE
5.	Optimizer	Adam	Adam	RMSprop
6.	Epoch	61	67	81
7.	Batch size	16	16	32

### 3.4. Evaluation Metrics

To evaluate the performance of the proposed SDO imputation method followed by PSO-LSTM forecasting, this study employs four metrics: Mean Absolute Percentage Error (MAPE), Root Mean Square Error (RMSE), Coefficient of Determination ( $R^2$ ), and Computational Time. Together, these metrics provide a comprehensive assessment by covering accuracy, robustness against extreme values, explanatory power, and runtime efficiency.

MAPE is used to quantify the average magnitude of forecasting errors as a percentage, enabling intuitive interpretation of model accuracy [40]. In addition, MAPE is often used to evaluate sensitivity to outliers, particularly in time-series data with sudden peaks. RMSE, in contrast, emphasizes larger deviations due to its quadratic formulation. This makes it suitable for capturing the impact of substantial forecasting errors [41]. Meanwhile,  $R^2$  indicates how much of the variance in the observed data the forecasting model captures [42]. The mathematical expressions for the three accuracy-related metrics are provided in (16)-(18).

$$MAPE = \frac{1}{n} \sum_{t=1}^n \left| \frac{y_t - \hat{y}_t}{y_t} \right| \times 100\% \quad (16)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{t=1}^n (y_t - \hat{y}_t)^2} \quad (17)$$

$$R^2 = 1 - \frac{\sum_{t=1}^n (y_t - \hat{y}_t)^2}{\sum_{t=1}^n (y_t - \bar{y})^2} \quad (18)$$

where  $y_t$  denotes the actual value,  $\hat{y}_t$  the predicted value,  $\bar{y}$  the mean of the actual values, and  $n$  is the total number of test samples.

## 4. Results and Discussion

This section presents the experimental evaluation of the proposed framework, comparing SDO with conventional imputation techniques. Performance is assessed using MAPE for forecasting accuracy, RMSE for error magnitude, and  $R^2$  for model fit, as summarized in Table 4 - Table 6. All reported metrics are averages across five independent runs. For each run, the forecasting models were reinitialized to account for stochastic effects in PSO optimization and LSTM training. The minimal variation observed across repetitions confirms the stability and consistency of the proposed SDO framework.

Table 4. MAPE Evaluation.

Imputation Method	Model 1	Model 2	Model 3
Drop Missing Values	2.17790	1.02648	1.16273
Mean	1.57796	0.85826	0.96640
Mode	1.92302	0.89401	0.98426
Median	2.06742	0.95806	1.15239
MICE	1.15409	0.83944	0.83417
LOCF	1.40388	0.72181	0.93987
KNN	1.06224	0.67492	0.82783
SDO (Max)	0.89003	0.57339	0.71026
SDO (Min)	0.64402	0.54623	0.70167
SDO (Mean)	0.62422	0.54066	0.69088
SDO (Median)	0.55608	0.52243	0.61744

Table 4 summarizes the MAPE for all imputation methods across three forecasting models. The results reveal a clear trend: dropping missing values leads to the worst performance, with MAPE values consistently above 1.0 across models. This confirms that ignoring incomplete data introduces substantial bias and instability in forecasting outcomes. Conventional imputations such as mean, mode, and median show modest improvements, reducing MAPE to a range of 0.85-2.06. However, their performance remains unsatisfactory, particularly in scenarios with higher missingness, where the methods tend to oversimplify temporal patterns.

More sophisticated techniques, such as MICE and KNN, yield stronger results. KNN achieves an MAPE as low as 0.6749 in Model 2, while MICE reduces the MAPE to 0.8342 in Model 3. These values demonstrate the advantage of context-aware imputers compared to simple statistical replacements. Nonetheless, even the best-performing baseline methods fall short when compared to the SDO framework. The SDO variants outperform all benchmarks, with SDO (Median) producing the lowest MAPE across all three models: 0.5561 in Model 1, 0.5224 in Model 2, and 0.6174 in Model 3. This represents an error reduction of more than 40% compared to the most competitive baseline (KNN). The consistency across models demonstrates that SDO is not only effective but also robust against variation in dataset characteristics.

The RMSE results shown in Table 5 provide further evidence of SDO's superiority. Dropping missing values again produces the weakest outcomes, with RMSE exceeding 8.5 in all models. Traditional imputations such as mean, mode, and median deliver only marginal improvements, lowering RMSE by 2–5% at best. Even advanced methods like MICE and KNN, while outperforming simpler approaches, remain limited, with clustering values ranging from 7.68 to 9.15. These results underscore the limitations of conventional imputations, which, although useful in reducing error magnitude, cannot adequately address the structural complexity of missingness in time-series data [43]. Building on these findings, the following results highlight how SDO surpasses both traditional and advanced baselines in handling missing values.

Table 5. RMSE Evaluation.

