# A REVIEW OF CHALLENGES AND SOLUTIONS FOR USING MACHINE LEARNING APPROACHES FOR MISSING DATA

Aasim Ayaz Wani
Department of Engineering
Cornell University, Ithaca, NY, USA

*Abstract—* **Missing data poses significant challenges to the reliability of statistical analyses and predictive modeling across diverse research fields. This paper provides an in-depth review of both traditional and machine learning imputation techniques, enabling researchers to navigate the complexities of missing data with greater efficacy. We evaluate simple imputation methods, such as mean, median, and mode, and delve into more sophisticated strategies including regression-based, hot and cold deck, and probabilistic models like Gaussian Mixture Models and K-Nearest Neighbors. Furthermore, the paper explores cutting-edge machine learning approaches like Random Forest, Multiple Imputation by Chained Equations, and deep learning models such as autoencoders and Generative Adversarial Networks. Our comprehensive analysis highlights the effectiveness of each method, tailored to various missing data mechanisms MCAR, MAR, and NMAR providing actionable insights for researchers to enhance data integrity and improve the outcomes of their studies.**

*Keywords—* **Missing Data Imputation, Machine Learning, Data Integrity, Predictive Modeling, Deep Learning, GANs, KNN, Random Forest, Autoencoders, Gaussian Mixture Models.**

## I. INTRODUCTION

The ubiquity of missing data in research datasets presents a persistent challenge that compromises the integrity of statistical analyses and predictive modeling. Missing data can arise from various sources such as non-response in surveys, errors in data collection, or unforeseen disruptions in data transmission Alam et al. (2023). The nature of missingness, whether random or systematic, can significantly influence the bias and variance of the estimates derived from such datasets. Thus, managing missing data through appropriate imputation methods is crucial to ensure the validity of research findings Newman (2014).

In this paper, we embark on a comprehensive survey of imputation methods tailored to address the complexities introduced by missing data in diverse datasets. Our exploration encompasses a range of traditional and advanced techniques, each suited to different types of missing data mechanisms—Missing Completely at Random (MCAR), Missing at Random (MAR), and Not Missing at Random (NMAR). We commence our discussion by delineating the conceptual framework that underpins missing data theory, highlighting the implications of each missing data mechanism on the validity of statistical inferences. Subsequently, we delve into simple imputation methods like mean, median, and mode imputation, which, despite their simplicity, are often limited by their assumption of randomness in missing data patterns.

To address the limitations of simple methods, we explore sophisticated techniques such as hot deck and cold deck imputation, regression imputation, and probabilistic models like Gaussian Mixture Models and KNN imputation. These methods offer nuanced approaches that account for the underlying relationships within the data, enhancing the accuracy of imputations under certain conditions. Moreover, we investigate the role of machine learning in imputation, with a focus on Random Forest and MICE. These methods utilize the inherent data patterns and correlations to generate more reliable and robust imputations, particularly when dealing with MAR data.

Our survey also extends to cutting-edge techniques such as imputation using deep learning models like Autoencoders and Generative Adversarial Networks, which represent the frontier of imputation methodology. These models are particularly adept at handling complex and high-dimensional data, offering innovative solutions that traditional methods may not provide. Through this detailed examination, our paper aims to equip researchers with the knowledge to select and apply the most appropriate imputation techniques, thereby enhancing the reliability of their analyses in the face of incomplete data.

This paper is structured as follows: In Section 2, we provide a conceptual framework of missing data mechanisms, elaborating on MCAR, MAR, and NMAR. Section 3 reviews traditional imputation methods, including mean, median, and mode imputation. Section 4 explores more imputation techniques such as hot deck, cold deck, and regression imputation, as well as probabilistic models like

GMM and KNN. Section 5 delves into machine learning-based imputation methods, focusing on Random Forest and MICE. Section 6 discusses cutting-edge imputation techniques involving deep learning models such as autoencoders and GANs. Finally, Section 7 concludes with a discussion on future research directions and the implications of our findings for the field of data imputation.

## 1.1 Rationale and Audience

In the era of data-driven decision-making, the integrity of data underpins the reliability of research outcomes across various disciplines; missing data is a pervasive challenge that compromises the accuracy of statistical analyses and decision-making processes, highlighting a critical need for robust imputation methods Pansara (2023). This study systematically evaluates both established and innovative machine learning-based imputation techniques to address the gaps in traditional methods that often fail to accommodate the nuances of modern datasets. Our research aims to enhance the reliability of research findings in data-rich environments, bridging the gap between theoretical models and practical applications. This comprehensive survey targets a diverse array of professionals grappling with incomplete data, including academic researchers, data scientists, and statisticians across fields such as biology, economics, and computer science. Additionally, industry practitioners who rely on accurate data for informed decision-making will find this analysis invaluable. By detailing the practical applications and limitations of each imputation technique, this paper serves as an essential resource for anyone tasked with ensuring data integrity and making informed decisions based on analysis

## 1.2. Search Methodology

We conducted a systematic search for imputation methods using Google Scholar focusing on terms like "mean imputation," "median imputation," "K-Nearest Neighbors imputation," and "autoencoder imputation." Boolean operators (AND, OR) refined the queries to target studies published in the last 15 years, peer-reviewed, and written in English. Papers not focused on imputation, older than 15 years unless seminal, non-peer-reviewed, or in other languages were excluded. After an initial screening of titles and abstracts, we reviewed full texts of relevant articles and checked their references for additional studies.

## 1.3 Family of Missing Values

Missing data is categorized into three primary mechanisms:

- Missing Completely at Random (MCAR): MCAR occurs when the missingness of data is independent of both observed and unobserved variables Rubin (1976). In this scenario, the absence of data does not introduce systematic bias into the analysis, allowing standard statistical methods to be applied without special adjustments for handling missingness. An example of MCAR is when laboratory samples are randomly lost due to logistical errors, meaning the missing data is unrelated to any patient characteristics or health conditions. However, while MCAR does not bias the results, extensive missing data under this mechanism can still reduce statistical power and diminish the effective sample size, potentially limiting the robustness and generalizability of the study's conclusions.
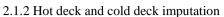
- Missing at Random (MAR): MAR describes a situation where the probability of missingness is systematically related to the observed data but is independent of the missing values themselves, conditional on the observed data Rubin (1976). This allows for the use of advanced imputation methods that leverage relationships between fully observed variables to estimate the missing data. A common example of MAR would be younger patients being less likely to complete certain clinical tests. In such cases, methods like Multiple Imputation by Chained Equations (MICE) can use the correlations among observed predictors to handle the missing data, thereby reducing bias and improving the accuracy of the analysis.

- Not Missing at Random (NMAR): NMAR occurs when the missingness is related to the unobserved data values themselves Rubin (1976). In this challenging scenario, the missing data mechanism depends on information that is not available, requiring assumptions about the relationship between the likelihood of missingness and the missing data itself. For instance, in clinical studies, patients with more severe symptoms may be less likely to report all their symptoms, leading to non-random missingness of critical health information. Effectively handling NMAR is crucial to avoid biased conclusions, especially in clinical trials and observational studies where the accuracy of conclusions can directly influence clinical decisions and policy-making

## II. SURVEY OF MISSING VALUE METHODS

### 2.1 Traditional Methods

#### 2.1.1 Simple Imputation Methods

This includes mean, median, and mode imputation, popular for their simplicity and computational efficiency Jadhav et al. (2019). These methods replace missing values with central tendency measures of observed data. They are effective under the MCAR assumption, however, when applied to MAR or NMAR data, they can introduce bias and distort data distribution, leading to inaccurate estimates and overly confident inferences Buczak et al. (2023). This distortion can significantly impact analyses like regression models or predictive algorithms, where maintaining the original data distribution is crucial. Despite these limitations this methods can be useful for preliminary analyses if their assumptions are carefully evaluated.

2.1.2 Hot deck and cold deck imputation
Hot deck and cold deck imputation are practical techniques for handling missing data using observed values, though they differ in source utilization. Hot deck imputation selects values from similar records within the same dataset, defined by key characteristics such as age or education, making it effective when clear patterns exist in the data. However, inaccuracies arise if donors are poorly matched, and this method underestimates standard errors, leading to biased inferences. It works well for both continuous and categorical data but struggles with complex data structures. Cold deck imputation, by contrast, uses external historical datasets to fill missing values. This method is effective when external data remains valid and reflective of current conditions, such as using prior years' test scores in longitudinal studies. However, it risks inaccuracies if the historical data doesn't align with contemporary situations. Both methods assume data is MCAR or MAR, maintaining original distributions and relationships, but they perform poorly under NMAR conditions, leading to bias.

