
Rethinking the Diffusion Models for Numerical Tabular Data Imputation from the Perspective of Wasserstein Gradient Flow

Zhichao Chen/Ziciu Can¹ Haoxuan Li² Fangyikang Wang¹ Odin Zhang¹ Hu Xu¹
Xiaoyu Jiang¹ Zhihuan Song¹ Eric H. Wang^{1*}
¹Zhejiang University ²Peking University

Abstract

Diffusion models (DMs) have gained attention in Missing Data Imputation (MDI), but there remain two long-neglected issues to be addressed: (1). Inaccurate Imputation, which arises from inherently sample-diversification-pursuing generative process of DMs. (2). Difficult Training, which stems from intricate design required for the mask matrix in model training stage. To address these concerns within the realm of numerical tabular datasets², we introduce a novel principled approach termed Kernelized Negative Entropy-regularized Wasserstein gradient flow Imputation (KnewImp). Specifically, based on Wasserstein gradient flow (WGF) framework, we first prove that issue (1) stems from the cost functionals implicitly maximized in DM-based MDI are equivalent to the MDI’s objective plus diversification-promoting non-negative terms. Based on this, we then design a novel cost functional with diversification-discouraging negative entropy and derive our KnewImp approach within WGF framework and reproducing kernel Hilbert space. After that, we prove that the imputation procedure of KnewImp can be derived from another cost functional related to the joint distribution, eliminating the need for the mask matrix and hence naturally addressing issue (2). Extensive experiments demonstrate that our proposed KnewImp approach significantly outperforms existing state-of-the-art methods.

1 Introduction

The imputation of missing values from observational data is crucial for constructing machine learning models with broad applications across various fields, including e-commerce [21, 20, 42], health-care [38], and process industry [24]. Recently, diffusion models (DMs) have emerged as a powerful tool for missing data imputation (MDI), celebrated for their excellent capability to model data distributions and generate high-quality synthetic data [29, 36, 49]. These models excel by approximating the (Stein) score function of the conditional distribution between missing and observed data, thereby reformulating the imputation problem as a generative task grounded in the learned score function.

Although DMs have shown considerable success in MDI tasks, they face significant challenges that result from model inference and training: (1). Inaccurate Imputation: While DM-based approaches treat MDI as a conditional generative task by sampling from the learned score function, it is important to note that the primary evaluation metric for MDI focuses on accuracy [13, 28], rather than the sample diversification typically emphasized in generative tasks. Consequently, the inference objectives implicitly pursued for DMs may not align well with the specific needs of the MDI task. (2). Difficult Training: The training of diffusion models is complicated by the unknown nature of the ground-truth values of missing data. Previous methods [38, 7] have sought to address this by masking parts of the

*Corresponding author.

²Datasets organized in a table format where each entry is a numerical value.

available observational data and then imputing these masked entries as a means to construct model training. In this procedure, the design of mask matrices is essential, as pointed out by reference [38], and hence results in a training difficulty due to complex mask matrix selection mechanism.

The cornerstone for mitigating issue (1) lies in identifying the cost functional that is ‘secretly maximized’ during the DM-based MDI procedure, delineating its relationship to the ‘vanilla’ MDI’s cost functional, refining it where it proves deficient, and redeveloping the corresponding imputation procedure for this enhanced functional. Building on this, the resolution to issue (2) involves bypassing the use of the conditional distribution in the imputation phase. In other words, maintaining the existing imputation procedures while transforming the cost functional into an equivalent form that merely contains the joint distribution. To tackle these challenges in the realm of numerical tabular data, we introduce our approach named Kernelized Negative Entropy-regularized Wasserstein Gradient Flow Imputation (KnewImp). Specifically, to address issue (1), we unify the DMs’ generative processes into the Wasserstein Gradient Flow (WGF) framework, recover their cost functionals, and validate their connections between MDI’s objective. Based on this analysis, we then introduce a novel negative entropy-regularized (NER) cost functional, and establish a new easy-to-implement and closed-form MDI procedure similar to DMs’ generative processes within the WGF framework and reproducing kernel Hilbert space (RKHS). After that, to circumvent the use of the mask matrix and address issue (2), we further develop a novel cost functional concerned with joint distribution, proving it serves as a lower bound to the NER functional, with a constant gap that preserves the imputation procedure within the WGF framework, and consequently streamlining the model training procedure.

In summary, the contributions of this manuscript can be summarized as follows:

- We elucidate the inaccurate imputation issue for DM-based MDI approaches by revealing the relationship between their cost functional and identical MDI’s objective. Based on this, we propose KnewImp by designing a novel NER functional and obtaining corresponding imputation procedure.
- We demonstrate that the conditional distribution-related NER functional, can be seamlessly transformed into another joint distribution-related NER functional, maintaining a constant gap. Consequently, we bypass the design of a mask matrix during the model training stage of KnewImp without explicitly altering its imputation procedure.
- We empirically validate the superiority of our proposed KnewImp method over state-of-the-art models through rigorous testing on various numerical tabular datasets.

2 Preliminaries

2.1 Missing Data Imputation

Building upon previous works in MDI [28, 54], our objective is articulated as follows: Consider a numerical tabular dataset represented by a matrix $\mathbf{X} \in \mathbb{R}^{N \times D}$, comprising N samples each of dimension D , where certain entries are missing. Accompanying \mathbf{X} is a binary mask matrix $\mathbf{M} \in \{0, 1\}^{N \times D}$, where entry $M[i, d]$ is set as 0 if $\mathbf{X}[i, d]$ is missing, thereby assigned as $\mathbf{X}[i, d] = \text{NaN}$ (denoting ‘not a number’), and 1 otherwise. Hence, the matrix \mathbf{X} is expressed as $\mathbf{X} = \mathbf{X}^{(\text{obs})} \odot \mathbf{M} + \text{NaN} \odot (\mathbb{1}_{N \times D} - \mathbf{M})$, where $\mathbf{X}^{(\text{obs})}$ represents the matrix with observed entries, \odot denotes the Hadamard product, and $\mathbb{1}_{N \times D}$ is a matrix of ones sized $N \times D$. The task at hand is to impute the missing entries in \mathbf{X} , yielding an estimation $\hat{\mathbf{X}}$ using the imputed matrix $\mathbf{X}^{(\text{imp})}$, formulated as $\hat{\mathbf{X}} = \mathbf{X}^{(\text{obs})} \odot \mathbf{M} + \mathbf{X}^{(\text{imp})} \odot (\mathbb{1}_{N \times D} - \mathbf{M})$. Notably, according to reference [31], missing data can be classified into three categories: Missing Completely at Random (MCAR), Missing at Random (MAR), and Missing Not at Random (MNAR)³, and we mainly restrict our discussion scope on numerical tabular data with MAR and MCAR settings.

2.2 Diffusion Models and its application for Missing Data Imputation

According to reference [36], DMs begin by corrupting data towards a tractable noise distribution, typically a standard Gaussian, and then reverse this process to generate samples. Specifically, the forward corruption or diffusion process can be described as a discretization of the stochastic differential equation (SDE) along time τ : $d\mathbf{X}_\tau = f(\mathbf{X}_\tau, \tau)d\tau + g_\tau dW_\tau$, where $f(\mathbf{X}_\tau, \tau)$ is drift

³Detailed information about these missing mechanisms is given in Appendix D.1.

term, g_τ is volatility term, and dW_τ is standard Wiener process. The solution to this SDE forms a continuous trajectory of random variables $\mathbf{X}_\tau|_{\tau=0}^T$. The density function q_τ of these variables is governed by the Fokker-Planck-Kolmogorov (FPK) equation $\frac{\partial q_\tau}{\partial \tau} = -\nabla \cdot (q_\tau f(\mathbf{X}_\tau, \tau)) + \frac{1}{2} g_\tau^2 \nabla \cdot \nabla q_\tau$, as per Theorem 5.4 in reference [34]. According to reference [3], the reverse process for sample generation is described by: $d\mathbf{X}_\tau = [f(\mathbf{X}_\tau, \tau) - g_\tau^2 \nabla \log p(\mathbf{X}_\tau)] d\tau + g_\tau dW_\tau$, where $\nabla \log p(\mathbf{X}_\tau)$ represents the score function and learned via neural networks during DM training phase.

Based on this, DMs approach MDI as a generative problem, and the score function $\nabla \log p(\mathbf{X})$ in the reverse process is replaced with conditional distribution $\nabla_{\mathbf{X}^{(\text{miss})}} \log p(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})})$ [38]. Therefore, the challenge in DM-based MDI is to obtain an estimation $\nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})})$ that effectively approximates $\nabla_{\mathbf{X}^{(\text{miss})}} \log p(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})})$. However, constructing model training remains challenging due to ground truth $\mathbf{X}^{(\text{miss})}$ is unknown. To alleviate this issue, previous DM-based MDI approaches necessitate the design of a mask matrix to obscure portions of the observational data; despite practical efficacy, the selection of mask mechanism, which determines the effectiveness of $\nabla \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})})$, remain challenges and hence may result in training difficulty.

2.3 Wasserstein Gradient Flow

The Wasserstein space $\mathcal{P}_2(\mathbb{R}^D) := \{p \in \mathcal{P}_2(\mathbb{R}^D) : \int \|x\|^2 dp(x) < \infty\}$ is a metric space where distances between probability distributions are quantified using the 2-Wasserstein distance, defined as $(\inf_{\gamma \in \Gamma(p, q)} \int \|x - y\|^2 d\gamma(x, y))^{\frac{1}{2}}$. In this space, gradient flows resemble the steepest descent curves similar to those in classical Euclidean spaces. Specifically, for a cost functional $\mathcal{F}_{\text{cost}} : \mathcal{P}_2(\mathbb{R}^D) \rightarrow \mathbb{R}$, a gradient flow in Wasserstein space is an absolute continuous trajectory $(q_\tau)_{\tau > 0}$ that seeks to minimize $\mathcal{F}_{\text{cost}}$ as efficiently as possible, as described in [33]. This dynamic process is governed by the celebrated continuity equation $\frac{\partial q_\tau}{\partial \tau} = -\nabla \cdot (v_\tau q_\tau)$, where $v_\tau : \mathbb{R}^D \rightarrow \mathbb{R}^D$ is a time-dependent velocity field [2]. Additionally, the evolution of sample \mathbf{X} over time τ in $\mathcal{P}_2(\mathbb{R}^D)$ can be delineated by the ordinary differential equation (ODE) expressed as $\frac{d\mathbf{X}}{d\tau} = v_\tau(\mathbf{X})$.

3 Proposed Approach

This section proposes our Kernelized Negative Entropy-regularized Wasserstein gradient flow Imputation (KnewImp) approach. We first define the MDI’s objective oriented towards maximization and unify the DM-based MDI approaches within WGF framework. In this procedure, we prove that addressing MDI through the generative processes of DMs maximizes a diversification-promoting upper bound of MDI’s objective. Building on this foundational analysis, we introduce a novel diversification-discouraging negative entropy-regularized (NER) cost functional that acts as a lower bound for the MDI’s objective, ensuring precise imputation through maximizing the MDI’s lower bound. Based on this, we then develop the imputation procedure that features a closed-form, easily implementable expression within the WGF framework and RKHS. Further, we establish that our NER functional, associated with the conditional distribution $p(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})})$, is equivalent to another functional concerned with the joint distribution $p(\mathbf{X}^{(\text{miss})}, \mathbf{X}^{(\text{obs})})$, adjusted by a constant. This equivalence maintains the same velocity field but effectively eliminates the need for a mask matrix in the training stage. Finally, we conclude this section by outlining the procedure of KnewImp approach that encapsulates these innovations.