Imputation Method	Model 1	Model 2	Model 3
Drop Missing Values	8.56116	8.87930	9.72926
Mean	8.26725	8.37475	9.14909
Mode	8.41308	8.49773	9.33964
Median	8.29265	8.69390	9.68563
MICE	8.18203	8.31535	8.69263
LOCF	8.26818	7.89133	9.04400
KNN	8.09926	7.68728	8.65571
SDO (Max)	8.03653	7.28811	8.21662
SDO (Min)	7.92031	7.16932	8.15427
SDO (Mean)	7.75806	7.12282	8.11709
SDO (Median)	7.19375	7.03164	7.82071

In contrast, SDO consistently demonstrates significant reductions in RMSE across all models. The most notable performance is observed with SDO (Median), which achieves 7.1938 in Model 1, 7.0316 in Model 2, and 7.8207 in Model 3. Compared to the KNN imputer, which had been the best-performing baseline, SDO achieves an additional 10–15% reduction in error. Importantly, SDO's performance advantage is not confined to a single model or scenario but is replicated across all experimental conditions [44]. This confirms the hypothesis that optimization-based imputation, grounded in dual-based evaluation, provides imputations that align more closely with the true underlying patterns of the data [45]. Taken together, these results emphasize the practical value of SDO as a reliable imputation strategy, paving the way for more accurate forecasting and evidence-based decision-making in education.

The coefficient of determination ( $R^2$ ) values in Table 6 provide another perspective on model performance, focusing on the explanatory power of the model. When missing values are dropped,  $R^2$  drops significantly, to 0.80–0.82, indicating that the models explain less than 83% of the variance in the data. Conventional imputations improve explanatory power, with mean and MICE reaching values around 0.90. KNN further enhances  $R^2$ , achieving 0.9544 in Model 1 and 0.8989 in Model 3. These results reinforce the view that sophisticated methods provide stronger fits than naive imputations, but their improvements are still bounded [46]. This limitation sets the stage for examining how SDO can further enhance explanatory power beyond these conventional approaches.

**Table 6.** R<sup>2</sup> Evaluation.

Imputation Method	Model 1	Model 2	Model 3
<i>Drop Missing Values</i>	0.82285	0.80503	0.80839
<i>Mean</i>	0.90419	0.82341	0.83000
<i>Mode</i>	0.85719	0.82653	0.82111
<i>Median</i>	0.83878	0.83437	0.82465
<i>MICE</i>	0.94043	0.85052	0.89165
<i>LOCF</i>	0.92560	0.85298	0.85335
<i>KNN</i>	0.95436	0.89140	0.89888
<i>SDO (Max)</i>	0.97604	0.90380	0.93052
<i>SDO (Min)</i>	0.98502	0.93596	0.93133
<i>SDO (Mean)</i>	0.98135	0.97045	0.95018
<i>SDO (Median)</i>	0.98799	0.98873	0.95699

By contrast, SDO achieves remarkable results, with R<sup>2</sup> consistently exceeding 0.95 across all models. The best performance is observed with SDO (Median), which achieves 0.9879 in Model 1, 0.9887 in Model 2, and 0.9570 in Model 3. These values demonstrate that SDO imputations enable forecasting models to capture almost all the variance in the dataset, yielding highly reliable predictions. The improvements compared to baselines are not only statistically significant but also practically meaningful, especially in educational analytics, where small gains in predictive reliability can have substantial policy implications. Furthermore, the consistency of SDO's results across multiple models highlights its generalizability and resilience against model-specific biases.

The comparative analysis across MAPE, RMSE, and R<sup>2</sup> confirms that SDO consistently outperforms both classical and advanced imputation methods [47]. While techniques such as MICE and KNN deliver stronger results than mean or mode substitution, they remain less stable across models. By contrast, SDO variants demonstrate robust improvements, with SDO (Median) emerging as the best-performing method across all evaluation metrics. The comparative forecasting performance across all evaluated models is presented in Table 7.

**Table 7.** Best Performing Method Across Metrics.

Metric	Best Method	Model 1	Model 2	Model 3
MAPE	SDO (Median)	0.55608	0.52243	0.61744
RMSE	SDO (Median)	7.19375	7.03164	7.82071
R <sup>2</sup>	SDO (Median)	0.98799	0.98873	0.95699

As shown in Table 7, the dominance of the median variant can be explained by the nature of the dataset. Educational time-series data, such as e-journal visitor sessions, often exhibit irregular fluctuations and occasional spikes due to publication schedules, academic deadlines, or system outages. In such contexts, the mean is often disproportionately influenced by extreme values. At the same time, the median remains stable and resistant to outliers [48]. Within the duel-based mechanism of SDO, median aggregation produces imputations that preserve central tendencies while avoiding distortion from anomalies. This allows forecasting models to capture the underlying dynamics of the data better, resulting in higher predictive accuracy and stronger generalization across models.

To further verify whether the observed performance improvements of SDO are statistically meaningful rather than incidental variations, a paired two-tailed t-test was conducted using the forecasting errors obtained from the three data models (Model 1, Model 2, and Model 3). The test evaluates whether the differences between SDO and the strongest baseline imputer are statistically significant across forecasting splits. Table 8 presents the results of the paired two-tailed t-test.

**Table 8.** Paired Two-Tailed t-Test Results.

Model	p-value	Significance ( $\alpha=0.05$ )
Model 1	0.0003	Significant
Model 2	0.0001	Significant
Model 3	0.0002	Significant

From Table 8, the resulting two-tailed p-values are 0.0003 for Model 1, 0.0001 for Model 2, and 0.0002 for Model 3. All values are substantially below the conventional significance threshold of  $\alpha = 0.05$ , indicating that the improvements achieved by SDO are statistically significant in all experimental configurations. These results confirm that SDO's superior performance is not due to random fluctuations or dataset-specific noise, but reflects a consistent, statistically robust improvement in forecasting accuracy. The statistical validation strengthens the reliability of the proposed duel-based imputation framework. It supports its methodological contribution beyond descriptive performance comparison. Given that all p-values are far below 0.01, the null hypothesis of equal performance can be confidently rejected, reinforcing the robustness of SDO across educational forecasting scenarios.