2.1.3 Forward and Backward Fill Imputation
Forward and backward fill imputations (FF & BF) are simple propagation techniques for handling missing data in time-series datasets by filling gaps with the nearest observed value. Forward fill replaces missing values with the most recent preceding value, while backward fill uses the next available value Ahn et al. (2022). These methods assume a MCAR or MAR pattern and maintain temporal continuity, making them effective for datasets with gradual trends Van Ginkel et al. (2020). However, they can introduce biases in rapidly changing or non-sequential data, especially when MNAR is present Molenberghs et al. (2014). Additionally, FF and BF can propagate errors from earlier data points, inflate variance by repeating values, and struggle with long sequences of missing data Kenward and Molenberghs (2009). Though useful for time-series data with gradual changes, they may not be suitable for more complex datasets, in which case methods like mean/median, hot deck, or model-based imputation are preferred Esling and Agon (2012).

2.2 Machine Learning Methods
2.2.1 Gaussian Mixture Models Imputation
GMMs are used for imputing missing values by modeling data as a mixture of several Gaussian distributions, each characterized by its own mean and covariance Asheri et al. (2021). At the core of this method is the Expectation-Maximization (EM) algorithm, an iterative procedure designed to find maximum likelihood estimates in datasets with incomplete data. The EM algorithm addresses missing data by alternating between two key steps. In the Expectation Step (E-step), the algorithm estimates the missing data using the observed data and the current parameter estimates. For each incomplete data point, the E-

step computes the expected value of the missing component based on the current model parameters. By leveraging the observed data, the E-step statistically infers the missing values, enhancing the imputation process Enders (2022). GMMs assume that both the observed and missing data are generated from a mixture of underlying Gaussian distributions. This assumption greatly impacts the success of the imputation, especially in datasets with multimodal distributions. An essential consideration is the choice of the number of Gaussian components, k. If too few components are selected, the model may oversimplify the data, failing to capture its complexity. On the other hand, too many components can lead to overfitting, particularly in datasets with limited data. Techniques such as the Bayesian Information Criterion (BIC) or Akaike Information Criterion (AIC) are frequently used to determine the optimal number of components Steele and Raftery (2010). The computational complexity of the EM algorithm in GMMs depends on both the number of components and the dimensionality of the data. Each iteration of the algorithm has a time complexity of $O(nkd + nk + kd2)$, where n is the number of data points, d is the number of features, and k is the number of components Asheri et al. (2021). The memory complexity includes storing the dataset $(O(nd))$ and the model parameters $(O(kd + kd2))$, resulting in a total memory requirement of $O(nd + kd + kd2)$.

GMMs are particularly effective for datasets with complex, multimodal distributions, as they can model multiple distinct subpopulations, each following a Gaussian distribution. This makes GMMs well-suited for heterogeneous data that contains different clusters with varying means and covariances. GMMs are also effective in handling high-dimensional datasets where traditional imputation methods may fail to capture intricate relationships between variables. Their ability to manage MAR or MCAR data using observed data patterns further underscores their suitability for imputing missing values in complex, high-dimensional, and heterogeneous datasets Cho et al. (2020).

GMM imputation ensures computational efficiency and accuracy through several stopping criteria, including a maximum number of iterations, monitoring log-likelihood changes, and checking for minimal changes in model parameters (e.g., means and covariances) Steele and Raftery (2010). These criteria balance computational resources and imputation quality, allowing the process to terminate efficiently while maintaining accuracy. By modeling data probabilistically, GMMs offer robust and reliable imputations, making them well-suited for managing missing data in complex datasets. Although computationally intensive and not always guaranteed to converge to a global optimum, the EM algorithm's precision makes GMMs a valuable tool in diverse applications. After the E-step, the Maximization Step (M-step) updates the means, covariances, and mixture weights to maximize the

likelihood of the observed and estimated missing data. These steps are repeated until convergence, typically determined by minimal changes in log-likelihood or parameter stabilization. In case of a two-dimensional image, after a DWT transform, the image is divided into four corners, upper left corner of the original image, lower left corner of the vertical details, upper right corner of the horizontal details, lower right corner of the component of the original image detail (high frequency). You can then continue to the low frequency components of the same upper left corner of the 2nd, 3rd inferior wavelet transform.

2.2.2 K-Nearest Neighbors (kNN)
KNN imputation substitutes missing values by leveraging the K most similar instances in the dataset Triguero et al. (2019). This method calculates similarity using a distance metric, such as Euclidean or Manhattan distance, and then imputes missing values based on the mean, median, or mode of the neighboring values Santos et al. (2020). The success of KNN imputation heavily depends on the proper selection of the distance metric and the number of neighbors (K), which must strike a balance between capturing local data patterns and avoiding the introduction of noise or over-smoothing. For example, Euclidean distance is often ideal for continuous data, while Manhattan distance may better accommodate categorical data. Moreover, the choice of K significantly impacts how well local patterns are preserved—smaller values of K may be more sensitive to noise Zhang (2012). KNN imputation excels at preserving local structures and is especially useful in datasets where similar instances are predictive of similar outcomes Cho et al. (2020). It assumes that the missing data mechanism is either MAR (Missing at Random) or MCAR (Missing Completely at Random), allowing for accurate inferences based on observed patterns in the data. However, KNN imputation requires substantial computational resources due to the need for calculating distances and retrieving neighbors. The computational complexity is given by $O(m \cdot (nd + n \log n + k))$, where n represents the number of samples, d is the number of features, and m is the number of missing values Beretta and Santaniello (2016). This, combined with the memory requirement $O(nd + n)$, makes KNN a resource-intensive approach.
Stopping criteria in KNN imputation are often based on performance evaluations, with the process typically concluding once the predefined number of neighbors (K) has been reached. Performance metrics such as mean squared error or accuracy are also commonly used to determine when further improvement has plateaued. In practice, computational resources and time constraints often dictate the stopping point for the imputation process. Despite its computational demands, optimizations like KD-trees have been introduced to mitigate the burden of processing Muja and Lowe (2014). Nonetheless, KNN imputation's effectiveness depends on careful selection of

the distance metric and K; poor choices can result in biased or inaccurate imputations, particularly in large datasets or those prone to outliers. Despite these challenges, KNN remains a flexible and reliable imputation method when properly tuned to the characteristics of the dataset and the available computational resources

2.2.3 Random Forest (RF) imputation
RF imputation leverages the power of ensemble learning through multiple decision trees to effectively handle missing data Tang and Ishwaran (2017). In this method, a large number of trees are constructed using different subsets of the available data, and the imputed value for each missing entry is determined based on the consensus prediction across all trees. This collective approach enhances the robustness and reliability of the imputations, especially in the presence of complex, nonlinear relationships and interactions between continuous and categorical variables Plaia et al. (2022). RF imputation is particularly useful in cases where the missing data mechanism is MAR or MCAR, as it can accommodate various data types and complex structures.
One key feature of RF imputation is that it typically does not rely on formal convergence criteria. Instead, the imputation process is terminated after a predetermined number of iterations, guided by performance metrics like mean squared error or by computational limits. While iterative imputation can lead to improved results, the optimal number of iterations often depends on the specific dataset and requires empirical validation. It is important to note that RF imputation assumes the missing data is MAR, which can introduce bias if the data is MNAR. Additionally, this method is computationally intensive, particularly for large datasets with significant amounts of missing data Kokla et al. (2019). The accuracy of imputations heavily relies on the RF model's ability to correctly identify important features, and any errors in feature selection can compromise the quality of the imputed values. Moreover, when too many trees are used relative to the size of the dataset, RF models may overfit, leading to overly complex models that generalize poorly Shah et al. (2014). Unlike multiple imputation methods, RF typically generates a single imputed dataset, potentially underestimating the uncertainty in the imputed values. The computational demands of RF imputation are substantial, with the complexity of training a random forest involving n samples, d features, and t trees amounting to $O(t \cdot n \log n \cdot d)$. This complexity arises primarily from sorting operations during bootstrap sampling. Furthermore, each missing value must be passed through all t trees, resulting in an imputation complexity of $O(t \cdot d)$ per missing value. For a dataset with m missing values, the overall complexity becomes $O(t \cdot n \log n \cdot d + m \cdot t \cdot d)$. Memory-wise, RF imputation is also demanding, requiring storage for both the dataset $(O(nd))$

and the decision trees (O(t ·n log n)) Rahman and Islam (2013).