3.1 Unifying DM-based MDI within WGF framework

Drawing from previous works [27, 38], the objective function for the MDI task can be defined as:

$$\mathbf{X}^{(\text{imp})} = \arg \max_{\mathbf{X}^{(\text{miss})}} \log p(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}). \quad (1)$$

From the perspective of generative models, $\mathbf{X}^{(\text{imp})}$ is considered as samples drawn from a certain distribution $r(\mathbf{X}^{(\text{miss})})$ (we name it ‘proposal distribution’), and results in the following reformulation:

$$\arg \max_{\mathbf{X}^{(\text{miss})}} \log p(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}), \mathbf{X}^{(\text{miss})} \sim r(\mathbf{X}^{(\text{miss})}) \Rightarrow \arg \max_{\mathbf{X}^{(\text{miss})}} \mathbb{E}_{r(\mathbf{X}^{(\text{miss})})} [\log p(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})})]. \quad (2)$$

Consequently, DM-based MDI approaches address this optimization problem by generating samples from the estimated conditional score function: $\nabla_{\mathbf{X}^{(miss)}} \log \hat{p}(\mathbf{X}^{(miss)} | \mathbf{X}^{(obs)})$.

Note that, according to Section 2.2, the generative process of DMs satisfies the FPK equation, a specific instance of the continuity equation that underpins the WGF according to Section 2.3. Meanwhile, WGF framework is central to functional optimization, indicating that the divergence of the objective between MDI and the DM-based MDI approaches can be effectively analyzed within the WGF framework. In support of this, we first give the following proposition, which elucidates the relationship between the objective of DM-based MDI approaches and MDI⁴:

Proposition 3.1. *Within WGF framework, DM-based MDI approaches can be viewed as finding the imputed values $\mathbf{X}^{(imp)}$ that maximize the following objective:*

$$\arg \max_{\mathbf{X}^{(miss)}} \mathbb{E}_{r(\mathbf{X}^{(miss)})} [\log \hat{p}(\mathbf{X}^{(miss)} | \mathbf{X}^{(obs)})] + \psi(\mathbf{X}^{(miss)}) + const, \quad (3)$$

where *const* is the abbreviation of constant, and $\psi(\mathbf{X}^{(miss)})$ is a scalar function determined by the type of SDE underlying the DMs.

- **VP-SDE:** $\psi(\mathbf{X}^{(miss)}) = \mathbb{E}_{r(\mathbf{X}^{(miss)})} \left\{ \frac{1}{4} [\mathbf{X}^{(miss)}]^\top [\mathbf{X}^{(miss)}] - \frac{1}{2} \log r(\mathbf{X}^{(miss)}) \right\}$
- **VE-SDE:** $\psi(\mathbf{X}^{(miss)}) = \mathbb{E}_{r(\mathbf{X}^{(miss)})} \left\{ -\frac{1}{2} \log r(\mathbf{X}^{(miss)}) \right\}$
- **sub-VP-SDE:** $\psi(\mathbf{X}^{(miss)}) = \mathbb{E}_{r(\mathbf{X}^{(miss)})} \left\{ \frac{1}{4\gamma_\tau} [\mathbf{X}^{(miss)}]^\top [\mathbf{X}^{(miss)}] - \frac{1}{2} \log r(\mathbf{X}^{(miss)}) \right\}$, where γ_τ is determined by noise scale β_τ : $\gamma_\tau := (1 - \exp(-2 \int_0^\tau \beta_s ds)) > 0, 0 < \beta_1 < \beta_2 < \dots < \beta_T < 1$.

It is important to note that in DMs, the condition $\psi(\mathbf{X}^{(miss)}) \geq 0$ consistently holds. This assertion is supported by the fact that the inner product $[\mathbf{X}^{(miss)}]^\top [\mathbf{X}^{(miss)}] \geq 0$, and the entropy function defined as $\mathbb{H}[r(\mathbf{X}^{(miss)})] := - \int r(\mathbf{X}^{(miss)}) \log r(\mathbf{X}^{(miss)}) d\mathbf{X}^{(miss)}$ is also non-negative.

Based on this proposition, it becomes evident that ‘inaccurate imputation’ issue may arise from the misalignment in the optimization objectives: By comparing Eqs. (2) and (3), we observe that DM-based MDI methods are optimizing an upper bound of the MDI. Additionally, it is important to note that the gaps involving inner products and entropy, which promote maximization, inherently encourage sample diversification [40]. This diversification is fundamentally at odds with the precision required in MDI tasks. To achieve a more accurate imputation, it is crucial to reformulate $\psi(\mathbf{X}^{(miss)})$ to ensure that $\psi(\mathbf{X}^{(miss)}) \leq 0$, thereby aligning the objective of the imputation procedure with the accuracy-oriented goal of MDI.

3.2 Negative Entropy Regularized & Closed-form Velocity Field Expression

Based on previous subsection, we adopt negative entropy as $\psi(\mathbf{X}^{(miss)})$ intuitively:

$$\psi(\mathbf{X}^{(miss)}) = \lambda \int r(\mathbf{X}^{(miss)}) \log r(\mathbf{X}^{(miss)}) d\mathbf{X}^{(miss)} = -\lambda \mathbb{H}[r(\mathbf{X}^{(miss)})], \lambda > 0, \quad (4)$$

where positive constant λ is a predefined regularization strength, and consequently we can formulate our NER cost functional for MDI as follows:

$$\mathcal{F}_{\text{NER}} := \mathbb{E}_{r(\mathbf{X}^{(miss)})} [\log \hat{p}(\mathbf{X}^{(miss)} | \mathbf{X}^{(obs)})] - \lambda \mathbb{H}[r(\mathbf{X}^{(miss)})]. \quad (5)$$

From a theoretical perspective, the NER term serves a critical role: The optimal $r(\mathbf{X}^{(miss)})$ inherently allows for infinite possibilities, potentially leading to a diversification of imputed values that could adversely affect accuracy. By incorporating this regularization term, we not only keep the objective direction but also reduce the diversification of samples by eliminating the entropy of $r(\mathbf{X}^{(miss)})$, which may result in an improvement in accuracy.

⁴according to reference [36], we mainly consider the Variance Preserving-SDE (VP-SDE), Variance Exploding-SDE (VE-SDE), and sub-Variance Preserving-SDE (sub-VP-SDE)

As demonstrated in [33, 53], we can directly incorporate Eq. (5) into the WGF framework and result in the following velocity field that drives the ODE in Section 2.3:

$$v_\tau = -\nabla_{\mathbf{X}^{(miss)}} \frac{\delta(-\mathcal{F}_{NER})}{\delta r(\mathbf{X}^{(miss)})} = [\nabla_{\mathbf{X}^{(miss)}} \log \hat{p}(\mathbf{X}^{(miss)} | \mathbf{X}^{(obs)}) + \lambda \nabla_{\mathbf{X}^{(miss)}} \log r(\mathbf{X}^{(miss)})], \quad (6)$$

where $\frac{\delta \mathcal{F}_{NER}}{\delta r(\mathbf{X}^{(miss)})}$ represents the first variation of \mathcal{F}_{NER} with respect to $r(\mathbf{X}^{(miss)})$. However, implementing this velocity field poses substantial challenges within both the ODE framework of the WGF and the SDE contexts.⁵ To alleviate this issue, we attempt to derive the expressions for model implementation based on the steepest ascent direction of functional gradient [10, 14, 8]. On this basis, we first derive the following ODE for the evolution of \mathcal{F}_{NER} along time τ :

Proposition 3.2. *The evolution of \mathcal{F}_{NER} along τ can be characterized by the following ODE, assuming that the boundary condition $\mathbb{E}_{r(\mathbf{X}^{(miss)}, \tau)} \{ \nabla_{\mathbf{X}^{(miss)}} \cdot [u(\mathbf{X}^{(miss)}, \tau) \log \hat{p}(\mathbf{X}^{(miss)} | \mathbf{X}^{(obs)})] \} = 0$ is satisfied for the velocity field $u(\mathbf{X}^{(miss)}, \tau)$:*

$$\frac{d\mathcal{F}_{NER}}{d\tau} = \mathbb{E}_{r(\mathbf{X}^{(miss)}, \tau)} [u^\top(\mathbf{X}^{(miss)}, \tau) \nabla_{\mathbf{X}^{(miss)}} \log \hat{p}(\mathbf{X}^{(miss)} | \mathbf{X}^{(obs)}) - \lambda \nabla_{\mathbf{X}^{(miss)}} \cdot u(\mathbf{X}^{(miss)}, \tau)]. \quad (7)$$

This boundary condition is achievable, for instance, when $\hat{p}(\mathbf{X}^{(miss)} | \mathbf{X}^{(obs)})$ is bounded, and the limit of the velocity field as the norm of $\mathbf{X}^{(miss)}$ approaches zero is zero ($\lim_{\|\mathbf{X}^{(miss)}\| \rightarrow 0} u(\mathbf{X}^{(miss)}, \tau) = 0$).

Despite the clarity provided by Eq. (7), practical model implementation remains a significant challenge due to the potential variability in the velocity field $u(\mathbf{X}^{(miss)}, \tau)$. To address this, we propose constraining the velocity within some specified function class \mathcal{V} [23, 47, 8], such that $u(\mathbf{X}^{(miss)}, \tau) \in \mathcal{V}$, which allows us to explore the steepest ascent direction systematically. *To obtain a closed-form and easily implementable expression*, we choose RKHS denoted by \mathcal{H} with RKHS norm $\|\cdot\|_{\mathcal{H}}$ to represent \mathcal{V} . On this basis, we have the following proposition for the expression of $u(\mathbf{X}^{(miss)}, \tau)$:

Proposition 3.3. *When the velocity field $u(\mathbf{X}^{(miss)}, \tau)$ is regularized by RKHS norm, the problem of finding the steepest gradient ascent direction can be formulated as follows:*

$$u(\mathbf{X}^{(miss)}, \tau) = \arg \max_{v(\mathbf{X}^{(miss)}, \tau) \in \mathcal{H}^d} \left\{ \mathbb{E}_{r(\mathbf{X}^{(miss)}, \tau)} [v^\top(\mathbf{X}^{(miss)}, \tau) \nabla_{\mathbf{X}^{(miss)}} \log \hat{p}(\mathbf{X}^{(miss)} | \mathbf{X}^{(obs)}) - \lambda \nabla_{\mathbf{X}^{(miss)}} \cdot v(\mathbf{X}^{(miss)}, \tau)] \right\} - \frac{1}{2} \|v(\mathbf{X}^{(miss)}, \tau)\|_{\mathcal{H}}^2. \quad (8)$$

The corresponding optimal solution is given by:

$$u(\mathbf{X}^{(miss)}, \tau) = \mathbb{E}_{r(\tilde{\mathbf{X}}^{(miss)}, \tau)} \left\{ \begin{array}{l} -\lambda \nabla_{\tilde{\mathbf{X}}^{(miss)}} \mathcal{K}(\mathbf{X}^{(miss)}, \tilde{\mathbf{X}}^{(miss)}) \\ + [\nabla_{\tilde{\mathbf{X}}^{(miss)}} \log \hat{p}(\tilde{\mathbf{X}}^{(miss)} | \mathbf{X}^{(obs)})]^\top \mathcal{K}(\mathbf{X}^{(miss)}, \tilde{\mathbf{X}}^{(miss)}) \end{array} \right\}, \quad (9)$$

where $\mathcal{K}(\mathbf{X}^{(miss)}, \tilde{\mathbf{X}}^{(miss)})$ is kernel function.