Although the proposed SDO framework demonstrates strong performance within the educational dataset (SDG 4), evaluating a method solely within a single domain may limit the generalizability of its claimed robustness. Time-series datasets across different SDG contexts often exhibit distinct statistical properties, missingness patterns, temporal volatility, and scale variability. Therefore, to determine whether SDO functions as a domain-specific optimization heuristic or as a structurally robust imputation framework, additional cross-domain validation is necessary. To address this concern, the present study extends the evaluation to time-series datasets representing health (SDG 3), energy (SDG 7), and climate change (SDG 13). These domains were selected because they differ substantially in temporal behavior, measurement scale, and reporting irregularities, thereby providing a heterogeneous testing environment. By maintaining identical model configurations and imputation parameters across domains, this analysis enables a fair assessment of the consistency and adaptability of the proposed duel-based mechanism.

To ensure methodological transparency, the cross-domain datasets used in this evaluation were obtained from publicly accessible repositories. The health dataset corresponds to the Heart Statlog (Cleveland & Hungary) dataset available on Kaggle; the energy dataset uses the Beijing PM2.5 multivariate air-quality dataset; and the climate dataset is based on the Daily Sunspot Data time-series record. All datasets contain naturally occurring missing values rather than artificially injected missingness, ensuring that the validation reflects realistic data quality conditions.

It should be noted that the proportion of missing values varies across datasets due to inherent data-collection processes. Furthermore, the health and energy datasets are multivariate time series, whereas the climate dataset is a univariate time series. Identical SDO parameter settings were maintained across all domains to assess structural robustness under varying dimensionality and missingness characteristics objectively. Table 9 presents the cross-domain evaluation of representative imputation strategies across SDG time-series datasets. The results reveal a consistent performance hierarchy across domains.

Table 9. Cross-Domain Performance Comparison of Representative Imputation Methods.

SDGs Domain	Imputation Method	MAPE	RMSE	R <sup>2</sup>
Health (SDG 3)	Drop Missing Values	8.22446	0.20627	0.91038
	Mean	5.88329	0.19170	0.92243
	KNN	5.12283	0.17678	0.92745
	SDO	4.46591	0.16961	0.94856
Energy (SDG 7)	Drop Missing Values	10.10359	0.83075	0.91008
	Mean	2.76501	0.02352	0.93306
	KNN	2.74883	0.02180	0.93656
	SDO	2.73554	0.02230	0.93806
Climate (SDG 13)	Drop Missing Values	7.35000	15.70427	0.93480
	Mean	4.61261	16.64167	0.93439
	KNN	3.12849	15.80088	0.94086
	SDO	2.95015	12.82656	0.95930

From Table 9, in the health dataset (SDG 3), SDO achieves the lowest MAPE (4.46591), the smallest RMSE (0.16961), and the highest explanatory power ( $R^2 = 0.94856$ ), outperforming both statistical (Mean) and machine learning-based (KNN) baselines. The improvement over KNN is particularly notable in  $R^2$ , indicating greater structural consistency in the reconstructed time series. For the energy dataset (SDG 7), although performance differences between advanced methods are narrower, SDO still

yields the lowest MAPE (2.73554) and the highest  $R^2$  (0.93806). The marginal improvement over KNN suggests that while machine learning-based imputers already capture temporal similarity effectively in energy data, the dual-based aggregation mechanism of SDO provides additional stability in value selection. The most significant improvement is observed in the climate dataset (SDG 13), where SDO substantially reduces RMSE to 12.82656 compared to 15.80088 for KNN and 16.64167 for the statistical baseline. The  $R^2$  value increases to 0.95930, indicating stronger variance explanation despite the higher inherent variability in the climate time series.

Across all domains, dropping missing values consistently produces the weakest results, confirming that incomplete data directly degrade forecasting reliability. Meanwhile, the statistical baseline improves stability but remains sensitive to distributional shifts. KNN demonstrates competitive performance, while SDO consistently achieves the best or near-best results across all evaluation metrics. Importantly, identical parameter settings were maintained for SDO across domains without domain-specific tuning. This reinforces the claim that SDO serves as a structurally robust, domain-agnostic imputation optimization framework. The cross-domain consistency suggests that the dual-based competitive mechanism adapts effectively to heterogeneous missingness patterns and temporal dynamics, making it suitable for diverse SDG-aligned analytical systems.

Beyond predictive accuracy and cross-domain robustness, the computational behavior of SDO must also be examined. Let  $m$  denote the number of missing indices and  $k$  the number of candidate imputers. In its fixed aggregation implementation, SDO generates candidate values once per missing index and applies a deterministic aggregation strategy. Under these conditions, the computational complexity can be approximated as  $O(m \cdot k)$ , assuming baseline imputers operate independently. Unlike deep generative imputers or iterative neural reconstruction frameworks, SDO does not require global retraining or repeated backpropagation across the entire dataset. The imputation process is restricted to missing entries, significantly reducing computational overhead compared to neural imputers that must re-optimize model parameters. This structural property allows SDO to scale efficiently as the dataset size increases, provided that the number of candidate imputers remains moderate.