**2.2.4 Multivariate Imputation by Chained Equations**
Multivariate Imputation by Chained Equations (MICE) employs an iterative regression technique to address missing data, creating multiple imputations Van Buuren and Oudshoorn (2000). The process models each variable with missing entries as a dependent function of other variables within the dataset. Predictions from these regression models are then used to impute missing values. This iterative cycle is repeated, with each round refining the imputations based on updated model estimations, until convergence is achieved Azur et al. (2011). MICE generates multiple complete datasets, capturing the inherent uncertainty of the missing data, which enables robust statistical inference. Each dataset is independently analyzed, and the results are aggregated to yield comprehensive statistical estimation Van Buuren and Groothuis-Oudshoorn (2011). This method operates under the assumption that the data is MAR or MCAR, proving particularly effective in preserving the statistical characteristics of the original dataset, such as variability and inter-variable correlations.

Despite these benefits, MICE has several limitations. If the data is MNAR, MICE may produce biased imputations. The computational demands are considerable, largely due to its iterative and multivariate nature. The complexity for imputing data across M complete datasets through i iterations for d features with n samples calculates to $O(M \cdot i \cdot d \cdot nd2)$, reflecting the repeated fitting of regression models and their application in imputation across all datasets. Additionally, the memory requirements include storing these M imputed datasets, accumulating to $O(M \cdot nd)$ Beesley and Taylor (2021). Convergence in MICE refers to the stabilization point in the imputation process where subsequent iterations no longer significantly alter the imputed values. Common criteria for assessing convergence include: (1) Maximum Iterations: A set number of imputation cycles are performed, prioritizing computational efficiency over precision; (2) Change in Imputed Values: Monitoring the absolute difference between imputed values across iterations, assuming convergence when this difference falls below a specified threshold, balancing efficiency and accuracy; and (3) Convergence Diagnostics: Utilizing sophisticated statistical measures like log-likelihood or hypothesis tests, though these methods are computationally intensive and often impractical for routine use. It's important to note that convergence alone doesn't guarantee the quality of the imputed data, as factors such as the suitability of imputation models and the complexity of missing data patterns also play crucial roles. Therefore, combining multiple convergence criteria is recommended to enhance both efficiency and accuracy.

The accuracy of the imputation model depends on the correct specification of the relationships between variables, and incorrect model specification can lead to biased imputations. The choice of imputation method within MICE can significantly influence the results, and the method is sensitive to the quality of observed data; errors or outliers in the observed data will propagate into the imputed values Aleryani et al. (2020). Despite its resource-intensive nature, MICE stands out as a highly flexible and robust method, offering comprehensive solutions for missing data treatment in complex datasets, ultimately leading to more accurate and reliable statistical conclusions. Therefore, it is crucial to carefully assess the suitability of MICE for a given dataset and consider alternative methods if necessary.

**2.3 Deep Learning Methods**
**2.3.1 Autoencoders**
Autoencoder imputation utilizes the unique capabilities of autoencoders—a type of neural network optimized for learning efficient data representations—to address missing values Pinaya et al. (2020). An autoencoder comprises two primary components: an encoder which compresses input data into a lower-dimensional latent space, and a decoder, which reconstructs the data from this compressed form. This system is trained on incomplete datasets to discern underlying patterns and structures, making it highly effective for complex, high-dimensional datasets where traditional imputation methods falter Pereira et al. (2020). By reconstructing input data from the latent space, autoencoders can predict and fill in missing values, assuming the data is MCAR or MAR Ma et al. (2020). Autoencoders are adept at capturing complex, nonlinear relationships and interactions between variables, preserving the integrity of the original data. Their ability to learn a compressed representation in the latent space not only aids in imputation but also proves beneficial for tasks like visualization and anomaly detection. However, ensuring that this latent space retains essential information for accurate reconstruction requires careful management. The flexibility of autoencoders to adapt to various data types through appropriate network architectures and activation functions also demands meticulous selection of autoencoder variants to suit specific data characteristics.

Despite their strengths, autoencoders face several challenges. They are prone to reconstructing observed values while potentially overlooking the underlying distribution of missing values, leading to reconstruction bias. Their effectiveness heavily depends on the chosen architecture and hyperparameters, which can make the optimization process both time-consuming and resource-intensive. Autoencoders typically assume MAR data, and addressing MNAR data remains problematic Costa et al. (2018). The training of deep autoencoder models is computationally demanding, especially for large datasets, and can lead to underfitting—where the model fails to capture complex data patterns—or overfitting, where the model memorizes training data at the expense of general

performance. Moreover, autoencoders usually generate a single imputed value for each missing entry, disregarding the inherent uncertainty of the imputation process. Despite these limitations, autoencoder imputation remains a robust and flexible technique, offering realistic and accurate imputations that closely reflect the true characteristics of the dataset. This makes autoencoders an invaluable tool in the repertoire of methods for managing missing data in complex scenarios.

2.3.2 Generative Adversarial Networks (GANS)
Generative Adversarial Networks (GANs) imputation is a deep learning method for handling missing data that leverages the generative capabilities of GANs. A GAN consists of two neural networks: a generator and a discriminator. The generator creates synthetic data that mimics the real data, while the discriminator evaluates the authenticity of the data, distinguishing between real and generated samples. During imputation, the generator is trained to produce plausible values for the missing data by learning the underlying distribution of the observed data. The discriminator, on the other hand, assesses the quality of the generated imputations, guiding the generator to improve its outputs. This adversarial process continues until the generator produces realistic imputations that the discriminator cannot easily distinguish from the real data. GANs imputation assumes that the data is MAR or MCAR, making it suitable for various types of data, including high-dimensional and complex datasets. However, while GANs have shown promise, they come with several challenges. They are notoriously difficult to train, often suffering from instability issues such as mode collapse or vanishing gradients, which can significantly impact the quality of imputed values. They typically require large amounts of complete data to effectively learn the underlying data distribution, posing a challenge for datasets with high missingness rates. Mode collapse can result in a lack of diversity in the imputed values, and achieving optimal performance often requires careful tuning of numerous hyperparameters, which is time-consuming and computationally expensive. Evaluating the performance of GANs for imputation is difficult due to the absence of ground truth for the missing values. Additionally, handling MNAR data remains a challenge for GANs. The computational cost of training GANs, especially for large datasets, further limits their applicability in some cases. Despite these challenges, GANs imputation remains a powerful and flexible approach that offers a sophisticated solution for missing data, ensuring that the imputed values are both accurate and realistic.

### III. CHALLENGES

3.1 Handling Complex Data Structures
3.1.1 Heterogeneous Data

Real-world datasets often exhibit considerable heterogeneity, encompassing a mix of numerical, categorical, ordinal, and textual data. Imputing missing values in such diverse datasets presents unique challenges, as traditional imputation methods are typically optimized for single data types and may fall short when confronted with multiple data forms. If handled inadequately, imputation in heterogeneous datasets can lead to inconsistencies, introduce bias, or even distort the relationships among variables, ultimately compromising the quality and validity of subsequent analyses or predictive models.Numerical Data: For numerical data, common imputation methods such as mean, median, or k-nearest neighbors (KNN) imputation are often used. These methods exploit mathematical relationships between values, making them suitable for filling in missing continuous variables. However, these methods can break down when applied to non-numerical data types, underscoring the need for type-specific approaches.
Categorical Data: Categorical data, which consists of discrete categories without inherent order, requires a different approach. Techniques like mode imputation or hot deck imputation are more appropriate, as they preserve the categorical nature of the data. Mode imputation replaces missing values with the most frequent category, while hot deck imputation selects replacement values from similar records within the dataset. Both methods maintain the categorical structure but can be limited if the data has rare or underrepresented categories. Ordinal Data: Ordinal data, which carries a natural order but unequal intervals between values (such as rankings), requires an approach that respects this order. Simply applying categorical or numerical imputation methods risks distorting the rank relationships. Methods like ordinal regression or stratified imputation are better suited, as they maintain the ranking structure and avoid creating inconsistencies between the ordinal variables.
Text Data: Imputing missing text data introduces further complexity, as text data cannot be treated like numerical or categorical variables. Missing text often requires more advanced imputation methods grounded in natural language processing (NLP). Techniques such as word embeddings (e.g., Word2Vec, GloVe), transformer-based models (e.g., BERT, GPT-3), or topic modeling are utilized to infer missing text based on the context within surrounding words, sentences, or even documents. These methods allow for more accurate imputation by capturing the semantic meaning and contextual relationships within the text.
Handling Heterogeneity: A common approach to managing heterogeneous datasets is to apply different imputation methods tailored to each data type. For example, regression-based or KNN imputation might be used for numerical data, hot deck or mode imputation for categorical variables, ordinal regression for ordinal data, and NLP techniques for text. However, when imputation methods are applied independently to each data type, careful consideration is

necessary to avoid introducing artificial correlations or inconsistencies between the imputed values. For instance, failing to account for relationships between numerical and categorical data may result in incongruous imputed values that do not align with the rest of the dataset.