Importantly, since the missing value dimension is undefined, we did not specify the type signature of $\mathcal{K}(\mathbf{X}^{(miss)}, \tilde{\mathbf{X}}^{(miss)})$, and the expectation term $\mathbb{E}_{r(\tilde{\mathbf{X}}^{(miss)}, \tau)}$ can be efficiently estimated using Monte Carlo approximation. Leveraging this approach, the velocity field as outlined in Eq. (9) does not require explicit computation of proposal distribution $r(\mathbf{X}^{(miss)}, \tau)$. Consequently, we finally derive an easily implementable and closed-form imputation procedure for our KnewImp approach.

3.3 Modeling $p(\mathbf{X}^{(miss)} | \mathbf{X}^{(obs)})$ by $p(\mathbf{X}^{(miss)}, \mathbf{X}^{(obs)})$

As discussed in the previous subsection, accurately defining the conditional distribution $p(\mathbf{X}^{(miss)} | \mathbf{X}^{(obs)})$ is crucial for effectively simulating the ODE in Eq. (7) using the velocity field specified in Eq. (9). However, as previously noted, precise modeling of $p(\mathbf{X}^{(miss)} | \mathbf{X}^{(obs)})$ presents substantial challenges due to the diverse choices of masking matrices, which critically influence the efficacy of model training.

⁵Detailed analysis regarding the implementation challenges is provided in Appendix B.1.

To circumvent this difficulty, a practical approach involves substituting the conditional distribution $p(\mathbf{X}^{(\text{miss})}|\mathbf{X}^{(\text{obs})})$ with the joint distribution $p(\mathbf{X}^{(\text{miss})}, \mathbf{X}^{(\text{obs})})$. By denoting $\mathbf{X}^{(\text{joint})} = (\mathbf{X}^{(\text{miss})}, \mathbf{X}^{(\text{obs})})$, we can accordingly redefine the velocity field u as follows:

$$u(\mathbf{X}^{(\text{joint})}, \tau) = \mathbb{E}_{r(\tilde{\mathbf{X}}^{(\text{joint})}, \tau)} \left\{ \begin{array}{l} -\lambda \nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \mathcal{K}(\mathbf{X}^{(\text{joint})}, \tilde{\mathbf{X}}^{(\text{joint})}) \\ + [\nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \log \hat{p}(\tilde{\mathbf{X}}^{(\text{joint})})]^\top \mathcal{K}(\mathbf{X}^{(\text{joint})}, \tilde{\mathbf{X}}^{(\text{joint})}) \end{array} \right\}, \quad (10)$$

where $\mathcal{K}(\mathbf{X}^{(\text{joint})}, \tilde{\mathbf{X}}^{(\text{joint})}) : \mathbb{R}^D \rightarrow \mathbb{R}^D$ is kernel function, and we use radial basis function kernel defined as $\mathcal{K}(\mathbf{X}, \tilde{\mathbf{X}}) := \exp(-\frac{\|\mathbf{X} - \tilde{\mathbf{X}}\|^2}{2h^2})$ with bandwidth h in this paper [56, 25, 26]. Based on this, $\nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{joint})})$ can be obtained according to the following equation:

$$\nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \log \hat{p}(\tilde{\mathbf{X}}^{(\text{joint})}) = \nabla_{\tilde{\mathbf{X}}^{(\text{joint})}} \log \hat{p}(\tilde{\mathbf{X}}^{(\text{joint})}) \odot (\mathbb{1}_{N \times D} - \mathbf{M}) + 0 \times \mathbf{M}. \quad (11)$$

As such, a pertinent question arises: What's the relationship between Eqs. (9) and (10)? Interestingly, these formulations are identical. In light of this, the remainder of this subsection is dedicated to demonstrating that the velocity field associated with the cost functional in Eq. (10) does not compromise the optimization of \mathcal{F}_{NER} . To support this assertion, we present the following proposition:

Proposition 3.4. *Assume that the proposal distribution $r(\mathbf{X}^{(\text{joint})})$ is factorized by $r(\mathbf{X}^{(\text{joint})}) := r(\mathbf{X}^{(\text{miss})})p(\mathbf{X}^{(\text{obs})})$. The cost functional associated with the joint distribution is defined by the following equation:*

$$\mathcal{F}_{\text{joint-NER}} := \mathbb{E}_{r(\mathbf{X}^{(\text{joint})})} [\log \hat{p}(\mathbf{X}^{(\text{joint})})] - \lambda \mathbb{H}[r(\mathbf{X}^{(\text{joint})})], \quad (12)$$

which leads to the velocity field delineated in Eq. (10) and establishes $\mathcal{F}_{\text{joint-NER}}$ as a lower bound for \mathcal{F}_{NER} , with the difference being a constant (i.e., $\mathcal{F}_{\text{joint-NER}} = \mathcal{F}_{\text{NER}} - \text{const}, \text{const} \geq 0$).

Based on this, the following corollary can be obtained:

Corollary 3.5. *The following equation holds:*

$$u(\mathbf{X}^{(\text{joint})}, \tau) = u(\mathbf{X}^{(\text{miss})}, \tau). \quad (13)$$

Building on these foundations, the imputed value is obtained by simulating the following ODE:

$$\frac{d\mathbf{X}^{(\text{miss})}}{d\tau} = u(\mathbf{X}^{(\text{joint})}, \tau). \quad (14)$$

For simplicity, in this paper, we simulate this ODE by forward Euler's method with step size η .⁶

To date, our primary objective has been to determine the estimation of score function $\nabla_{\mathbf{X}^{(\text{joint})}} \log \hat{p}(\mathbf{X}^{(\text{joint})})$. To achieve this, we employ Denoising Score Matching (DSM) [12, 39] to parameterize $\nabla_{\mathbf{X}^{(\text{joint})}} \log \hat{p}(\mathbf{X}^{(\text{joint})})$ using a neural network with θ as parameter set. We design the learning objective to minimize the discrepancy between the actual score and the model's predicted score after introducing Gaussian noise to the clean $\mathbf{X}^{(\text{joint})}$ as $\hat{\mathbf{X}}^{(\text{joint})}$:

$$\mathcal{L}_{\text{DSM}} := \frac{1}{2} \mathbb{E}_{q_\sigma(\hat{\mathbf{X}}^{(\text{joint})}|\mathbf{X}^{(\text{joint})})} [\|\nabla_{\hat{\mathbf{X}}^{(\text{joint})}} \log \hat{p}(\mathbf{X}^{(\text{joint})}) - \nabla_{\hat{\mathbf{X}}^{(\text{joint})}} \log q_\sigma(\hat{\mathbf{X}}^{(\text{joint})}|\mathbf{X}^{(\text{joint})})\|^2]. \quad (15)$$

Notably, σ is variance scale, $\hat{\mathbf{X}}^{(\text{joint})}$ is obtained by $\hat{\mathbf{X}}^{(\text{joint})} = \mathbf{X}^{(\text{joint})} + \epsilon, \epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$, and $\nabla_{\hat{\mathbf{X}}^{(\text{joint})}} \log q_\sigma(\hat{\mathbf{X}}^{(\text{joint})}|\mathbf{X}^{(\text{joint})}) = -\frac{\hat{\mathbf{X}}^{(\text{joint})} - \mathbf{X}^{(\text{joint})}}{\sigma^2}$. Once $\nabla_{\mathbf{X}^{(\text{joint})}} \log \hat{p}(\mathbf{X}^{(\text{joint})})$ is trained, we can obtain the imputation value by simulating the differential equation based on Eq. (10).

3.4 Overall Architecture of KnewImp

Fig. 1 presents the architecture of our KnewImp approach, which consists of two parts namely 'Impute' and 'Estimate'. The 'Impute' part alleviates the missing data imputation as an ODE simulation problem within WGF framework, and the imputed matrix is obtained by simulating the velocity field as per Eq. (10). Meanwhile, since the velocity field requires the computation of the score function of the joint data, the 'Estimate' part serves for estimating the score function. By alternatively repeating these two parts, we can finally obtain the imputed value. To better delineate the KnewImp approach, we summarize the corresponding algorithms in Appendix C.2.

⁶please see Appendix C.1 for more detailed information about forward Euler's method.

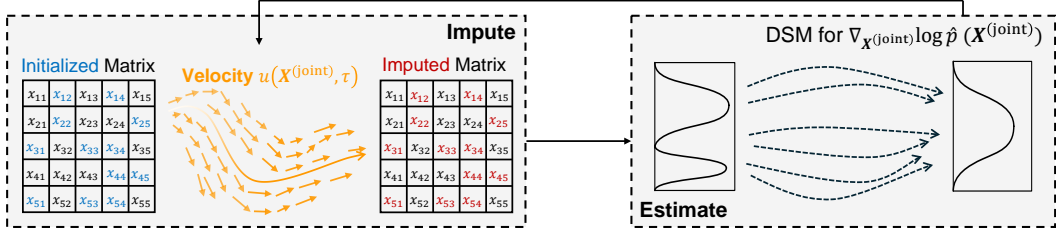


Figure 1: The illustration of KnewImp. The left part indicates we impute the missing value by WGF, and the right part indicates we use DSM to estimate $\log p(\mathbf{X}^{(\text{miss})})$.

4 Experiments

4.1 Experimental Setup

Datasets: Eight real-world datasets from [UCI repository](#) are chosen to validate the efficacy of our KnewImp approach. Detailed information for these datasets and the missing scenario simulation method is provided in [Appendix D.1](#).

Baselines: For fairness, we mainly consider the following models as baseline models: DMs-based approaches: conditional score-based diffusion models for Tabular Data (CSDI_T) [38], MissDiff [29]; Non-DMs: Optimal Transport Imputer (Sink) [28], Transform Distribution Matching (TDM, state-of-the-art) [54], Generative Adversarial Imputation Nets (GAIN) [51], Missing Data Importance-Weighted Autoencoder (MIWAE) [27], and Missing data Imputation Refinement And Causal Learning (MIRACLE) [19]. Concerning experimental details are given in [Appendix D.2](#).