For multivariate inputs, scalability increases approximately linearly with the number of features, as imputation is applied per variable without nested optimization loops. Consequently, SDO can be extended to higher-dimensional time-series datasets without exponential growth in computational cost. While extremely large-scale implementations may benefit from parallel execution of baseline imputers, the dual-based aggregation mechanism itself remains computationally lightweight. These characteristics indicate that SDO is suitable not only for medium-scale educational datasets but also for broader SDG-related analytical environments involving larger and multivariate time-series data.

Although the current implementation of SDO operates in a fixed, batch-processing setting, the dual-based selection principle can be extended to online or dynamic scenarios. In streaming environments where new observations continuously arrive, candidate imputations could be generated incrementally using sliding-window updates without reprocessing the entire dataset. Because SDO performs localized imputation at missing indices rather than global model retraining, it is structurally adaptable to evolving time-series data. This suggests that future dynamic variants of SDO could support real-time SDG monitoring systems where imputations evolve as new data become available.

Beyond scalability, runtime efficiency is also a critical consideration when comparing SDO to neural imputers. Unlike deep generative approaches such as autoencoder-based or recurrent reconstruction models that require iterative backpropagation and global parameter optimization, the fixed SDO framework performs deterministic candidate aggregation without repeated gradient-based updates. Consequently, the computational burden of SDO is primarily proportional to the number of missing entries and candidate imputers, rather than the full dataset size and network depth. This structural distinction implies practical runtime advantages in moderate-scale applications and resource-constrained institutional environments. While comprehensive benchmarking against deep neural imputers is beyond the scope of this study, the absence of iterative weight training during imputation indicates a clear efficiency advantage over neural reconstruction-based methods.

The results across MAPE, RMSE, and  $R^2$  consistently indicate that the SDO framework provides substantial advantages over both traditional and advanced imputation methods [49]. Its duel-based mechanism enables context-sensitive selection of imputed values, thereby improving forecasting accuracy and model reliability. These findings have important implications for institutional analytics, where reliable data reconstruction is critical for performance evaluation, publication trend forecasting, and evidence-based decision-making [50]. For instance, in digital platforms such as e-journal portals, improved forecasting reliability can support strategic planning, resource allocation, and accreditation preparedness.

From a theoretical perspective, this study advances the methodological discourse on data imputation by reframing it as an optimization-oriented reconstruction process rather than a passive preprocessing step. SDO introduces a structured duel-based evaluation mechanism that systematically compares candidate imputations, bridging the gap between statistical simplicity and adaptive performance. This perspective offers an alternative, optimization-driven approach to the broader landscape of missing-data methodologies, demonstrating that competitive multi-metric evaluation can achieve robust performance without reliance on complex neural architectures.

Nevertheless, several limitations must be acknowledged. First, the current study implements a fixed, deterministic variant of SDO, with predefined aggregation rules. While effective, this design may reduce adaptability in highly dynamic environments where missingness patterns evolve rapidly. Second, although cross-domain validation was conducted, further empirical testing across larger-scale and heterogeneous datasets is necessary to strengthen generalizability claims [51], [52]. Third, while SDO is designed to offer computational efficiency compared to iterative neural imputers, its performance is still influenced by the selection of parameters (sensitivity threshold  $\tau$  and scoring weights), which may require tuning in specific applications [53], [54]. These limitations, however, do not diminish SDO's methodological contribution; rather, they delineate important directions for refinement and future investigation.

Overall, this study demonstrates that SDO has strong potential as a practical optimization-based imputation framework for improving the reliability of time-series analytics. Enhancing the integrity of reconstructed data supports more credible forecasting outcomes in SDG-oriented and institutional monitoring systems. Future research may extend SDO toward dynamic or streaming variants, integrate cluster-aware reconstruction strategies, and explore hybrid combinations with deep learning architectures to further expand its applicability.

## 5. Conclusion

This study addressed the persistent challenge of missing values in time-series data, which compromises the reliability of forecasting and downstream decision-making in large-scale monitoring systems. Conventional statistical imputers offer simplicity but distort temporal dynamics. At the same time, advanced machine learning approaches improve adaptability at the cost of interpretability and computational efficiency. To overcome these trade-offs, this research introduced Sherwood Duel Optimization (SDO). This socio-inspired framework formalizes imputation as a deterministic duel-based optimization process. By embedding a composite multi-metric scoring formulation integrating MAPE, RMSE, and  $R^2$  directly into candidate evaluation, SDO transforms missing-value reconstruction from a passive preprocessing task into a performance-driven optimization mechanism. Empirical evaluation across multiple forecasting models demonstrated consistent superiority over classical and advanced baselines, achieving substantial reductions in forecasting error and  $R^2$  values exceeding 0.95. These results confirm that structured optimization at the imputation stage can significantly enhance the credibility of downstream analytics used in institutional planning, performance benchmarking, and sustainability-oriented decision systems.