Advanced Methods: Advanced imputation techniques such as Multiple Imputation by Chained Equations (MICE) offer a solution for handling heterogeneous data by modeling each variable conditional on the others. MICE iteratively imputes missing values for each variable using the observed data from other variables, thus preserving the relationships across mixed data types. Another approach is Random Forest imputation, which naturally handles heterogeneous data types by constructing decision trees that can partition the data using the most informative splits, regardless of the data type.

Conclusion: Effectively handling heterogeneous data during imputation requires a flexible, adaptive approach that respects the specific characteristics of each data type. The chosen imputation method should be carefully matched to the data type—whether numerical, categorical, ordinal, or text—to minimize bias and maintain the accuracy of the dataset. By employing methods that account for the diversity and complexity of the data, practitioners can ensure that imputation not only fills in missing values but also preserves the integrity and relationships within the dataset, resulting in more reliable and robust analyses.

### 3.1.2 Interdependent Variables

In complex datasets, variables often exhibit intricate interdependencies, particularly in time-series and spatial data. These interdependencies pose substantial challenges during imputation, as simplistic methods that ignore the relationships between variables can disrupt the data's natural structure and distort downstream analyses. Imputing missing values without accounting for these interdependencies risks introducing significant biases, misrepresenting trends, and degrading the performance of predictive models. In time-series datasets, values at a given time point are frequently influenced by preceding observations, reflecting underlying temporal dynamics such as trends, seasonality, or autocorrelation. For instance, in economic or financial data, variables like stock prices or consumer spending often depend on prior values, making it critical to preserve this continuity during imputation. Using simple techniques such as mean imputation can sever these temporal relationships, resulting in artificially smoothed data that fails to capture important fluctuations or trends. More sophisticated approaches, like autoregressive imputation or incorporating lag variables, take the sequential nature of time-series data into account by using past observations as predictors for missing values. This ensures that imputed values align with the temporal structure of the dataset, thereby maintaining data integrity and improving model accuracy. Spatial datasets, which

often contain locational dependencies, present similar challenges. In spatial data, observations at proximate locations are typically more correlated than those at distant ones, as seen in environmental datasets where variables like temperature or pollution levels exhibit spatial continuity. Simple imputation methods, such as global averages, overlook these spatial dependencies, potentially creating unrealistic or spatially inconsistent estimates. Spatial imputation techniques like Kriging or spatial autoregressive models address this by incorporating spatial correlation structures during imputation. These methods account for the relationships between neighboring points, ensuring that the imputed values are consistent with the geographical patterns in the data. This not only preserves spatial integrity but also enhances the accuracy and reliability of the analysis.

Beyond time-series and spatial data, many other domains feature datasets with interdependent variables, such as multi-sensor data, medical time-series,and geospatial-temporal datasets. In such cases, preserving the relationships between variables during imputation is crucial for maintaining the dataset's coherence and ensuring that downstream models reflect the true underlying patterns. Ultimately, imputation methods that respect variable interdependencies—whether temporal, spatial, or even functional—are vital for producing accurate and reliable results. Without accounting for these complex relationships, imputation can introduce distortions that weaken model performance and compromise the validity of conclusions. Therefore, selecting imputation techniques that maintain interdependencies is essential in data-driven applications involving structured datasets, ensuring that the true nature of the data is preserved and faithfully reflected in analytical outcomes.

### 3.2 Bias in Imputation
### 3.2.1 Imbalance in Missing Data Patterns

Imputation in imbalanced datasets presents significant challenges that distinctly affect the outcomes and reliability of classification and regression tasks. Imbalanced datasets, where some classes or outcomes are significantly underrepresented, can also lead to overfitting. In such cases, models might overfit to the majority class or the more frequently observed outcomes because there is insufficient data to learn about the minority classes or less common outcomes accurately Ali et al. (2019). This can skew the imputation results, making them biased towards the dominant data points. Understanding the unique dynamics of these tasks is crucial for implementing effective imputation strategies that maintain data integrity and model accuracy. These risks stem from the inherent skewness in the data distribution, which poses unique challenges for accurately estimating missing values. Techniques such as resampling the data, using anomaly detection methods to identify and adjust for rare events, or employing cost-sensitive learning where the model pays more attention to

the minority class, are strategies that can help mitigate overfitting in these contexts.

In classification tasks, methods like mean, median, or mode imputation introduce significant bias in imbalanced datasets, disproportionately influenced by the majority class. These methods often skew imputations, failing to adequately represent minority classes. For instance, in datasets with a 90% prevalence of one class, mean imputation likely predicts this dominant class for missing values, exacerbating existing class imbalances and distorting classifier performance. This leads to decreased sensitivity and an increased rate of false negatives. Addressing this issue requires imputation techniques that account for class distribution or employing methods like multiple imputation, which captures data variability more effectively. Integrating class-aware imputation strategies, which consider underlying class proportions during imputation, can also help mitigate bias and improve predictive performance. Imputation is especially problematic in the context of minority classes. Techniques reliant on proximity or similarity, such as k-NN, struggle in sparse datasets where finding a sufficient number of similar cases within the minority class can be challenging Das et al. (2018). This often results in less reliable imputation, as the method may have to use proxies from the majority class to fill gaps, further diluting the characteristics of the minority class. The position and pattern of missing values significantly affect the risk associated with imputation in imbalanced datasets. Missing values more frequently occurring in the minority class can worsen the imbalance and increase bias.
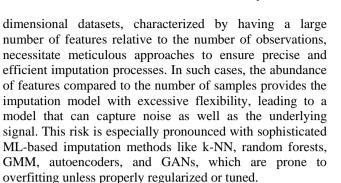
The position and pattern of missing values significantly impact the risk associated with imputation in imbalanced data. Missing values that occur more frequently in the minority class or certain ranges of the target variable can worsen the imbalance and increase bias. For example, if missing values are more common in the minority class, imputation methods may further reduce the representation of this class, making it even harder for the model to learn from it. Similarly, if missing values are concentrated in underrepresented ranges of the target variable, the imputation method may fail to capture the true data, leading to biased estimates.

The risk of bias increases if missing values are not randomly distributed but are more frequent in certain classes or data segments. This pattern can influence the effectiveness of the imputation method, potentially exacerbating existing imbalances. It is crucial to ensure the chosen imputation method is suitable for the specific type of imbalance. For instance, methods that incorporate class weights or sample balancing techniques can help mitigate bias. Similarly, methods that model the entire distribution of the target variable, rather than relying on central tendencies, can provide more accurate imputations in regression tasks.

Practical Recommendations
1. To effectively address the imbalance before imputation, the Synthetic Minority Over-sampling Technique (SMOTE) is employed, which enhances dataset balance by generating synthetic samples Chawla et al. (2002). SMOTE interpolates between existing instances of the minority class to synthesize new samples. For each minority class instance, SMOTE selects one or more of its nearest neighbors and synthesizes new samples along the line segments joining the original instance and its neighbors. By augmenting the minority class with artificial yet plausible examples, SMOTE enriches the dataset's diversity Fernandez´ et al. (2018).

2. Applied prior to imputation, SMOTE ensures the process does not disproportionately favor the majority class, thus preserving the unique characteristics of the minority class and contributing to more accurate and unbiased modeling. By balancing the dataset, SMOTE not only improves the effectiveness of imputation but also enhances the reliability of classification tasks, where underrepresented classes might otherwise be marginalized, Shin et al. (2021).

3. Regularization techniques such as L1 (lasso) and L2 (ridge) play a crucial role in preventing overfitting in models dealing with imbalanced data. L1 regularization promotes sparsity in the model coefficients, beneficial when the dataset contains irrelevant or redundant features, while L2 regularization manages the magnitude of the coefficients, enhancing stability Vidaurre et al. (2013).