Evaluation Metric: We choose the mean absolute error (abbreviated to ‘MAE’) and squared Wasserstein distance (abbreviated to ‘Wass’) as evaluation metrics.⁷

4.2 Baseline Comparison Results

Baseline comparison results are given in [Table 1](#), and the following observations can be given

- Models with neural architectures such as MIRACLE, MIWAE, and TDM demonstrate superior performance compared to models lacking such architectures. This observation suggests that integrating neural networks into MDI tasks can significantly enhance model performance.
- DM-based imputation approaches generally perform worse than other MDI methods. This outcome indicates that despite the incorporation of complex nonlinear neural architectures to boost performance, employing diversification-oriented generative approaches may not align well with the precision requirements of MDI tasks.
- Our proposed KnewImp method consistently ranks as the best or second-best across most comparisons. Notably, KnewImp significantly outperforms other DM-based MDI approaches, underscoring the effectiveness of our analytical enhancements and innovations in [Sections 3.1 to 3.3](#).

4.3 Ablation Study Results

In this subsection, we conduct the ablation study to assess the contributions of two key components in our KnewImp approach: the NER term and the joint modeling strategy (referred to as ‘Joint’). The results of this study are detailed in [Table 2](#). Analysis of the data between the second and last rows of [Table 2](#) reveals that, in the absence of the NER, the proposal distribution $r(\mathbf{X}^{(\text{miss})})$ may become pathological, leading to diminished model performance. Additionally, when comparing results from the first, third, and last rows, it becomes evident that modeling the joint distribution directly, rather than inferring it from the conditional distribution, significantly enhances model performance. This finding underscores the effectiveness of the strategies we have implemented, as discussed in [Section 3.3](#). Overall, the ablation study underscores the critical roles of both the NER term and the joint distribution learning strategy in promoting the performance of KnewImp.

⁷Detailed information about these two metrics is given in [Appendix D.3](#). We also report concerning results under the MNAR scenario in [Appendix E](#) for completeness.

Table 1: Performance of MAE and Wass metrics at 30% missing rate, and ‘**’ marks that KnewImp outperforms significantly at p -value < 0.05 over paired samples t -test. Best results are **bolded** and the second best results are underliend. Other results like standard deviation are given in appendix.

Scenario	Model	BT		BCD		CC		CBV		IS		PK		QB		WQW	
		MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass
MAR	CSDL T	0.93*	3.44*	0.92*	18.2*	0.85*	2.82*	0.81*	3.86*	0.70*	16.9*	0.99*	15.9*	0.65*	20.1*	0.77*	4.13*
	MissDiff	0.85*	2.20*	0.91*	16.5*	0.87*	1.59*	0.83*	3.87*	0.72*	13.3*	0.92*	17.1*	0.63*	26.3*	0.75*	6.88*
	GAIN	0.75*	0.65*	0.54*	1.64*	0.75*	0.67*	0.68*	0.68*	0.56*	1.88*	0.59*	1.90*	0.65*	5.05*	0.68*	0.87*
	MIRACLE	0.62*	0.38	0.55*	1.92*	0.43	0.25	0.55*	0.46*	3.39*	35.1*	4.14*	34.1*	0.46	2.87*	0.51*	0.56
	MIWAE	0.64	0.53	0.52*	1.54*	0.76*	0.64*	0.82*	0.92*	0.50*	1.87*	0.65*	1.98*	0.55*	5.05*	0.62*	0.75*
	Sink	0.87*	0.92*	0.92*	3.84*	0.88*	0.83*	0.84*	0.98*	0.75*	2.43*	0.94*	3.61*	0.65*	4.71*	0.76*	1.04*
	TDM	0.83*	0.89*	0.83*	3.47*	0.81*	0.73*	0.76*	0.85*	0.62*	1.96*	0.86*	3.36*	0.59*	4.46*	0.73*	0.99*
	KnewImp	0.52	0.38	0.34	0.82	<u>0.35</u>	<u>0.25</u>	0.31	0.20	0.39	1.31	0.44	1.21	0.45	<u>3.50</u>	0.46	0.55
MCAR	CSDL T	0.73*	1.93*	0.73*	15.5*	0.85*	2.71*	0.83*	3.79*	0.76*	15.2*	0.72*	12.4*	0.57*	19.9*	0.78*	4.11*
	MissDiff	0.72*	1.62*	0.73*	14.4*	0.84*	1.23*	0.82*	3.31*	0.75*	13.01*	0.71*	14.1*	0.56*	19.7*	0.76*	4.95*
	GAIN	0.72*	0.39*	0.38*	1.41*	0.78*	0.73*	0.72*	0.99*	0.57*	3.72*	0.46*	1.70	0.42*	3.62	0.73*	1.14*
	MIRACLE	0.52	0.15*	0.44*	1.94*	0.53*	0.35	0.61*	0.72*	2.99*	52.9*	3.38*	42.8*	0.35	2.71*	0.56*	0.75
	MIWAE	0.58*	0.24	0.50*	2.55*	0.76*	0.69*	0.83*	1.24*	0.64*	4.95*	0.51*	2.05*	0.48*	5.87*	0.67*	0.95*
	Sink	0.73*	0.48*	0.75*	4.39*	0.84*	0.85*	0.82*	1.27*	0.75*	4.94*	0.74*	3.36*	0.61*	5.92*	0.76*	1.25*
	TDM	0.68*	0.42*	0.63*	3.57*	0.77*	0.75*	0.77*	1.15*	0.66*	4.20*	0.64*	2.89*	0.52*	5.34*	0.74*	1.20*
	KnewImp	0.48	<u>0.18</u>	0.25	0.80	0.47	0.34	0.42	0.44	0.44	3.05	0.32	1.01	0.34	3.66	0.53	<u>0.76</u>

Table 2: Ablation Study Results with missing rate at 30%, and ‘**’ marks that KnewImp outperforms significantly at p -value < 0.05 over paired samples t -test. Best results are **bolded**.

Missing	NER	Joint	BT		BCD		CC		CBV		IS		PK		QB		WQW	
			MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass
MAR	✗	✗	0.96*	3.82*	1.05*	20.2*	1.04*	5.47*	0.86*	5.81*	0.67*	20.2*	1.06*	15.6*	0.72*	22.5*	0.79*	6.49*
	✗	✓	0.54	0.42	0.34	0.82	0.61*	0.40*	0.58*	0.47*	0.43*	1.34	0.46*	1.25*	0.47*	3.56*	0.55*	0.64*
	✓	✓	0.96*	3.83*	1.05*	20.3*	1.04*	5.49*	0.86*	5.83*	0.67*	20.2*	1.06*	15.7*	0.72*	22.5*	0.79*	6.51*
	✓	✓	0.52	0.38	0.34	0.82	0.35	0.25	0.31	0.20	0.39	1.31	0.44	1.21	0.45	3.50	0.46	0.55
MCAR	✗	✗	0.72*	2.11*	0.74*	16.7*	0.85*	3.72*	0.83*	5.22*	0.74*	18.4*	0.71*	12.7*	0.58*	20.1*	0.76*	5.57*
	✗	✓	0.52*	0.17*	0.25	0.79	0.62*	0.46*	0.61*	0.71*	0.46	3.05	0.34	1.09	0.36*	3.74*	0.58*	0.82*
	✓	✓	0.72*	2.12*	0.73*	16.8*	0.86*	3.73*	0.83*	5.24*	0.74*	18.4*	0.71*	12.8*	0.58*	20.1*	0.76*	5.60*
	✓	✓	0.48	0.18	0.25	0.80	0.47	0.34	0.42	0.44	0.44	3.05	0.32	1.01	0.34	3.66	0.53	0.76

4.4 Sensitivity Analysis

In this subsection, we analyze the impact of key hyperparameters within the KnewImp approach, including the bandwidth h of the RBF kernel function, the hidden units HU_{score} in the score network, the weight λ of the NER term, and the discretization step size η for simulating the ODE defined in Eq. (10). The profound influence of these hyperparameters on learning objectives and overall performance is substantiated by the experimental results presented in Figure 2. Initially, we explore the effects of varying the bandwidth h . We observe that an increase in bandwidth correlates with a decrease in imputation accuracy. For instance, as the bandwidth increases from 0.5 to 2.0, the MAE and Wasserstein distance escalate from 0.35 and 0.25 to 0.82 and 0.74, respectively. This trend suggests that excessive bandwidth can lead to an over-smoothed velocity field, expanding the exploration space of the distribution $r(\mathbf{X}^{(joint)})$ excessively and failing to adequately ‘concentrate’ this distribution, ultimately diminishing performance. Subsequently, we examine changes in the score network’s hidden units. Increasing the hidden units from 256 to 512 appears to decrease imputation accuracy, likely due to overfitting issues associated with larger neural networks. Next, we adjust the strength of the NER term and find that increasing its intensity generally improves imputation accuracy. This supports the necessity of the NER term, further validating its effectiveness. Lastly, we investigate the discretization step size for the ODE. We find that accuracy initially increases with smaller step sizes but then decreases. This pattern is consistent with ODE simulation behavior, where smaller step sizes require longer to converge, potentially resulting in lower accuracy within a predefined time. Conversely, larger step sizes increase discretization errors, adversely affecting accuracy as well.

5 Related Works

5.1 Diffusion Models for Missing Data Imputation

The impressive ability of DMs to synthesize data has inspired extensive research into their application for MDI tasks [43, 50]. Among the pioneering efforts, the Conditional Score-based Diffusion models for Imputation (CSDI) [38] was the first to adapt diffusion models for time-series MDI,

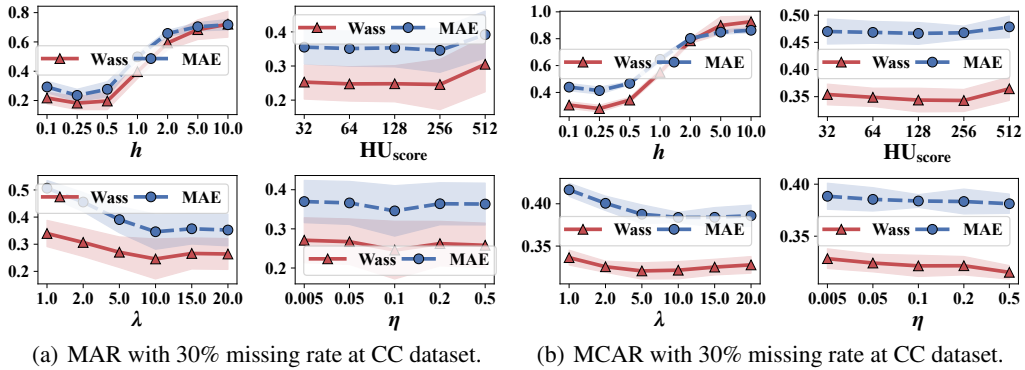


Figure 2: Parameter sensitivity of KnewImp on bandwidth for kernel function (h), hidden unit of score network HU_{score} , NER weight λ , and discretization step η for Eq. (10) on CC dataset. Mean values and one standard deviations from mean are represented by scatters and shaded area, respectively.

substituting the score function with a conditional distribution and pioneering a novel model training strategy by masking parts of the observational data. Building on this, to address categorical data in tabular datasets, CSDI_T [55] introduced an embedding layer within the feature extractor. To enhance inference efficiency, the conditional Schrödinger bridge method for probabilistic time series imputation proposed modeling the diffusion process as a Schrödinger bridge [6]. Meanwhile, MissDiff [29] sought to bypass the masking mechanism typically used during score function training by using missing data information as a mask matrix to improve the training process.