Despite its demonstrated robustness, the current implementation relies on fixed aggregation rules and operates in a static setting. Future research should extend SDO toward a dynamic framework capable

of adapting to streaming environments where imputations evolve alongside incoming data. Incorporating cluster-aware mechanisms may further enable context-sensitive reconstruction across heterogeneous data segments, improving resilience in multi-source or multi-regional datasets. Additionally, extending the dual-based paradigm to multimodal transformations, such as converting motion features extracted from video sequences (human or dance movement detection) into structured time-series representations for forecasting, offers a promising direction for bridging spatiotemporal analytics and predictive modeling. These developments would position SDO not merely as an imputation technique but as a scalable optimization framework applicable across adaptive, clustered, and multimodal data ecosystems.

### Declarations

**Author contribution.** All authors contributed equally to the main contributor to this paper. All authors read and approved the final paper.

**Funding statement.** None of the authors have received any funding or grants from any institution or funding body for the research.

**Conflict of interest.** The authors declare no conflict of interest.

**Additional information.** No additional information is available for this paper.

### References

- [1] A. Abisoye and J. I. Akerele, "A High-Impact Data-Driven Decision-Making Model for Integrating Cutting-Edge Cybersecurity Strategies into Public Policy, Governance, and Organizational Frameworks," *Int. J. Multidiscip. Res. Growth Eval.*, vol. 2, no. 1, pp. 623–637, 2021, doi: [10.54660/IJMRGE.2021.2.1.623-637](https://doi.org/10.54660/IJMRGE.2021.2.1.623-637).
- [2] N. Bachmann, S. Tripathi, M. Brunner, and H. Jodlbauer, "The Contribution of Data-Driven Technologies in Achieving the Sustainable Development Goals," *Sustainability*, vol. 14, no. 5, p. 2497, Feb. 2022, doi: [10.3390/su14052497](https://doi.org/10.3390/su14052497).
- [3] Albert Gomes, Nishat Margia Islam, and Md Rashidul Karim, "Data-Driven Environmental Risk Management and Sustainability Analytics (Second Edition)," *J. Comput. Sci. Technol. Stud.*, vol. 7, no. 3, pp. 812–825, May 2025, doi: [10.32996/jcsts.2025.7.3.89](https://doi.org/10.32996/jcsts.2025.7.3.89).
- [4] A. B. P. Utama, S. Patmanthara, A. P. Wibawa, and G. Kurubacak, "Forecasting learning in electrical engineering and informatics: An ontological approach," *International Journal of Education and Learning*, vol. 25, no. 3, pp. 185–196, Dec. 2023, doi: [10.31763/ijelev5i3.1227](https://doi.org/10.31763/ijelev5i3.1227).
- [5] W. Ben Gunawan, "Revisiting the Sustainable Development Goal 4 'Quality Education': Insights, Prospects, and Recommendations," *SAKAGURU: Journal of Pedagogy and Creative Teacher*, vol. 2, no. 1, p. 12–36, May. 2025, doi: [10.70211/sakaguru.v2i1.202](https://doi.org/10.70211/sakaguru.v2i1.202).
- [6] M. Alabadla *et al.*, "Systematic Review of Using Machine Learning in Imputing Missing Values," *IEEE Access*, vol. 10, pp. 44483–44502, 2022, doi: [10.1109/ACCESS.2022.3160841](https://doi.org/10.1109/ACCESS.2022.3160841).
- [7] S. M. Piryonesi and T. E. El-Diraby, "Data Analytics in Asset Management: Cost-Effective Prediction of the Pavement Condition Index," *J. Infrastruct. Syst.*, vol. 26, no. 1, p. 04019036, Mar. 2020, doi: [10.1061/\(ASCE\)IS.1943-555X.0000512](https://doi.org/10.1061/(ASCE)IS.1943-555X.0000512).
- [8] F. Kong, Z. Song, and Q. Liu, "The frontiers of intelligent health services: cardiovascular disease prediction using novel machine learning methods and metaheuristic algorithm," *Comput. Methods Biomech. Biomed. Engin.*, pp. 1–19, May 2025, doi: [10.1080/10255842.2025.2502823](https://doi.org/10.1080/10255842.2025.2502823).
- [9] H. Hewamalage, K. Ackermann, and C. Bergmeir, "Forecast evaluation for data scientists: common pitfalls and best practices," *Data Min. Knowl. Discov.*, vol. 37, no. 2, pp. 788–832, Mar. 2023, doi: [10.1007/s10618-022-00894-5](https://doi.org/10.1007/s10618-022-00894-5).
- [10] A. S. Tejani, Y. S. Ng, Y. Xi, and J. C. Rayan, "Understanding and Mitigating Bias in Imaging Artificial Intelligence," *RadioGraphics*, vol. 44, no. 5, p. 13, May 2024, doi: [10.1148/rg.230067](https://doi.org/10.1148/rg.230067).
- [11] T. T. Khoei and A. Singh, "Data reduction in big data: a survey of methods, challenges and future directions," *Int. J. Data Sci. Anal.*, vol. 20, no. 3, pp. 1643–1682, Sep. 2025, doi: [10.1007/s41060-024-00603-z](https://doi.org/10.1007/s41060-024-00603-z).