4. The integration of ensemble methods like bagging, boosting, and stacking significantly bol-sters the robustness and accuracy of imputations in imbalanced datasets. These methods leverage the strengths of various models to improve overall prediction quality, reducing vari-ance (bagging), iteratively correcting errors (boosting), and optimizing predictions through a meta-learner (stacking) Liu and Zhou (2013).

5. Monitoring the quality of imputation is crucial to identify biases and inaccuracies, especially in datasets where minority classes are underrepresented. Techniques such as histogram and density comparisons, statistical tests like the Kolmogorov-Smirnov, and performance metrics including confusion matrices, precision, recall, and F1-scores, are instrumental in assessing the effectiveness of the imputation Gaudreault et al. (2021). Cross-validation and sensitivity analyses provide additional validation, uncovering any instabilities or inconsistencies that may arise during the imputation process.

### 3.2.2 High Dimensionality

Imputing missing values in high-dimensional datasets presents significant challenges, primarily due to the risks of overfitting and the increased computational demands. High-

dimensional datasets, characterized by having a large number of features relative to the number of observations, necessitate meticulous approaches to ensure precise and efficient imputation processes. In such cases, the abundance of features compared to the number of samples provides the imputation model with excessive flexibility, leading to a model that can capture noise as well as the underlying signal. This risk is especially pronounced with sophisticated ML-based imputation methods like k-NN, random forests, GMM, autoencoders, and GANs, which are prone to overfitting unless properly regularized or tuned.

Overfitting Risks in High-Dimensional Contexts One of the primary concerns in high-dimensional settings is the 'curse of dimensionality', which exacerbates overfitting risks Fan and Li (2006). As the number of features increases, the distance between data points grows, diluting the meaning of "proximity" or "similarity" that is crucial for methods like k-NN. This phenomenon complicates the task of identifying meaningful patterns, reducing the efficacy of distance-dependent imputation methods. The 'curse of dimensionality' further exacerbates overfitting risks. This phenomenon complicates the task of identifying meaningful patterns, reducing the efficacy of imputation methods reliant on distance metrics Fan and Li (2006). Similarly, random forest imputation might develop overly complex trees that are overly specific to the training data, capturing spurious correlations that do not represent true data.

Computational Challenges: Beyond the risk of overfitting, high-dimensional datasets impose significant computational challenges. The complexity of imputation methods escalates with the increase in the number of features, leading to prolonged processing times and escalated memory demands. For instance, the computational cost of k-NN imputation expands as it involves computing distances between each data point, increasing with both the number of features and observations. Similarly, random forest imputation, which constructs multiple decision trees, experiences a heightened computational load as the feature count rises. MICE, which iteratively models each feature with missing values based on other features, becomes increasingly resource-intensive with more features. This iterative nature significantly amplifies the computational burden, making it less feasible for extremely high-dimensional datasets without substantial computational resources. GMM, which models the data as a mixture of Gaussians, faces increased complexity with additional features, necessitating the estimation of more parameters and lengthening convergence times.

Practical Recommendations Handling missing values in high-dimensional datasets requires careful consideration of the overfitting risks and computational constraints. Effective imputation demands selecting appropriate methods that balance complexity with the ability to generalize well to new data, alongside managing the substantial computational resources needed for high-dimensional data processing.

1. Dimensionality Reduction Methods: Dimensionality reduction techniques such as Princi-pal Component Analysis (PCA) and t-Distributed Stochastic Neighbor Embedding (t-SNE) are crucial for managing high-dimensional data, Hasan and Abdulazeez (2021). PCA reduces dimensions by projecting data onto principal components that capture the most variance, simplifying the structure and reducing the number of dimensions Wold et al. (1987). This reduction allows for more accurate imputation in a lower-dimensional space, which can then be transformed back to the original space. t-SNE, on the other hand, focuses on preserving the local structure of data while mapping it to a lower-dimensional space Cieslak et al. (2020). t-SNE is useful for visualizing high-dimensional data and can aid in understanding complex patterns before performing imputation Van der Maaten and Hinton (2008).

2. Sparse Imputation Methods: Sparse imputation methods like Singular Value Decompo-sition (SVD) and Lasso regression effectively handle high-dimensional data by focusing on the most significant components or enforcing sparsity in model coefficients Tay et al. (2021). SVD decomposes the high-dimensional data matrix into simpler matrices, facilitating imputation in a lower-dimensional space while retaining essential data characteristics Kalman (1996). Lasso regression applies L1 regularization, which shrinks less important coefficients to zero, simplifying the model and improving imputation accuracy by reducing the influence of irrelevant or redundant features Vidaurre et al. (2013).

3. Cluster-Based Imputation: Clustering methods such as K-means and GMM segment the data into clusters based on similarity, making imputation more manageable. In K-means clustering, the data is partitioned into k clusters, with each data point assigned to the cluster with the nearest mean Zalik (2008); Wani (2024). Imputation within each cluster is more accurate as similar data points are likely to have similar missing values Zhang et al. (2008). GMMs use a probabilistic approach, assuming the data is generated from a mixture of Gaussian distributions; each component represents a cluster, and imputation is performed based on the parameters (mean and covariance) of these distributions, offering a robust method for dealing with high-dimensional imputation tasks Yu et al. (2015).

In summary, addressing the challenges of high dimensionality in missing data imputation involves a careful balance of methodological rigor and computational strategy. By leveraging the above methods, and by being cognizant of the inherent limitations of various imputation methods,

researchers can achieve more accurate and reliable imputations. This careful approach ensures that the resultant data analyses are both robust and insightful, leading to more informed conclusions and interpretations.

### 3.3 Evaluation and Validation

In real-world datasets, the absence of ground truth for missing values presents a significant challenge in evaluating imputation quality. Without access to the true values, it becomes difficult to directly assess how accurately the imputation method has filled in the gaps. Standard metrics, like Root Mean Squared Error (RMSE), often fall short because they focus on pointwise error between imputed and actual values when true values are available. However, in many cases, the goal of imputation is not just to minimize error but to preserve the underlying data relationships and improve the performance of downstream tasks, such as classification, regression, or clustering Sofi and Wani (2021). To overcome this limitation, alternative evaluation strategies must be adopted. One common approach is to artificially introduce missing values in a dataset where the ground truth is known, allowing for the calculation of metrics like RMSE on the artificially incomplete data. While this provides insight into performance, it may not fully capture the method's ability to maintain complex variable relationships or task-specific accuracy in a real-world scenario.

Another approach involves evaluating the impact of imputation on downstream tasks. For example, the performance of a machine learning model (e.g., classification accuracy, R-squared in regression) trained on the imputed data can serve as an indirect measure of imputation quality. If the model performs well despite missing data, it suggests that the imputation method has effectively preserved the important relationships between variables. Additionally, metrics that assess the structure of the data, such as cluster integrity (e.g., Silhouette Score) or temporal pattern consistency (e.g., Dynamic Time Warping), can provide further validation when ground truth is unavailable. Ultimately, in the absence of ground truth, the effectiveness of imputation should be judged by how well it supports the dataset's intended analytical purpose, emphasizing the need for task-specific evaluation metrics rather than solely relying on pointwise error metrics.

### 3.4 Adaptability to Dynamic Data

In real-time and streaming data environments, data distributions are not static but evolve over time, a phenomenon known as concept drift. Concept drift occurs when the statistical properties of the input features or target variables shift, making static imputation methods inadequate. These traditional imputation techniques often assume that relationships within the data remain constant. However, in dynamic environments such as financial markets, sensor networks, or real-time health monitoring,

these assumptions quickly become outdated, leading to poor imputation and degraded model performance.

The challenge lies in how static imputation methods, which rely on fixed patterns from historical data, struggle to keep pace with changing trends. For instance, imputing missing values in a time series of sensor data based on outdated information may miss critical new patterns or anomalies, compromising decision-making. As a result, the inability to adapt to dynamic changes can significantly impair downstream machine learning models, reducing their predictive accuracy and reliability. To effectively handle concept drift, online learning algorithms provide a solution by continuously updating imputation models as new data arrives. These algorithms adjust their internal parameters incrementally, enabling them to adapt in real time. For example, online versions of k-Nearest Neighbors (k-NN), Support Vector Machines (SVMs), and Gaussian Mixture Models (GMMs) can dynamically adjust to shifts in data patterns, ensuring that imputed values reflect the latest trends rather than outdated correlations.