Despite these advancements from the perspective of feature extraction module [1, 48], loss function [29], and model inference approach [44], the reconciliation of the inherent diversity-seeking nature of DMs’ generative processes and the accuracy-centric demands of MDI task remains underexplored. To our knowledge, this paper is the first to elucidate the relationship between DM generative processes and MDI tasks from an optimization perspective (Section 3.1). Based on these insights, we further propose our KnewImp approach, which prioritizes MDI accuracy (Section 3.2).

5.2 Modeling Conditional Distribution by Joint Distribution

Modeling conditional distribution as joint distribution remains an opening question and has a broad potential for application [52, 4, 37]. Conditional sliced WGF [9] first empirically validated that the velocity field of joint distribution and conditional distribution are identical when choosing sliced Wasserstein distance as cost functional. After that, reference [17] extended this relationship and derived the relationship between conditional and joint distribution in various discrepancy metrics like f-divergence, Wasserstein distance, and integral probability metrics. On this basis, reference [11] further theoretically proved the equivalence of velocity fields for conditional and joint distribution.

However, the objective of KnewImp does not belong to any kind of discrepancy metric [17]. The most similar discrepancy metric is Kullback Leiber (KL) divergence ($-\int r(x) \frac{\log r(x)}{p(x)} dx$). Notably, KL divergence contains diversification-encouraging ‘positive’ entropy $\mathbb{H}[r(x)]$ as the regularization term, and the regularization term in our study is diversification-discouraging ‘negative’ entropy (i.e., $-\mathbb{H}[r(x)]$), and thus more than directly applying these results to our research is needed. On this basis, our theoretical contribution proves that this joint distribution modeling approach can still be applied when the functional is regularized by the negative entropy (Section 3.3).

6 Conclusions

Existing DM-based MDI approaches face two critical issues that may hinder model performance: inaccurate imputation and difficult training. The first issue arises from the inherent conflict between the diversification-oriented generative process of DMs and the accuracy-focused demands of the MDI tasks. The second issue stems from the selection complexities of the masking matrix to facilitate conditional distribution between missing and observed data. To this end, this study initially applied the WGF framework to analyze DM-based MDI tasks, elucidating the relationship between the optimization objectives of DMs’ generative process and the MDI task, and answered the reason for inaccurate imputation issue from the perspective of optimization. On this basis, we proposed our KnewImp approach by redesigning a novel effect cost functional and developing the corresponding DM-like imputation procedure within WGF and RKHS. Furthermore, we proved that another joint-

distribution-related cost functional can result in the same imputation procedure, which naturally copes with the need for a masking matrix during model training. Finally, we conducted extensive experiments and demonstrated that the KnewImp approach can mitigate the abovementioned issues and achieve better performance than prevalent baseline models.

References

- [1] Juan Lopez Alcaraz and Nils Strodthoff. Diffusion-based time series imputation and forecasting with structured state space models. Transactions on Machine Learning Research, 2022.
- [2] Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré. Gradient flows: in metric spaces and in the space of probability measures. Springer Science & Business Media, 2005.
- [3] Brian DO Anderson. Reverse-time diffusion equation models. Stochastic Processes and their Applications, 12(3):313–326, 1982.
- [4] Jannis Chemseddine, Paul Hagemann, Christian Wald, and Gabriele Steidl. Conditional wasserstein distances with applications in bayesian ot flow matching. arXiv preprint arXiv:2403.18705, 2024.
- [5] Changyou Chen, Ruiyi Zhang, Wenlin Wang, Bai Li, and Liqun Chen. A unified particle-optimization framework for scalable bayesian sampling. arXiv preprint arXiv:1805.11659, 2018.
- [6] Yu Chen, Wei Deng, Shikai Fang, Fengpei Li, Nicole Tianjiao Yang, Yikai Zhang, Kashif Rasul, Shandian Zhe, Anderson Schneider, and Yuriy Nevmyvaka. Provably convergent schrödinger bridge with applications to probabilistic time series imputation. In International Conference on Machine Learning, pages 4485–4513. PMLR, 2023.
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Tamar Solorio, editors, Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.
- [8] Hanze Dong, Xi Wang, LIN Yong, and Tong Zhang. Particle-based variational inference with preconditioned functional gradient flow. In The Eleventh International Conference on Learning Representations, pages 1–26, 2022.
- [9] Chao Du, Tianbo Li, Tianyu Pang, Shuicheng Yan, and Min Lin. Nonparametric generative modeling with conditional sliced-wasserstein flows. In International Conference on Machine Learning, pages 8565–8584. PMLR, 2023.
- [10] Jerome H Friedman. Greedy function approximation: a gradient boosting machine. Annals of statistics, pages 1189–1232, 2001.
- [11] Paul Hagemann, Johannes Hertrich, Fabian Altekrüger, Robert Beinert, Jannis Chemseddine, and Gabriele Steidl. Posterior sampling based on gradient flows of the mmd with negative distance kernel. arXiv preprint arXiv:2310.03054, 2023.
- [12] Aapo Hyvärinen and Peter Dayan. Estimation of non-normalized statistical models by score matching. Journal of Machine Learning Research, 6(4), 2005.
- [13] Daniel Jarrett, Bogdan C Cebere, Tennison Liu, Alicia Curth, and Mihaela van der Schaar. Hyperimpute: Generalized iterative imputation with automatic model selection. In International Conference on Machine Learning, pages 9916–9937. PMLR, 2022.
- [14] Rie Johnson and Tong Zhang. A framework of composite functional gradient methods for generative adversarial models. IEEE Transactions on Pattern Analysis and Machine Intelligence, 43(1):17–32, 2021. doi: 10.1109/TPAMI.2019.2924428.
- [15] Richard Jordan, David Kinderlehrer, and Felix Otto. The variational formulation of the fokker-planck equation. SIAM journal on mathematical analysis, 29(1):1–17, 1998.
- [16] Valentin Khrulkov, Gleb Ryzhakov, Andrei Chertkov, and Ivan Oseledets. Understanding ddpm latent codes through optimal transport. In The Eleventh International Conference on Learning Representations, pages 1–15, 2022.

- [17] Young-geun Kim, Kyungbok Lee, and Myunghee Cho Paik. Conditional wasserstein generator. IEEE Transactions on Pattern Analysis and Machine Intelligence, 45(6):7208–7219, 2023. doi: 10.1109/TPAMI.2022.3220965.
- [18] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In International Conference on Learning Representations (ICLR), pages 1–8, 2015.
- [19] Trent Kyono, Yao Zhang, Alexis Bellot, and Mihaela van der Schaar. Miracle: Causally-aware imputation via learning missing data mechanisms. Advances in Neural Information Processing Systems, 34:23806–23817, 2021.
- [20] Haoxuan Li, Quanyu Dai, Yuru Li, Yan Lyu, Zhenhua Dong, Xiao-Hua Zhou, and Peng Wu. Multiple robust learning for recommendation. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 37, pages 4417–4425, 2023.
- [21] Haoxuan Li, Kunhan Wu, Chunyuan Zheng, Yanghao Xiao, Hao Wang, Zhi Geng, Fuli Feng, Xiangnan He, and Peng Wu. Removing hidden confounding in recommendation: a unified multi-task learning approach. Advances in Neural Information Processing Systems, 36:1–13, 2024.
- [22] Yingzhen Li and Richard E. Turner. Gradient estimators for implicit models. In International Conference on Learning Representations (ICLR), 2018. URL <https://openreview.net/forum?id=SJi9W0eRb>.
- [23] Chang Liu, Jingwei Zhuo, Pengyu Cheng, Ruiyi Zhang, and Jun Zhu. Understanding and accelerating particle-based variational inference. In International Conference on Machine Learning, pages 4082–4092. PMLR, 2019.
- [24] Diju Liu, Yalin Wang, Chenliang Liu, Xiaofeng Yuan, Kai Wang, and Chunhua Yang. Scope-free global multi-condition-aware industrial missing data imputation framework via diffusion transformer. IEEE Transactions on Knowledge and Data Engineering, pages 1–12, 2024. doi: 10.1109/TKDE.2024.3392897.
- [25] Qiang Liu and Dilin Wang. Stein variational gradient descent: A general purpose bayesian inference algorithm. Advances in neural information processing systems, 29, 2016.
- [26] Weiming Liu, Xiaolin Zheng, Jiajie Su, Longfei Zheng, Chaochao Chen, and Mengling Hu. Contrastive proxy kernel stein path alignment for cross-domain cold-start recommendation. IEEE Transactions on Knowledge and Data Engineering, 35(11):11216–11230, 2023. doi: 10.1109/TKDE.2022.3233789.
- [27] Pierre-Alexandre Mattei and Jes Frelsen. MIWAE: Deep generative modelling and imputation of incomplete data sets. In International conference on machine learning, pages 4413–4423. PMLR, 2019.
- [28] Boris Muzellec, Julie Josse, Claire Boyer, and Marco Cuturi. Missing data imputation using optimal transport. In International Conference on Machine Learning, pages 7130–7140. PMLR, 2020.
- [29] Yidong Ouyang, Liyan Xie, Chongxuan Li, and Guang Cheng. Missdiff: Training diffusion models on tabular data with missing values. In ICML 2023 Workshop on Structured Probabilistic Inference & Generative Modeling, 2023.
- [30] Prajit Ramachandran, Barret Zoph, and Quoc V Le. Searching for activation functions. arXiv preprint arXiv:1710.05941, 2017.
- [31] Donald B Rubin. Inference and missing data. Biometrika, 63(3):581–592, 1976.
- [32] Walter Rudin et al. Principles of mathematical analysis, volume 3. McGraw-hill New York, 1964.
- [33] Filippo Santambrogio. {Euclidean, Metric, and Wasserstein} gradient flows: an overview. Bulletin of Mathematical Sciences, 7:87–154, 2017.