- [12] S. N. P. Sreeramana Aithal, Shubhrajyotsna Aithal, "Future of Higher Education through Technology Prediction and Forecasting," ResearchGate. Mar. 02, 2026, doi: [10.5281/zenodo.11903348](https://doi.org/10.5281/zenodo.11903348).
- [13] A. A. Wani and F. Abeer, "Application of machine learning techniques for warfarin dosage prediction: a case study on the MIMIC-III dataset," *PeerJ Comput. Sci.*, vol. 11, p. e2612, Jan. 2025, doi: [10.7717/peerj-cs.2612](https://doi.org/10.7717/peerj-cs.2612).
- [14] M. Afkanpour, D. Tehrani Dehkordy, M. Momeni, and H. Tabesh, "Conceptual framework as a guide to choose appropriate imputation method for missing values in a clinical structured dataset," *BMC Med. Res. Methodol.*, vol. 25, no. 1, p. 43, Feb. 2025, doi: [10.1186/s12874-025-02496-3](https://doi.org/10.1186/s12874-025-02496-3).
- [15] D. Adhikari *et al.*, "A Comprehensive Survey on Imputation of Missing Data in Internet of Things," *ACM Comput. Surv.*, vol. 55, no. 7, pp. 1–38, Jul. 2023, doi: [10.1145/3533381](https://doi.org/10.1145/3533381).
- [16] I. D. Mienye and T. G. Swart, "A Comprehensive Review of Deep Learning: Architectures, Recent Advances, and Applications," *Information*, vol. 15, no. 12, p. 755, Nov. 2024, doi: [10.3390/info15120755](https://doi.org/10.3390/info15120755).
- [17] M. Y. Shakor and M. Ibrahim Khaleel, "Modern Deep Learning Techniques for Big Medical Data Processing in Cloud," *IEEE Access*, vol. 13, pp. 62005–62028, 2025, doi: [10.1109/ACCESS.2025.3556327](https://doi.org/10.1109/ACCESS.2025.3556327).
- [18] S. M. Alhammad, M. M. Eid, E. A. Mattar, and E.-S. M. El-Kenawy, "Optimization-Driven Learning for Leakage-Controlled Geospatial Modeling of Antenna Structure Registration Data," *IEEE Access*, vol. 14, pp. 15273–15310, 2026, doi: [10.1109/ACCESS.2026.3657224](https://doi.org/10.1109/ACCESS.2026.3657224).
- [19] H. Lee, D. Kim, H. Cho, G. Song, and J. Yoon, "Evaluation of data imputation models for building-integrated photovoltaic systems with practical performance and reproducibility," *Sol. Energy*, vol. 308, no. April, p. 114428, Apr. 2026, doi: [10.1016/j.solener.2026.114428](https://doi.org/10.1016/j.solener.2026.114428).
- [20] S. Zahmatkesh and P. Zech, "Spatio-Temporal Missing Data Imputation: A Systematic Literature Review with a Focus on Statistical and Machine Learning-Based Approaches," *ACM Comput. Surv.*, vol. 0, p. 37, Feb. 2026, doi: [10.1145/3797903](https://doi.org/10.1145/3797903).
- [21] F. Wunderlich *et al.*, "Assessing machine learning and data imputation approaches to handle the issue of data sparsity in sports forecasting," *Mach. Learn.*, vol. 114, no. 2, p. 48, Feb. 2025, doi: [10.1007/s10994-024-06651-7](https://doi.org/10.1007/s10994-024-06651-7).
- [22] A. Vehtari, A. Gelman, D. Simpson, B. Carpenter, and P.-C. Bürkner, "Rank-Normalization, Folding, and Localization: An Improved  $\hat{R}$  for Assessing Convergence of MCMC (with Discussion)," *Bayesian Anal.*, vol. 16, no. 2, pp. 667–718, Jun. 2021, doi: [10.1214/20-BA1221](https://doi.org/10.1214/20-BA1221).
- [23] H. Anahideh, P. Haghghat, N. Nezami, and D. G`andara, "Auditing the Imputation Effect on Fairness of Predictive Analytics in Higher Education," in *Computers and Society (cs.CY)*, Dec. 2022, p. 48. doi: [10.48550/arXiv.2109.07908](https://doi.org/10.48550/arXiv.2109.07908).
- [24] J. Koehler and C. Kuenzer, "Forecasting Spatio-Temporal Dynamics on the Land Surface Using Earth Observation Data—A Review," *Remote Sens.*, vol. 12, no. 21, p. 3513, Oct. 2020, doi: [10.3390/rs12213513](https://doi.org/10.3390/rs12213513).
- [25] M. Dhilsath Fathima, R. Hariharan, and S. P. Raja, "Multiple Imputation by Chained Equations– K - Nearest Neighbors and Deep Neural Network Architecture for Kidney Disease Prediction," *Int. J. Image Graph.*, vol. 23, no. 02, Mar. 2023, doi: [10.1142/S0219467823500146](https://doi.org/10.1142/S0219467823500146).
- [26] S. van Buuren and K. Groothuis-Oudshoorn, "mice : Multivariate Imputation by Chained Equations in R," *J. Stat. Softw.*, vol. 45, no. 3, p. 67, 2011, doi: [10.18637/jss.v045.i03](https://doi.org/10.18637/jss.v045.i03).
- [27] M. Abdelsattar, M. A. Azim, A. AbdelMoety, and A. Emad-Eldeen, "Comparative analysis of deep learning architectures in solar power prediction," *Sci. Rep.*, vol. 15, no. 1, p. 31729, Aug. 2025, doi: [10.1038/s41598-025-14908-x](https://doi.org/10.1038/s41598-025-14908-x).
- [28] M. Zhang, R. Zhao, C. Wang, L. Jing, and D. Li, "Real-Time Imputation Model for Missing Sensor Data Based on Alternating Attention Mechanism," *IEEE Sens. J.*, vol. 25, no. 5, pp. 8962–8974, Mar. 2025, doi: [10.1109/JSEN.2024.3519370](https://doi.org/10.1109/JSEN.2024.3519370).
- [29] S. Dhanka, A. Sharma, A. Kumar, S. Maini, and H. Vundavilli, "Advancements in Hybrid Machine Learning Models for Biomedical Disease Classification Using Integration of Hyperparameter-Tuning and