Sliding windows offer another method of handling dynamic data. By focusing on the most recent observations within a defined window and discarding older data, sliding windows ensure that the imputation is aligned with the latest patterns in the dataset. This approach is particularly effective in scenarios where older data becomes irrelevant due to significant concept drift, as seen in rapidly changing fields like financial forecasting or online user behavior analysis. Furthermore, can be adapted to address concept drift in real-time settings. By combining predictions from multiple models, ensembles can effectively balance historical data with current trends. For example, an ensemble of models trained on different subsets of the data or at different times can produce more robust imputations, accounting for both long-term patterns and short-term shifts in the data. In summary, real-time imputation in dynamic environments requires methods that can quickly adapt to changing data distributions. Online learning algorithms, sliding windows, and dynamic ensemble methods are crucial strategies for maintaining imputation accuracy and ensuring the ongoing reliability of machine learning models. These adaptive techniques enable systems to keep pace with concept drift, preserving the integrity of predictions and decisions in ever-evolving data landscapes.

### IV.    SOLUTIONS

#### 4.1 Domain-Specific Imputation Methods

Imputation methods should be customized to fit the specific characteristics of the data, ensuring both accuracy and the preservation of intrinsic data patterns.  For datasets containing mixed data types numerical and categorical, decision-tree-based methods like Random Forest imputation are particularly effective. These techniques excel at modeling complex, non-linear relationships between

variables by splitting the data into subgroups based on similarity, allowing them to handle different types of data simultaneously. This makes Random Forest imputation ideal for datasets where traditional methods might struggle to account for variable interactions across mixed data types. In time-series data, where preserving temporal continuity is crucial, domain-specific methods such as Kalman filters and autoregressive models are more appropriate. Kalman filters work by estimating missing values through a recursive process that smooths noise and tracks underlying system dynamics over time. This makes them highly effective for applications involving time-varying processes. Autoregressive models, meanwhile, predict missing values based on the linear dependence of current data points on prior observations, capturing trends and seasonal patterns that are common in time-series data. Both methods are designed to preserve the sequential nature of time-series data, ensuring that temporal dependencies and fluctuations are accurately represented. By aligning imputation methods with the specific characteristics of the data—whether through decision-tree-based approaches for mixed data or time-series techniques that account for temporal relationships—domain-specific imputation minimizes bias, reduces error, and ensures the integrity of the imputed dataset.

### 4.1.1 Feature Engineering for Dependencies

Feature engineering plays a vital role in preserving the relationships between interdependent variables during imputation. This is particularly important in datasets where variables are not independent, such as in time-series or spatial data, where imputation methods need to account for inherent dependencies to avoid distorting the underlying patterns. In time-series datasets, one common technique is the creation of lag variables. Lag variables are created by shifting the time-series data by one or more time steps, allowing the imputation method to use past or future values as predictors for missing data. This approach helps to preserve the temporal relationships within the dataset.

In imputation, these lag variables are incorporated into the model, ensuring that the missing values are filled in a way that maintains the temporal structure of the data. In spatial data, similar principles apply, but the focus shifts to capturing spatial dependencies. Spatial imputation methods, such as Kriging or spatial autoregressive models, leverage the spatial relationships between observations. These methods use the values of neighboring observations to predict missing data, based on the assumption that spatially close points tend to exhibit similar characteristics. For example, Kriging models the spatial correlation using a variogram function, which quantifies the spatial dependency between observations based on their geographic distance. These feature engineering techniques ensure that the interdependencies between variables are preserved during imputation, leading to more accurate and contextually

consistent results. In datasets with strong temporal or spatial dependencies, such engineered features can significantly improve the quality of imputation by aligning it with the natural relationships in the data. Moreover, these techniques provide a robust framework for handling missing data in complex, structured datasets, making them invaluable in fields such as economics, geostatistics, environmental science, and any domain that relies on structured time-series or spatial observations.

### 4.2 Reducing Bias in Imputation
### 4.2.1 Class-Aware Imputation

Class-aware imputation techniques, such as stratified imputation, aim to preserve class distributions while handling missing data, making them particularly valuable for imbalanced datasets. These methods ensure that missing values are imputed separately within each class, preventing the over-representation of majority classes in the imputed data. By stratifying imputation by class, the inherent distribution of both minority and majority classes is maintained, reducing biases that could distort subsequent analysis or model performance. This is especially important in classification tasks where accurate representation of all classes, particularly the minority classes, is crucial for model fairness and effectiveness. Class-aware imputation methods help mitigate the risk of the majority class dominating the imputation process, which could otherwise lead to skewed predictions and unfair performance evaluations, particularly in highly imbalanced datasets.

### 4.3 Addressing Computational Constraints
### 4.3.1 Approximate and Efficient Methods

To address computational constraints, approximate imputation methods provide a balance between accuracy and computational efficiency. These techniques are designed to reduce resource consumption, particularly in environments with limited processing power or memory. Two prominent approaches are mini-batch processing and low-rank matrix factorization, such as Singular Value Decomposition (SVD). Mini-batch Processing: Mini-batch processing tackles the imputation problem by dividing large datasets into smaller, manageable subsets (mini-batches). Instead of processing the entire dataset in a single pass, the imputation algorithm iterates over mini-batches, updating imputed values incrementally. The size of the mini-batches is a critical factor; smaller batches reduce memory usage but may lead to noisier estimates, while larger batches improve accuracy but require more memory.

This method reduces the computational complexity from $O(n)$ to $O(m)$ per iteration, where n is the size of the dataset and m is the size of the mini-batch (with $m \ll n$). Low-Rank Matrix Factorization (SVD): Low-rank matrix factorization techniques, such as SVD, aim to approximate the original data matrix by decomposing it into lower-dimensional

components. This technique assumes that the missing values can be imputed by capturing the underlying low-dimensional structure of the data. The SVD decomposes the data matrix X into three matrices: The computational complexity of SVD is $O(n2k)$, which is significantly less than full matrix imputation methods when k is small compared to the matrix dimensions. While SVD offers computational savings, it is particularly effective when the data exhibits low-rank properties, which allows for reasonable approximations with minimal loss in imputation accuracy.

Trade-offs: Both mini-batch processing and low-rank matrix factorization present trade-offs between accuracy and computational efficiency. Mini-batch processing accelerates computation but may result in slightly noisier imputations due to reduced data exposure in each iteration. Low-rank matrix factorization simplifies data by assuming an underlying low-dimensional structure, which works well in many real- world datasets but may not capture more complex patterns or interactions in high-rank datasets. By utilizing these approximate methods, it is possible to impute large-scale datasets more efficiently, making them well-suited for environments with limited computational resources, such as edge devices, embedded systems, or cloud infrastructure with cost constraints.

### 4.4 Improved Evaluation and Validation Techniques
#### 4.4.1 Task-Specific Evaluation Metrics

Evaluating imputation methods requires more than simply calculating general metrics like Root Mean Squared Error (RMSE), as these often fail to reflect how well imputation preserves critical relationships in the data. Instead, task-specific metrics provide a clearer picture of imputation quality, helping researchers better assess how imputations impact downstream tasks. In classification tasks, metrics such as the F1-score and Area Under the Receiver Operating Char- acteristic Curve (AUC) are more appropriate than simple accuracy, especially for imbalanced datasets. For example, in a medical dataset with a rare disease class, the F1-score helps ensure that the imputation doesn't compromise the model's ability to correctly identify this minority class, while AUC provides a comprehensive measure of how well the imputed data supports distinguishing between classes across varying thresholds. For multi-class classification, macro-averaged and micro-averaged F1-scores offer complementary perspectives. In a plant species classification dataset, macro-averaged scores ensure that rare species are considered equally, while micro-averaged scores reflect overall classification performance, weighted by class size. In regression tasks, metrics like Mean Absolute Error (MAE) and Mean Absolute Percentage Error (MAPE) are often better suited than RMSE. For instance, in a financial forecasting model, MAE directly measures the average difference between imputed and actual values, giving a clear sense of how imputation

affects predictions. MAPE is particularly valuable when relative accuracy matters, such as when imputing missing sales figures for a product line. In clustering tasks, imputation can affect the coherence of clusters. For example, in customer segmentation, the Silhouette Score measures how well imputed data maintains the distinction between clusters, while the Davies-Bouldin Index quantifies intra-cluster similarity versus inter-cluster separation, ensuring the imputation process doesn't blur important groupings. For time-series tasks, preserving temporal patterns during imputation is critical. Dynamic Time Warping (DTW) distance measures how well imputation captures temporal dynamics, which is particularly important in applications like energy consumption forecasting, where shifts in time must be aligned.