- [34] Simo Särkkä and Arno Solin. Applied stochastic differential equations, volume 10. Cambridge University Press, 2019.
- [35] Yang Song, Sahaj Garg, Jiaxin Shi, and Stefano Ermon. Sliced score matching: A scalable approach to density and score estimation. In Uncertainty in Artificial Intelligence, pages 574–584. PMLR, 2020.
- [36] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In International Conference on Learning Representations, pages 1–36, 2020.
- [37] Ryan Szeto, Ximeng Sun, Kunyi Lu, and Jason J. Corso. A temporally-aware interpolation network for video frame inpainting. IEEE Transactions on Pattern Analysis and Machine Intelligence, 42(5):1053–1068, 2020. doi: 10.1109/TPAMI.2019.2951667.
- [38] Yusuke Tashiro, Jiaming Song, Yang Song, and Stefano Ermon. CSDI: Conditional score-based diffusion models for probabilistic time series imputation. Advances in Neural Information Processing Systems, 34:24804–24816, 2021.
- [39] Pascal Vincent. A connection between score matching and denoising autoencoders. Neural computation, 23(7):1661–1674, 2011.
- [40] Dilin Wang and Qiang Liu. Nonlinear stein variational gradient descent for learning diversified mixture models. In International Conference on Machine Learning, pages 6576–6585. PMLR, 2019.
- [41] Fangyikang Wang, Huminhao Zhu, Chao Zhang, Hanbin Zhao, and Hui Qian. Gad-pvi: A general accelerated dynamic-weight particle-based variational inference framework. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 38, pages 15466–15473, 2024.
- [42] Hao Wang, Tai-Wei Chang, Tianqiao Liu, Jianmin Huang, Zhichao Chen, Chao Yu, Ruopeng Li, and Wei Chu. Escm2: entire space counterfactual multi-task model for post-click conversion rate estimation. In Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, pages 363–372, 2022.
- [43] Jun Wang, Wenjie Du, Wei Cao, Keli Zhang, Wenjia Wang, Yuxuan Liang, and Qingsong Wen. Deep learning for multivariate time series imputation: A survey. arXiv preprint arXiv:2402.04059, 2024.
- [44] Xu Wang, Hongbo Zhang, Pengkun Wang, Yudong Zhang, Binwu Wang, Zhengyang Zhou, and Yang Wang. An observed value consistent diffusion model for imputing missing values in multivariate time series. In Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, pages 2409–2418, 2023.
- [45] Yifei Wang and Wuchen Li. Accelerated information gradient flow. Journal of Scientific Computing, 90:1–47, 2022.
- [46] Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In Proceedings of the 28th international conference on machine learning (ICML-11), pages 681–688. Citeseer, 2011.
- [47] Veit David Wild, Sahra Ghalebikesabi, Dino Sejdinovic, and Jeremias Knoblauch. A rigorous link between deep ensembles and (variational) bayesian methods. Advances in Neural Information Processing Systems, 36, 2024.
- [48] Jingwen Xu, Fei Lyu, and Pong C Yuen. Density-aware temporal attentive step-wise diffusion model for medical time series imputation. In Proceedings of the 32nd ACM International Conference on Information and Knowledge Management, pages 2836–2845, 2023.
- [49] Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications. ACM Computing Surveys, 56(4):1–39, 2023.

- [50] Yiyuan Yang, Ming Jin, Haomin Wen, Chaoli Zhang, Yuxuan Liang, Lintao Ma, Yi Wang, Chenghao Liu, Bin Yang, Zenglin Xu, et al. A survey on diffusion models for time series and spatio-temporal data. arXiv preprint arXiv:2404.18886, 2024.
- [51] Jinsung Yoon, James Jordon, and Mihaela Schaar. GAIN: Missing data imputation using generative adversarial nets. In International conference on machine learning, pages 5689–5698. PMLR, 2018.
- [52] Shipeng Yu, Kai Yu, Volker Tresp, Hans-Peter Kriegel, and Mingrui Wu. Supervised probabilistic principal component analysis. In Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining, pages 464–473, 2006.
- [53] Chao Zhang, Zhijian Li, Xin Du, and Hui Qian. Dpvi: A dynamic-weight particle-based variational inference framework. In Lud De Raedt, editor, Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22, pages 4900–4906. International Joint Conferences on Artificial Intelligence Organization, 7 2022. doi: 10.24963/ijcai.2022/679. URL <https://doi.org/10.24963/ijcai.2022/679>. Main Track.
- [54] He Zhao, Ke Sun, Amir Dezfouli, and Edwin V Bonilla. Transformed distribution matching for missing value imputation. In International Conference on Machine Learning, pages 42159–42186. PMLR, 2023.
- [55] Shuhan Zheng and Nontawat Charoenphakdee. Diffusion models for missing value imputation in tabular data. In NeurIPS 2022 First Table Representation Workshop, 2022.
- [56] Ding-Xuan Zhou. Derivative reproducing properties for kernel methods in learning theory. Journal of computational and Applied Mathematics, 220(1-2):456–463, 2008.

Appendix Contents

A Detailed Preliminaries of Wasserstein Gradient Flow	16
B Theoretical Analysis	17
B.1 Implementation Difficulty of Velocity Field	17
B.2 Proof & Discussions of Propositions & Corollaries	18
C Detailed Information for KnewImp Implementation	25
C.1 Forward Euler’s Method for ODE Simulation	25
C.2 Algorithms for KnewImp	26
D Detailed Information for Experiments	27
D.1 Background & Simulation of Missing Data	27
D.2 Training Protocols of Different Models	28
D.3 Evaluation Protocols	28
E Additional Empirical Evidence	29
E.1 Additional Experimental Results with MNAR Scenario	29
E.2 Time Complexity Analysis	31
E.3 Convergence Analysis	32
E.4 Baseline Comparison Vary Different Missing Rates and Scenarios	36
F Limitations & Future Directions and Broader Impact	43
F.1 Limitations & Future Directions	43
F.2 Broader Impact Statement	43

Appendix A Detailed Preliminaries of Wasserstein Gradient Flow

In this section, we want to introduce the WGF technique and its application scenarios to better understand this paper. Before introduction, the following concepts are listed to better understand the WGF framework:

1. **Wasserstein Metric:** Let $\mathcal{P}_2(\mathbb{R}^D)$ represent the space of probability measures on \mathbb{R}^D that possess finite second moments. Formally, this is expressed as $\mathcal{P}_2(\mathbb{R}^D) = \{\mu \in \mathcal{M}(\mathbb{R}^D) \mid \int \|x\|^2 d\mu(x) < \infty\}$, where $\mathcal{M}(\mathbb{R}^D)$ denotes the set of all probability measures on \mathbb{R}^D . Considering any two probability measures $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^D)$, we define the Wasserstein- p distance between them as follows:

$$\mathcal{W}_p = \left(\inf_{\pi \in \Gamma(\mu, \nu)} \int_{\mathbb{R}^D \times \mathbb{R}^D} \|x - y\|^p d\pi(x, y) \right)^{\frac{1}{p}}. \quad (\text{A.1})$$

Here, $\Gamma(\mu, \nu)$ represents the collection of all joint distributions (couplings) between μ and ν . For every joint distribution $\pi \in \Gamma(\mu, \nu)$, it holds that $\mu(x) = \int_{\mathbb{R}^D} \pi(x, y) dy$ and $\nu(y) = \int_{\mathbb{R}^D} \pi(x, y) dx$. The integral on the right-hand side encapsulates the transportation cost in the optimal transport (OT) problem, framed by Kantorovich's formulation, where π^* denotes the optimal transportation plan.

Furthermore, leveraging Jensen's inequality facilitates demonstrating the monotonicity of the Wasserstein- p distance, affirming that for $1 \leq p \leq q$, the relationship $\mathcal{W}_p(\mu, \nu) \leq \mathcal{W}_q(\mu, \nu)$ invariably holds. Building on this principle, we can articulate the inner product within the measurable space $(\mathcal{P}_2(\mathbb{R}^D), \mathcal{W})$ as delineated below:

$$\langle \mu, \nu \rangle_{\mu_\tau} = \int_{\mathbb{R}^D} \langle \mu, \nu \rangle_{\mathbb{R}^D} d\mu_\tau \quad (\text{A.2})$$

2. **Gradient Flow in Wasserstein Space:** Consider a functional \mathcal{F} associated with $\mu \in \mathcal{P}_2(\mathbb{R}^D)$. Our objective is to identify the optimal μ that minimizes \mathcal{F} :

$$\min_{\mu \in \mathcal{P}_2(\mathbb{R}^D)} \mathcal{F}(\mu) + \text{const}. \quad (\text{A.3})$$

To facilitate the decrease of $\mathcal{F}(\mu)$, we introduce a velocity field $v_{\mu_\tau} : \mathbb{R}^D \rightarrow \mathbb{R}^D$ designed to expedite the reduction of $\mathcal{F}(\mu)$ as μ evolves under this field. Utilizing the chain rule yields:

$$\frac{d\mathcal{F}(\mu)}{d\tau} = \int \left\langle \nabla \frac{\delta \mathcal{F}}{\delta \mu_\tau}, v_{\mu_\tau} \right\rangle d\mu_\tau, \quad (\text{A.4})$$

where δ represents the first variation operator. To ensure the decrease of $\mathcal{F}(\mu)$, i.e., $\frac{d\mathcal{F}(\mu)}{d\tau} \leq 0$, the velocity field is defined as:

$$v_{\mu_\tau} = -\nabla \frac{\delta \mathcal{F}}{\delta \mu_\tau}. \quad (\text{A.5})$$

The decline of $\mathcal{F}(\mu)$ aligns with the continuity equation:

$$\frac{\partial \mu_\tau}{\partial \tau} = -\nabla \cdot (\mu_\tau v_{\mu_\tau}). \quad (\text{A.6})$$

Hence, the continuity equation Eq. (A.6), coupled with the velocity field articulated in Eq. (A.5), is recognized as the *Wasserstein Gradient Flow*, delineating the steepest descent in the Wasserstein space.

3. **Simulation of WGF & Sampling:** There are primarily two discretization techniques for the WGF: the forward scheme and the backward scheme.

- **Forward Scheme:** The forward scheme applies gradient descent within the Wasserstein space to identify the direction of the steepest descent. For an energy functional $\mathcal{F}(\mu_\tau)$ with a specified step size η , the update rule in the forward scheme is formulated as:

$$\mu_{t+1} = (\text{Id} - \nabla \frac{\delta \mathcal{F}}{\delta \mu_\tau})_{\#} \mu_\tau, \quad (\text{A.7})$$

facilitating an intuitive and direct update mechanism that emulates the gradient flow in the Euclidean space but transposed into the Wasserstein space.

- **Backward Scheme:** Conversely, the backward scheme, often referred to as the Jordan-Kinderlehrer-Otto (JKO) scheme [15], represents a more implicit discretization approach. It defines the subsequent distribution $\mu_{\tau+1}$ by solving an optimization problem that balances the energy decrease and the transportation cost. This scheme is mathematically denoted as:

$$\mu_{\tau+1} = \arg \min_{\mu \in \mathcal{P}_2(\mathbb{R}^D)} \mathcal{F}(\mu) + \frac{1}{2\eta} \mathcal{W}_2^2(\mu, \mu_\tau), \quad (\text{A.8})$$

thereby integrating the energy minimization and transport efficiency into a single variational problem that reflects the inherent structure of the Wasserstein space.

These schemes provide distinct yet complementary approaches to discretizing the dynamics defined by WGFs, offering different perspectives and tools for the analysis and computation of these flows.