- Feature Selection Methodologies: A Comprehensive Review,” *Arch. Comput. Methods Eng.*, vol. 33, no. 1, pp. 289–324, Jan. 2026, doi: [10.1007/s11831-025-10309-5](https://doi.org/10.1007/s11831-025-10309-5).
- [30] P. C. Chiu, A. Selamat, O. Krejcar, K. K. Kuok, S. D. A. Bujang, and H. Fujita, “Missing Value Imputation Designs and Methods of Nature-Inspired Metaheuristic Techniques: A Systematic Review,” *IEEE Access*, vol. 10, pp. 61544–61566, 2022, doi: [10.1109/ACCESS.2022.3172319](https://doi.org/10.1109/ACCESS.2022.3172319).
- [31] S. Batool, J. Rashid, M. W. Nisar, J. Kim, H.-Y. Kwon, and A. Hussain, “Educational data mining to predict students’ academic performance: A survey study,” *Educ. Inf. Technol.*, vol. 28, no. 1, pp. 905–971, Jan. 2023, doi: [10.1007/s10639-022-11152-y](https://doi.org/10.1007/s10639-022-11152-y).
- [32] M. Afkanpour, E. Hosseinzadeh, and H. Tabesh, “Identify the most appropriate imputation method for handling missing values in clinical structured datasets: a systematic review,” *BMC Med. Res. Methodol.*, vol. 24, no. 1, p. 188, Aug. 2024, doi: [10.1186/s12874-024-02310-6](https://doi.org/10.1186/s12874-024-02310-6).
- [33] A. P. Wibawa, A. B. P. Utama, H. Elmunsyah, U. Pujiyanto, F. A. Dwiyanto, and L. Hernandez, “Time-series analysis with smoothed Convolutional Neural Network,” *J. Big Data*, vol. 9, no. 1, p. 44, Dec. 2022, doi: [10.1186/s40537-022-00599-y](https://doi.org/10.1186/s40537-022-00599-y).
- [34] A. P. Wibawa, “Mean-Median Smoothing Backpropagation Neural Network to Forecast Unique Visitors Time Series of Electronic Journal,” *Journal of Applied Data Sciences*, vol. 4, no. 3, pp. 163–174, Sep. 2023, doi: [10.47738/jads.v4i3.97](https://doi.org/10.47738/jads.v4i3.97).
- [35] J. Yang, Y. Wang, Y. Yang, K. Ding, C. Na, and Y. Yang, “Effects of single and multiple imputation strategies on addressing over-fitting issues caused by imbalanced data from various scenarios,” *Appl. Intell.*, vol. 54, no. 3, pp. 2812–2830, Feb. 2024, doi: [10.1007/s10489-024-05295-3](https://doi.org/10.1007/s10489-024-05295-3).
- [36] G. S. Ramnath, R. Harikrishnan, S. M. Muyeen, A. Kukker, S. D. Pohekar, and K. Kotecha, “A peer-and self-group competitive behavior-based socio-inspired approach for household electricity conservation,” *Sci. Rep.*, vol. 14, no. 1, p. 17245, Jul. 2024, doi: [10.1038/s41598-024-56926-1](https://doi.org/10.1038/s41598-024-56926-1).
- [37] A. P. Wibawa *et al.*, “Deep Learning Approaches with Optimum Alpha for Energy Usage Forecasting,” *Knowledge Engineering and Data Science*, vol. 6, no. 2, p. 5, Oct. 2023, doi: [10.17977/um018v6i22023p170-187](https://doi.org/10.17977/um018v6i22023p170-187).
- [38] A. W. Saputra, A. P. Wibawa, U. Pujiyanto, A. B. Putra Utama, and A. Nafalski, “LSTM-based Multivariate Time-Series Analysis: A Case of Journal Visitors Forecasting,” *ILKOM Jurnal Ilmiah*, vol. 14, no. 1, pp. 57–62, Apr. 2022, doi: [10.33096/ilkom.v14i1.1106.57-62](https://doi.org/10.33096/ilkom.v14i1.1106.57-62).
- [39] A. B. Putra Utama, A. P. Wibawa, M. Muladi, and A. Nafalski, “PSO based Hyperparameter tuning of CNN Multivariate Time- Series Analysis,” *Jurnal Online Informatika*, vol. 7, no. 2, pp. 193–202, 2022, doi: [10.15575/join.v7i2.858](https://doi.org/10.15575/join.v7i2.858).
- [40] W. Zhou, Z. Yan, and L. Zhang, “A comparative study of 11 non-linear regression models highlighting autoencoder, DBN, and SVR, enhanced by SHAP importance analysis in soybean branching prediction,” *Sci. Rep.*, vol. 14, no. 1, p. 5905, Mar. 2024, doi: [10.1038/s41598-024-55243-x](https://doi.org/10.1038/s41598-024-55243-x).
- [41] A. Pak, A. K. Rad, M. J. Nematollahi, and M. Mahmoudi, “Application of the Lasso regularisation technique in mitigating overfitting in air quality prediction models,” *Sci. Rep.*, vol. 15, no. 1, p. 547, Jan. 2025, doi: [10.1038/s41598-024-84342-y](https://doi.org/10.1038/s41598-024-84342-y).
- [42] P. N. Sharma, G. Shmueli, M. Sarstedt, N. Danks, and S. Ray, “Prediction-Oriented Model Selection in Partial Least Squares Path Modeling,” *Decis. Sci.*, vol. 52, no. 3, pp. 567–607, Jun. 2021, doi: [10.1111/dec.12329](https://doi.org/10.1111/dec.12329).
- [43] E. Afrifa-Yamoah, U. A. Mueller, S. M. Taylor, and A. J. Fisher, “Missing data imputation of high-resolution temporal climate time series data,” *Meteorol. Appl.*, vol. 27, no. 1, p. e1873, Jan. 2020, doi: [10.1002/met.1873](https://doi.org/10.1002/met.1873).
- [44] J. S. Prince, I. Charest, J. W. Kurzawski, J. A. Pyles, M. J. Tarr, and K. N. Kay, “Improving the accuracy of single-trial fMRI response estimates using GLMsingle,” *Elife*, vol. 11, p. 28, Nov. 2022, doi: [10.7554/eLife.77599](https://doi.org/10.7554/eLife.77599).
- [45] J. S. Joswig *et al.*, “Imputing missing data in plant traits: A guide to improve gap-filling,” *Glob. Ecol. Biogeogr.*, vol. 32, no. 8, pp. 1395–1408, Aug. 2023, doi: [10.1111/geb.13695](https://doi.org/10.1111/geb.13695).