Symmetric Mean Absolute Percentage Error (SMAPE) further refines this by mitigating the effects of extreme values, offering a more robust assessment for volatile time-series data. In anomaly detection, precision-focused metrics are vital. For instance, in fraud detection, Precision at K measures the ability to correctly identify the top suspected fraud cases, while the Precision-Recall AUC ensures that imputation doesn't inflate false positives or obscure rare but critical events. Using task-specific metrics allows for a targeted evaluation of how imputation methods impact downstream tasks, ensuring that they support the specific goals of the analysis rather than just minimizing error.

#### 4.4.2 Sensitivity and Robustness Testings

To ensure that imputation methods are reliable and don't introduce bias or inconsistencies into downstream models, sensitivity and robustness testing are essential. This involves testing how different imputation strategies affect model outcomes under various scenarios, providing insights into their stability and effectiveness. Monte Carlo simulations are a common method for sensitivity testing. In practice, this involves generating multiple imputed datasets under varying parameters or assumptions. For instance, in a healthcare dataset, multiple versions of the dataset can be created with different imputation techniques (e.g., mean imputation, k-nearest neighbors) or with different settings for the same method. These datasets are then used to train models, and their performance is compared. By analyzing the variance in performance metrics, such as classification accuracy or mean error, researchers can assess how sensitive the downstream models are to changes in the imputation process. Low variance indicates that the imputation method is robust, while high variance suggests instability. A more concrete step-by-step approach to Monte Carlo testing includes: 1. Generate N imputed datasets by applying different imputation methods or varying hyperparameters. 2. Train models on each imputed dataset using a consistent modeling approach. 3. Evaluate performance across all imputed datasets using relevant metrics (e.g., F1-score, MAE, AUC). 4. Analyze variance in the results to assess

imputation robustness, noting how sensitive model performance is to changes in the imputation method. Additionally, perturbation analysis introduces small, controlled variations in the imputation process to observe their effects on outcomes. For example, by slightly altering input data (e.g., shifting a few key variables) before and after imputation, researchers can evaluate whether the imputation method handles minor fluctuations robustly. This approach helps identify methods that may be overly sensitive to noise or small variations in the data, further refining the choice of an appropriate imputation technique. Together, these sensitivity and robustness tests provide a robust empirical foundation for selecting imputation methods. They ensure that the chosen imputation technique not only fills in missing data but also supports reliable and consistent results across different scenarios, making the analysis more generalizable and trustworthy.

## V.    CONCLUSION

In this paper, we explored the various challenges associated with missing data imputation and proposed solutions to address these issues. By applying domain-specific imputation methods, reducing bias, scaling computational techniques, and using appropriate evaluation metrics, we can improve the quality and applicability of imputed data in real-world scenarios.

## VI.    REFERENCE

[1]    Ahn, H., Sun, K., Kim, K. P., et al. "Comparison of missing data imputation methods in time series forecasting." Computers, Materials & Continua, vol. 70, no. 1, 2022, pp. 767–779. DOI: 10.32604/cmc.2022.019369.

[2]    Alam, S., Ayub, M. S., Arora, S., and Khan, M. A. "An investigation of the imputation techniques for missing values in ordinal data enhancing clustering and classification analysis validity." Decision Analytics Journal, vol. 9, 2023, pp. 100341. DOI: 10.1016/j.dajour.2023.100341.

[3]    Aleryani, A., Wang, W., and De La Iglesia, B. "Multiple imputation ensembles (MIE) for dealing with missing data." SN Computer Science, vol. 1, no. 3, 2020, pp. 134. DOI: 10.1007/s42979-020-00136-0.

[4]    Ali, H., Salleh, M. M., Saedudin, R., Hussain, K., and Mushtaq, M. F. "Imbalance class problems in data mining: A review." Indonesian Journal of Electrical Engineering and Computer Science, vol. 14, no. 3, 2019, pp.1560–1571.DOI: 10.11591/ijeecs.v14.i3.pp1560-1571.

[5]    Asheri, H., Hosseini, R., and Araabi, B. N. "A new EM algorithm for flexibly tied GMMs with a large number of components." Pattern Recognition, vol. 114, 2021, pp. 107836. DOI: 10.1016/j.patcog.2021.107836.

[6]    Azur, M. J., Stuart, E. A., Frangakis, C., and Leaf, P. J. "Multiple imputation by chained equations: what is it and how does it work?" International Journal of Methods in Psychiatric Research, vol. 20, no. 1, 2011, pp. 40–49. DOI: 10.1002/mpr.329.

[7]    Beesley, L. J., and Taylor, J. M. "A stacked approach for chained equations multiple imputation incorporating the substantive model." Biometrics, vol. 77, no. 4, 2021, pp. 1342–1354.

[8]    Beretta, L., and Santaniello, A. "Nearest neighbor imputation algorithms: a critical evaluation." BMC Medical Informatics and Decision Making, vol. 16, 2016, pp. 74. DOI: 10.1186/s12911-016-0318-z.

[9]    Buczak, P., Chen, J.-J., and Pauly, M. "Analyzing the effect of imputation on classification performance under MCAR and MAR missing mechanisms." Entropy, vol. 25, no. 3, 2023, pp. 521. DOI: 10.3390/e25030521.

[10]    Chawla, N. V., Bowyer, K. W., Hall, L. O., and Kegelmeyer, W. P. "SMOTE: synthetic minority over-sampling technique." Journal of Artificial Intelligence Research, vol. 16, 2002, pp. 321–357. DOI: 10.1613/jair.953.

[11]    Cho, B., Dayrit, T., Gao, Y., Wang, Z., Hong, T., Sim, A., and Wu, K. "Effective missing value imputation methods for building monitoring data." IEEE International Conference on Big Data, 2020, pp. 2866–2875.DOI: 10.1109/BigData50022.2020.9378230

[12]    .

[13]    Cieslak, M. C., Castelfranco, A. M., Roncalli, V., Lenz, P. H., and Hartline, D. K. "t-distributed stochastic neighbor embedding (t-SNE): a tool for eco-physiological transcriptomic analysis." Marine Genomics, vol. 51, 2020, pp. 100723, DOI:10.1016/j.margen.2019.100723

[14]    Costa, A. F., Santos, M. S., Soares, J. P., and Abreu, P. H. "Missing data imputation via denoising autoencoders: the untold story." Advances in Intelligent Data Analysis XVII: 17th International Symposium, IDA 2018, pp. 87–98, Springer, 2018. DOI: 10.1007/978-3-030-01768-2_8.

[15]    Das, S., Datta, S., and Chaudhuri, B. B. "Handling data irregularities in classification: foundations, trends, and future challenges." Pattern Recognition, vol. 81, 2018, pp. 674–693, DOI: 10.1016/j.patcog.2018.03.008

[16]    Enders, C. K. Applied Missing Data Analysis. Guilford Publications, 2022.

[17]    Esling, P., and Agon, C. "Time-series data mining." ACM Computing Surveys (CSUR), vol. 45, no. 1, 2012, pp. 1–34. DOI: 10.1145/2379776.2379788.

[18]    Fan, J., and Li, R. "Statistical challenges with high dimensionality: Feature selection in knowledge

discovery." arXiv preprint, math/0602133, 2006. DOI: 10.48550/arXiv.math/0602133.

[19] Fernández, A., García, S., Herrera, F., and Chawla, N. V. "SMOTE for learning from imbalanced data: progress and challenges, marking the 15-year anniversary." Journal of Artificial Intelligence Research, vol. 61, 2018, pp. 863–905. DOI: 10.1613/jair.1.11192.

[20] Gaudreault, J.-G., Branco, P., and Gama, J. "An analysis of performance metrics for imbalanced classification." Discovery Science: 24th International Conference, Springer, 2021, pp. 67–77.

[21] Hasan, B. M. S., and Abdulazeez, A. M. "A review of principal component analysis algorithm for dimensionality reduction." Journal of Soft Computing and Data Mining, vol. 2, no. 1, 2021, pp. 20–30. DOI: 10.30880/jscdm.2021.02.01.003.