Leveraging the WGF framework, if we designate the functional \mathcal{F} to be the KL divergence, it yields a particular formulation for the velocity field.

$$v_{\mu_\tau} = -\nabla \frac{\delta \mathbb{D}_{\text{KL}}(\mu_\tau \| p)}{\delta \mu_\tau} = \nabla \log p - \nabla \log \mu_\tau. \quad (\text{A.9})$$

On this basis, plug Eq. (A.9) into Eq. (A.6), we can get the following PDE:

$$\frac{\partial \mu_\tau}{\partial \tau} = -\nabla \cdot (\mu_\tau \nabla \log p) + \nabla \cdot \nabla \mu_\tau. \quad (\text{A.10})$$

According to Theorem 5.4 in reference [34], denote the random sample from distribution p as x , we can obtain the following SDE called Langevin equation [46] for implementing this gradient flow easily:

$$dx = \nabla_x \log p(x) d\tau + \sqrt{2} dW_\tau, \quad (\text{A.11})$$

where dW_τ is the Wiener process (also known as Brownian motion).

Appendix B Theoretical Analysis

B.1 Implementation Difficulty of Velocity Field

The difficulty of implementing velocity can be given from two perspectives, namely ODE-based implementation and SDE-based implementation, to the best of our knowledge. In this subsection, we want to discuss these two implementation approaches in detail.

ODE-based Implementation:

1. **WGF framework:** According to the continuity equation, we can obtain the following velocity field:

$$\frac{d\mathbf{X}^{(\text{miss})}}{d\tau} \stackrel{\text{(i)}}{=} v_\tau(\mathbf{X}^{(\text{miss})}) \stackrel{\text{(ii)}}{=} -[\nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) + \lambda \nabla_{\mathbf{X}^{(\text{miss})}} \log r(\mathbf{X}^{(\text{miss})})], \quad (\text{B.1})$$

where (i) is based on Section 2.3, and (ii) is based on Eq. (6). The expression of velocity field involves the computation of density term $r(\mathbf{X}^{(\text{miss})})$ [22, 5], which is intractable during practice. Based on this, we conclude that implementing this velocity field within the WGF framework is difficult.

2. **Probability flow ODE:** According to reference [36], if we directly plug Eq. (6) into the FPK equation, we can get the following PDE:

$$\begin{aligned}
& \frac{\partial r(\mathbf{X}^{(\text{miss})})}{\partial \tau} \\
&= -(\mathbf{X}^{(\text{miss})} r(\mathbf{X}^{(\text{miss})})) \\
&= -\left[\nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) r(\mathbf{X}^{(\text{miss})}) \right] - \lambda \nabla \cdot \nabla r(\mathbf{X}^{(\text{miss})}) \\
&\quad - \left[\nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) r(\mathbf{X}^{(\text{miss})}) \right] - \lambda \nabla \cdot \nabla r(\mathbf{X}^{(\text{miss})}) \\
&= \qquad \qquad \qquad + \frac{1}{2} \sigma_\tau^2 \nabla \cdot \nabla r(\mathbf{X}^{(\text{miss})}) - \frac{1}{2} \sigma_\tau^2 \nabla \cdot \nabla r(\mathbf{X}^{(\text{miss})}) \\
&= -\left\{ \left[\nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) + \left(\lambda + \frac{1}{2} \sigma_\tau^2 \right) \nabla \log r(\mathbf{X}^{(\text{miss})}) \right] r(\mathbf{X}^{(\text{miss})}) \right\} \\
&\quad + \frac{1}{2} \sigma_\tau^2 \nabla \cdot \nabla r(\mathbf{X}^{(\text{miss})}).
\end{aligned} \tag{B.2}$$

When we set σ_τ as 0, we can find that the corresponding ODE is Eq. (6), where we are obliged to compute the intractable density $r(\mathbf{X}^{(\text{miss})})$.

SDE-based Implementation:

If we plug Eq. (6) into the FPK equation, the corresponding PDE can be given as follows:

$$\begin{aligned}
& \frac{\partial r(\mathbf{X}^{(\text{miss})})}{\partial \tau} \\
&= -(v_\tau(\mathbf{X}^{(\text{miss})}) r(\mathbf{X}^{(\text{miss})})) \\
&= -\left[\nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) r(\mathbf{X}^{(\text{miss})}) \right] - \lambda \nabla \cdot \nabla r(\mathbf{X}^{(\text{miss})}),
\end{aligned} \tag{B.3}$$

where the coefficient before Laplacian operator $\nabla \cdot \nabla$ is -1 . To the best of our knowledge, this structure makes deriving a corresponding SDE impossible by current approaches.

B.2 Proof & Discussions of Propositions & Corollaries

Proposition (3.1). *Within WGF framework, DM-based MDI approaches can be viewed as finding the imputed values $\mathbf{X}^{(\text{imp})}$ that maximize the following objective:*

$$\arg \max_{\mathbf{X}^{(\text{miss})}} \mathbb{E}_{r(\mathbf{X}^{(\text{miss})})} [\log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})})] + \psi(\mathbf{X}^{(\text{miss})}) + \text{const}, \tag{B.4}$$

where *const* is the abbreviation of constant, and $\psi(\mathbf{X}^{(\text{miss})})$ is a scalar function determined by the type of SDE underlying the DMs.

- **VP-SDE:** $\psi(\mathbf{X}^{(\text{miss})}) = \mathbb{E}_{r(\mathbf{X}^{(\text{miss})})} \left\{ \frac{1}{4} [\mathbf{X}^{(\text{miss})}]^\top [\mathbf{X}^{(\text{miss})}] - \frac{1}{2} \log r(\mathbf{X}^{(\text{miss})}) \right\}$
- **VE-SDE:** $\psi(\mathbf{X}^{(\text{miss})}) = \mathbb{E}_{r(\mathbf{X}^{(\text{miss})})} \left\{ -\frac{1}{2} \log r(\mathbf{X}^{(\text{miss})}) \right\}$
- **sub-VP-SDE:** $\psi(\mathbf{X}^{(\text{miss})}) = \mathbb{E}_{r(\mathbf{X}^{(\text{miss})})} \left\{ \frac{1}{4\gamma_\tau} [\mathbf{X}^{(\text{miss})}]^\top [\mathbf{X}^{(\text{miss})}] - \frac{1}{2} \log r(\mathbf{X}^{(\text{miss})}) \right\}$, where γ_τ is determined by noise scale β_τ : $\gamma_\tau := (1 - \exp(-2 \int_0^\tau \beta_s ds)) > 0, 0 < \beta_1 < \beta_2 < \dots < \beta_T < 1$.

It is important to note that in DMs, the condition $\psi(\mathbf{X}^{(\text{miss})}) \geq 0$ consistently holds. This assertion is supported by the fact that the inner product $[\mathbf{X}^{(\text{miss})}]^\top [\mathbf{X}^{(\text{miss})}] \geq 0$, and the entropy function defined as $\mathbb{H}[r(\mathbf{X}^{(\text{miss})})] := - \int r(\mathbf{X}^{(\text{miss})}) \log r(\mathbf{X}^{(\text{miss})}) d\mathbf{X}^{(\text{miss})}$ is also non-negative.

Proof. Since there are various approaches for reversing the sampling procedure of DMs, for simplicity, we mainly consider the VP-SDE, VE-SDE, and sub-VP-SDE as analysis objects.

- **VP-SDE:** According to reference [36], the density evolution of the generative process for VP-SDE can be delineated by the following PDE:

$$\begin{aligned} \frac{\partial r(\mathbf{X}^{(\text{miss})})}{\partial \tau} = & -\nabla_{\mathbf{X}^{(\text{miss})}} \cdot \left\{ r(\mathbf{X}^{(\text{miss})}) [\beta_\tau] \left[\frac{1}{2} \mathbf{X}^{(\text{miss})} + \nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) \right] \right\} \\ & + \frac{\beta_\tau}{2} \nabla_{\mathbf{X}^{(\text{miss})}} \cdot \nabla_{\mathbf{X}^{(\text{miss})}} r(\mathbf{X}^{(\text{miss})}) \end{aligned} \quad (\text{B.5})$$

where $\beta_\tau \in (0, 1)$ is the noise scale. On this basis, by changing the variable as $d\tau := \frac{\beta_\tau}{2} d\tau$ [16], we can get the following equation:

$$\frac{\partial r(\mathbf{X}^{(\text{miss})})}{\partial \tau} = -\nabla_{\mathbf{X}^{(\text{miss})}} \cdot \left\{ \begin{aligned} & r(\mathbf{X}^{(\text{miss})}) \left[\frac{1}{2} \mathbf{X}^{(\text{miss})} + \nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) \right] \\ & - \frac{1}{2} \nabla_{\mathbf{X}^{(\text{miss})}} \log r(\mathbf{X}^{(\text{miss})}) \end{aligned} \right\}. \quad (\text{B.6})$$

Comparing Eq. (B.6) with Eqs. (A.5) and (A.6), the cost functional to be minimized of this simulation procedure can be given as follows:

$$\begin{aligned} \mathcal{F}_{\text{VP-SDE}} = & -\int r(\mathbf{X}^{(\text{miss})}) \left\{ \begin{aligned} & \frac{1}{4} [\mathbf{X}^{(\text{miss})}]^\top [\mathbf{X}^{(\text{miss})}] + \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) \\ & - \frac{1}{2} \log r(\mathbf{X}^{(\text{miss})}) + \text{const} \end{aligned} \right\} d\mathbf{X}^{(\text{miss})} \\ = & -\mathbb{E}_{r(\mathbf{X}^{(\text{miss})})} \left\{ \begin{aligned} & \frac{1}{4} [\mathbf{X}^{(\text{miss})}]^\top [\mathbf{X}^{(\text{miss})}] + \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) \\ & - \frac{1}{2} \log r(\mathbf{X}^{(\text{miss})}) + \text{const} \end{aligned} \right\}. \end{aligned} \quad (\text{B.7})$$

Note that $\frac{1}{4} [\mathbf{X}^{(\text{miss})}]^\top [\mathbf{X}^{(\text{miss})}] \geq 0$ and $-\frac{1}{2} \int r(\mathbf{X}^{(\text{miss})}) \log r(\mathbf{X}^{(\text{miss})}) d\mathbf{X}^{(\text{miss})} \geq 0$ hold, and thus the proposition for VP-SDE is proved by taking the negative of the abovementioned equation.

- **VE-SDE:** Similarly, based on reference [36], the following PDE can be given to delineate the density evolution of the generative process for VE-SDE:

$$\begin{aligned} \frac{\partial r(\mathbf{X}^{(\text{miss})})}{\partial \tau} = & -\nabla_{\mathbf{X}^{(\text{miss})}} \cdot \left\{ r(\mathbf{X}^{(\text{miss})}) \left[-\frac{d\sigma_\tau^2}{d\tau} \right] \nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) \right\} \\ & + \frac{1}{2} \frac{d\sigma_\tau^2}{d\tau} \nabla_{\mathbf{X}^{(\text{miss})}} \cdot \nabla_{\mathbf{X}^{(\text{miss})}} r(\mathbf{X}^{(\text{miss})}), \end{aligned} \quad (\text{B.8})$$

where σ_τ^2 is a time varying noise scale.