- [46] J. Zhu, X. Zhao, Y. Sun, S. Song, and X. Yuan, "Relational Data Cleaning Meets Artificial Intelligence: A Survey," *Data Sci. Eng.*, vol. 10, no. 2, pp. 147–174, Jun. 2025, doi: [10.1007/s41019-024-00266-7](https://doi.org/10.1007/s41019-024-00266-7).
- [47] H. Karnati, A. Soma, A. Alam, and B. Kalaavathi, "Comprehensive analysis of various imputation and forecasting models for predicting PM2.5 pollutant in Delhi," *Neural Comput. Appl.*, vol. 37, no. 17, pp. 11441–11458, Jun. 2025, doi: [10.1007/s00521-025-11047-2](https://doi.org/10.1007/s00521-025-11047-2).
- [48] V. V. Golovko, "Robust Method for Confidence Interval Estimation in Outlier-Prone Datasets: Application to Molecular and Biophysical Data," *Biomolecules*, vol. 15, no. 5, p. 704, May 2025, doi: [10.3390/biom15050704](https://doi.org/10.3390/biom15050704).
- [49] Y. S. Mohammed, H. Abdelkader, P. Pławiak, and M. Hammad, "A novel model to optimize multiple imputation algorithm for missing data using evolution methods," *Biomed. Signal Process. Control*, vol. 76, no. July, p. 103661, Jul. 2022, doi: [10.1016/j.bspc.2022.103661](https://doi.org/10.1016/j.bspc.2022.103661).
- [50] S. Mpfu and D. Chasokela, "Data-Informed Decision-Making," IGI Global Scientific Publishing, Nov. 2024, pp. 103–138. doi: [10.4018/979-8-3693-6967-8.ch004](https://doi.org/10.4018/979-8-3693-6967-8.ch004).
- [51] I. M. Wirawan, A. P. Wibawa, and T. Widiyanintyas, "Photovoltaic Energy Anomaly Detection using Transformer Based Machine Learning," *International Journal of Robotics and Control Systems*, vol. 4, no. 3, pp. 1337–1352, Aug. 2024, doi: [10.31763/ijrcs.v4i3.1260](https://doi.org/10.31763/ijrcs.v4i3.1260).
- [52] D. V. Ogunkan and S. K. Ogunkan, "Exploring big data applications in sustainable urban infrastructure: A review," *Urban Gov.*, vol. 5, no. 1, pp. 54–68, Mar. 2025, doi: [10.1016/j.ugj.2025.02.003](https://doi.org/10.1016/j.ugj.2025.02.003).
- [53] X. Jiang, Y. Yao, S. Liu, F. Shen, L. Nie, and X.-S. Hua, "Dual Dynamic Threshold Adjustment Strategy," *ACM Trans. Multimed. Comput. Commun. Appl.*, vol. 20, no. 7, pp. 1–18, Jul. 2024, doi: [10.1145/3656047](https://doi.org/10.1145/3656047).
- [54] B. I. Chigbu and S. L. Makapela, "Data-Driven Leadership in Higher Education: Advancing Sustainable Development Goals and Inclusive Transformation," *Sustainability*, vol. 17, no. 7, p. 3116, Apr. 2025, doi: [10.3390/su17073116](https://doi.org/10.3390/su17073116).