[22] Jadhav, A., Pramod, D., and Ramanathan, K. "Comparison of performance of data imputation methods for numeric dataset." Applied Artificial Intelligence, vol. 33, no. 10, 2019, pp. 913–933. DOI: 10.1080/08839514.2019.1637138.

[23] Kalman, D. "A singularly valuable decomposition: the SVD of a matrix." The College Mathematics Journal, vol. 27, no. 1, 1996, pp. 2–23. DOI: 10.1080/07468342.1996.11973744.

[24] Kenward, M. G., and Molenberghs, G. "Last observation carried forward: a crystal ball?" Journal of Biopharmaceutical Statistics, vol. 19, no. 5, 2009, pp. 872–888. DOI: 10.1080/10543400903105406.

[25] Kokla, M., Virtanen, J., Kolehmainen, M., Paananen, J., and Hanhineva, K. "Random forest-based imputation outperforms other methods for imputing LC-MS metabolomics data: a comparative study." BMC Bioinformatics, vol. 20, 2019, pp. 1–11. DOI: 10.1186/s12859-019-3110-0.

[26] Liu, X.-Y., and Zhou, Z.-H. "Ensemble methods for class imbalance learning." Imbalanced Learning: Foundations, Algorithms, and Applications, Wiley, 2013, pp. 61–82. DOI: 10.1002/9781118646106.ch4.

[27] Ma, Q., Lee, W.-C., Fu, T.-Y., Gu, Y., and Yu, G. "MiDIA: Exploring denoising autoencoders for missing data imputation." Data Mining and Knowledge Discovery, vol. 34, 2020, pp. 1859–1897, DOI: 10.1007/s10618-020-00706-8

[28] Molenberghs, G., Fitzmaurice, G., Kenward, M. G., Tsiatis, A., and Verbeke, G. Handbook of Missing Data Methodology. CRC Press, 2014. DOI: 10.1201/b17622

[29] Muja, M., and Lowe, D. G. "Scalable nearest neighbor algorithms for high dimensional data." IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 36, no. 11, 2014, pp. 2227–2240. DOI: 10.1109/TPAMI.2014.2321376.

[30] Newman, D. A. "Missing data: Five practical guidelines." Organizational Research Methods, vol. 17, no. 4, 2014, pp. 372–411. DOI: 10.1177/1094428114548590.

[31] Pansara, R. "Cultivating data quality to strategies, challenges, and impact on decision-making." International Journal of Management Education for Sustainable Development, vol. 6, no. 6, 2023, pp. 24–33.

[32] Pereira, R. C., Santos, M. S., Rodrigues, P. P., and Abreu, P. H. "Reviewing autoencoders for missing data imputation: Technical trends, applications and outcomes." Journal of Artificial Intelligence Research, vol. 69, 2020, pp. 1255–1285. DOI: 10.1613/jair.1.12312.

[33] Pinaya, W. H. L., Vieira, S., Garcia-Dias, R., and Mechelli, A. "Autoencoders." Machine Learning, Elsevier, 2020, pp. 193–208.

[34] Plaia, A., Buscemi, S., Furnkranz, J., and Mencía, E. L. "Comparing boosting and bagging for decision trees of rankings." Journal of Classification, vol. 39, no. 1, 2022, pp. 78–99. DOI: 10.1007/s00357-021-09400-0.

[35] Rahman, M. G., and Islam, M. Z. "Missing value imputation using decision trees and decision forests by splitting and merging records: Two novel techniques." Knowledge-Based Systems, vol. 53, 2013, pp. 51–65. DOI: 10.1016/j.knosys.2013.08.023.

[36] Rubin, D. B. "Inference and missing data." Biometrika, vol. 63, no. 3, 1976, pp. 581–592. DOI: 10.1093/biomet/63.3.581.

[37] Santos, M. S., Abreu, P. H., Wilk, S., and Santos, J. "How distance metrics influence missing data imputation with k-nearest neighbours." Pattern Recognition Letters, vol. 136, 2020, pp. 111–119. DOI: 10.1016/j.patrec.2020.05.032.

[38] Shah, A. D., Bartlett, J. W., Carpenter, J., Nicholas, O., and Hemingway, H. "Comparison of random forest and parametric imputation models for imputing missing data using MICE: a CALIBER study." American Journal of Epidemiology, vol. 179, no. 6, 2014, pp. 764–774. DOI: 10.1093/aje/kwt312.

[39] Shin, K., Han, J., and Kang, S. "Mi-MOTE: Multiple imputation-based minority oversampling technique for imbalanced and incomplete data classification." Information Sciences, vol. 575, 2021, pp. 80–89, DOI: 10.1016/j.ins.2021.06.043.

[40] Sofi, S. A., and Wani, A. A. "Predicting material stability using machine learning." Applications of Advanced Computing in Systems: Proceedings of International Conference on Advances in Systems, Control and Computing, Springer, 2021, pp. 203–209. DOI: 10.1007/978-981-33-4862-2_21.

[41] Steele, R. J., and Raftery, A. E. "Performance of Bayesian model selection criteria for Gaussian mixture models." Frontiers of Statistical Decision Making and Bayesian Analysis, vol. 2, 2010, pp. 113–130. DOI: 10.1007/978-1-4419-6944-6_10.

[42] Tang, F., and Ishwaran, H. "Random forest missing data algorithms." Statistical Analysis and Data Mining: The ASA Data Science Journal, vol. 10, no. 6, 2017, pp. 363–377. DOI: 10.1002/sam.11348.

[43] Tay, J. K., Friedman, J., and Tibshirani, R. "Principal component-guided sparse regression." Canadian Journal of Statistics, vol. 49, no. 4, 2021, pp. 1222–1257.

[44] Triguero, I., García-Gil, D., Maillo, J., Luengo, J., García, S., and Herrera, F. "Transforming big data into smart data: An insight on the use of the k-nearest neighbors algorithm to obtain quality data." Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, vol. 9, no. 2, 2019, pp. e1289. DOI: 10.1002/widm.1289.

[45] Van Buuren, S., and Groothuis-Oudshoorn, K. "mice: Multivariate imputation by chained equations in R." Journal of Statistical Software, vol. 45, no. 3, 2011, pp. 1–67. DOI: 10.18637/jss.v045.i03.

[46] Van der Maaten, L., and Hinton, G. "Visualizing data using t-SNE." Journal of Machine Learning Research, vol. 9, no. 11, 2008, pp. 2579–2605..

[47] Van Ginkel, J. R., Linting, M., Rippe, R. C. A., and Van Der Voort, A. "Rebutting existing misconceptions about multiple imputation as a method for handling missing data." Journal of Personality Assessment, vol. 102, no. 3, 2020, pp. 297–308. DOI: 10.1080/00223891.2018.1530680.

[48] Vidaurre, D., Bielza, C., and Larrañaga, P. "A survey of L1 regression." International Statistical Review, vol. 81, no. 3, 2013, pp. 361–387, DOI:10.1111/insr.12023

[49] Wani, A. A. "Comprehensive analysis of clustering algorithms: exploring limitations and innovative solutions." PeerJ Computer Science, vol. 10, 2024, pp. e2286. DOI: 10.7717/peerj-cs.2286.

[50] Wold, S., Esbensen, K., and Geladi, P. "Principal component analysis." Chemometrics and Intelligent Laboratory Systems, vol. 2, nos. 1–3, 1987, pp. 37–52. DOI: 10.1016/0169-7439(87)80084-9.

[51] Yu, D., Deng, L., Yu, D., and Deng, L. "Gaussian mixture models." Automatic Speech Recognition: A Deep Learning Approach, Elsevier, 2015, pp. 13–21. DOI: 10.1016/B978-0-12-802398-4.00002-6.

[52] Zalik, K. R. "An efficient k-means clustering algorithm." Pattern Recognition Letters, vol. 29, no. 9, 2008, pp. 1385–1391, DOI:10.1016/j.patrec.2008.02.014

[53] Zhang, S. "Nearest neighbor selection for iteratively KNN imputation." Journal of Systems and Software, vol. 85, no. 11, 2012, pp. 2541–2552, DOI:10.1016/j.jss.2012.05.073.

[54] Zhang, S., Zhang, J., Zhu, X., Qin, Y., and Zhang, C. "Missing value imputation based on data clustering." Transactions on Computational Science I, Springer, 2008, pp. 128–138. DOI: 10.1007/978-3-540-78491-8_7.