As such, by changing the variable as $d\tau := \left[\frac{d\sigma_\tau^2}{d\tau} \right] d\tau$ [16], Eq. (B.8) can be reformulated as follows:

$$\frac{\partial r(\mathbf{X}^{(\text{miss})})}{\partial \tau} = -\nabla_{\mathbf{X}^{(\text{miss})}} \cdot \left\{ r(\mathbf{X}^{(\text{miss})}) \left[\nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) - \frac{1}{2} \nabla_{\mathbf{X}^{(\text{miss})}} \log r(\mathbf{X}^{(\text{miss})}) \right] \right\}. \quad (\text{B.9})$$

Comparing Eq. (B.9) with Eqs. (A.5) and (A.6), the cost functional to be minimized of this simulation procedure can be given as follows:

$$\begin{aligned} \mathcal{F}_{\text{VE-SDE}} = & \int r(\mathbf{X}^{(\text{miss})}) \left\{ \frac{1}{2} \log r(\mathbf{X}^{(\text{miss})}) - \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) + \text{const} \right\} d\mathbf{X}^{(\text{miss})} \\ = & -\mathbb{E}_{r(\mathbf{X}^{(\text{miss})})} \left\{ -\frac{1}{2} \log r(\mathbf{X}^{(\text{miss})}) + \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) + \text{const} \right\}. \end{aligned} \quad (\text{B.10})$$

Note that the entropy function $-\frac{1}{2} \int r(\mathbf{X}^{(\text{miss})}) \log r(\mathbf{X}^{(\text{miss})}) d\mathbf{X}^{(\text{miss})} \geq 0$ holds, and thus the proposition for VE-SDE is proved by taking the negative of the abovementioned equation.

- **sub-VP-SDE:** Based on reference [36], the following PDE can be given to delineate the density evolution of the generative process for sub-VP-SDE:

$$\begin{aligned} \frac{\partial r(\mathbf{X}^{(\text{miss})})}{\partial \tau} = & -\nabla_{\mathbf{X}^{(\text{miss})}} \cdot \left\{ r(\mathbf{X}^{(\text{miss})}) [\beta_\tau] \left[\frac{1}{2} \mathbf{X}^{(\text{miss})} + \gamma_\tau \nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) \right] \right\} \\ & + \frac{\beta_\tau}{2} \gamma_\tau \nabla_{\mathbf{X}^{(\text{miss})}} \cdot \nabla_{\mathbf{X}^{(\text{miss})}} r(\mathbf{X}^{(\text{miss})}), \end{aligned} \quad (\text{B.11})$$

where $\gamma_\tau := (1 - \exp(-2 \int_0^\tau \beta_s ds)) > 0$. On this basis, by changing the variable as $d\tau := \frac{\beta_\tau}{2} d\tau$, we can get the following equation:

$$\frac{\partial r(\mathbf{X}^{(miss)})}{\partial \tau} = -\nabla_{\mathbf{X}^{(miss)}} \cdot \left\{ \begin{array}{l} r(\mathbf{X}^{(miss)}) \left[\frac{1}{2} \mathbf{X}^{(miss)} + \gamma_\tau \nabla_{\mathbf{X}^{(miss)}} \log \hat{p}(\mathbf{X}^{(miss)} | \mathbf{X}^{(obs)}) \right] \\ - \frac{\gamma_\tau}{2} \nabla_{\mathbf{X}^{(miss)}} \log r(\mathbf{X}^{(miss)}) \end{array} \right\}. \quad (\text{B.12})$$

Comparing Eq. (B.12) with Eqs. (A.5) and (A.6), the cost functional to be minimized of this simulation procedure can be given as follows:

$$\begin{aligned} \mathcal{F}_{\text{sub-VP-SDE}} &= - \int r(\mathbf{X}^{(miss)}) \left\{ \begin{array}{l} \frac{1}{4} [\mathbf{X}^{(miss)}]^\top [\mathbf{X}^{(miss)}] + \gamma_\tau \log \hat{p}(\mathbf{X}^{(miss)} | \mathbf{X}^{(obs)}) \\ - \frac{\gamma_\tau}{2} \log r(\mathbf{X}^{(miss)}) + \text{const} \end{array} \right\} d\mathbf{X}^{(miss)} \\ &= -\mathbb{E}_{r(\mathbf{X}^{(miss)})} \left\{ \begin{array}{l} \frac{1}{4} [\mathbf{X}^{(miss)}]^\top [\mathbf{X}^{(miss)}] + \gamma_\tau \log \hat{p}(\mathbf{X}^{(miss)} | \mathbf{X}^{(obs)}) \\ - \frac{\gamma_\tau}{2} \log r(\mathbf{X}^{(miss)}) + \text{const} \end{array} \right\} \\ &= -\mathbb{E}_{r(\mathbf{X}^{(miss)})} \left\{ \begin{array}{l} \frac{1}{4\gamma_\tau} [\mathbf{X}^{(miss)}]^\top [\mathbf{X}^{(miss)}] + \log \hat{p}(\mathbf{X}^{(miss)} | \mathbf{X}^{(obs)}) \\ - \frac{1}{2} \log r(\mathbf{X}^{(miss)}) + \text{const} \end{array} \right\}. \end{aligned} \quad (\text{B.13})$$

Note that $\frac{1}{4\gamma_\tau} [\mathbf{X}^{(miss)}]^\top [\mathbf{X}^{(miss)}] \geq 0$ and $-\frac{1}{2} \int r(\mathbf{X}^{(miss)}) \log r(\mathbf{X}^{(miss)}) d\mathbf{X}^{(miss)} \geq 0$ hold, and thus the proposition for sub-VP-SDE is proved by taking the negative of the abovementioned equation.

In summary, the regularization term $\psi(\mathbf{X}^{(miss)})$ for VP-SDE is $\mathbb{E}_{\mathbf{X}^{(miss)} \sim r(\mathbf{X}^{(miss)})} \left\{ \frac{1}{4} [\mathbf{X}^{(miss)}]^\top [\mathbf{X}^{(miss)}] - \frac{1}{2} \log r(\mathbf{X}^{(miss)}) \right\}$, for VE-SDE is $\frac{1}{2} \mathbb{H}(r(\mathbf{X}^{(miss)}))$, and for sub-VP-SDE is $\mathbb{E}_{\mathbf{X}^{(miss)} \sim r(\mathbf{X}^{(miss)})} \left\{ \frac{1}{4\gamma_\tau} [\mathbf{X}^{(miss)}]^\top [\mathbf{X}^{(miss)}] - \frac{1}{2} \log r(\mathbf{X}^{(miss)}) \right\}$. \square

Proposition (3.2). *The evolution of \mathcal{F}_{NER} along τ can be characterized by the following ODE, assuming that the boundary condition $\mathbb{E}_{r(\mathbf{X}^{(miss)}, \tau)} \{ \nabla_{\mathbf{X}^{(miss)}} \cdot [u(\mathbf{X}^{(miss)}, \tau) \log \hat{p}(\mathbf{X}^{(miss)} | \mathbf{X}^{(obs)})] \} = 0$ is satisfied for the velocity field $u(\mathbf{X}^{(miss)}, \tau)$:*

$$\frac{d\mathcal{F}_{NER}}{d\tau} = \mathbb{E}_{r(\mathbf{X}^{(miss)}, \tau)} [u^\top(\mathbf{X}^{(miss)}, \tau) \nabla_{\mathbf{X}^{(miss)}} \log \hat{p}(\mathbf{X}^{(miss)} | \mathbf{X}^{(obs)}) - \lambda \nabla_{\mathbf{X}^{(miss)}} \cdot u(\mathbf{X}^{(miss)}, \tau)]. \quad (\text{B.14})$$

This boundary condition is achievable, for instance, when $\hat{p}(\mathbf{X}^{(miss)} | \mathbf{X}^{(obs)})$ is bounded, and the limit of the velocity field as the norm of $\mathbf{X}^{(miss)}$ approaches zero is zero ($\lim_{\|\mathbf{X}^{(miss)}\| \rightarrow 0} u(\mathbf{X}^{(miss)}, \tau) = 0$).

Proof. Before proving this proposition, we should recognize that the evolution of $\mathbf{X}^{(miss)}$ should promise the probability density function $r(\mathbf{X}^{(miss)}, \tau)$ unchanged. In other words, the following continuity equation should be satisfied during the optimization of $r(\mathbf{X}^{(miss)}, \tau)$:

$$\frac{\partial r(\mathbf{X}^{(miss)}, \tau)}{\partial \tau} = -\nabla_{\mathbf{X}^{(miss)}} \cdot [r(\mathbf{X}^{(miss)}, \tau) u(\mathbf{X}^{(miss)}, \tau)]. \quad (\text{B.15})$$

On this basis, the evolution of \mathcal{F}_{NER} along time τ , $\frac{d\mathcal{F}_{\text{NER}}}{d\tau}$, can be given as follows based on the chain rule:

$$\begin{aligned}
& \frac{d\mathcal{F}_{\text{NER}}}{d\tau} \\
&= \int \frac{\partial r(\mathbf{X}^{(\text{miss})}, \tau)}{\partial \tau} \left[\log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) + \lambda \log r(\mathbf{X}^{(\text{miss})}, \tau) + \lambda \right] d\mathbf{X}^{(\text{miss})} \\
&= \int -\{\nabla_{\mathbf{X}^{(\text{miss})}} \cdot [r(\mathbf{X}^{(\text{miss})}, \tau) u(\mathbf{X}^{(\text{miss})}, \tau)]\} [\log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) + \lambda \log r(\mathbf{X}^{(\text{miss})}, \tau) + \lambda] d\mathbf{X}^{(\text{miss})} \\
&\stackrel{(i)}{=} \int [r(\mathbf{X}^{(\text{miss})}, \tau) u(\mathbf{X}^{(\text{miss})}, \tau)]^\top \nabla_{\mathbf{X}^{(\text{miss})}} [\log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) + \lambda \log r(\mathbf{X}^{(\text{miss})}, \tau) + \lambda] d\mathbf{X}^{(\text{miss})} \\
&= \int [r(\mathbf{X}^{(\text{miss})}, \tau) u(\mathbf{X}^{(\text{miss})}, \tau)]^\top \{\nabla_{\mathbf{X}^{(\text{miss})}} [\log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) + \lambda \log r(\mathbf{X}^{(\text{miss})}, \tau)]\} d\mathbf{X}^{(\text{miss})} \\
&= \int [u(\mathbf{X}^{(\text{miss})}, \tau)]^\top [r(\mathbf{X}^{(\text{miss})}, \tau) \nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) + \lambda r(\mathbf{X}^{(\text{miss})}, \tau) \nabla_{\mathbf{X}^{(\text{miss})}} \log r(\mathbf{X}^{(\text{miss})}, \tau)] d\mathbf{X}^{(\text{miss})} \\
&= \int [u(\mathbf{X}^{(\text{miss})}, \tau)]^\top [r(\mathbf{X}^{(\text{miss})}, \tau) \nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) + \lambda \nabla_{\mathbf{X}^{(\text{miss})}} r(\mathbf{X}^{(\text{miss})}, \tau)] d\mathbf{X}^{(\text{miss})} \\
&\stackrel{(ii)}{=} \int r(\mathbf{X}^{(\text{miss})}, \tau) [u^\top(\mathbf{X}^{(\text{miss})}, \tau) \nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) - \lambda \nabla_{\mathbf{X}^{(\text{miss})}} \cdot u(\mathbf{X}^{(\text{miss})}, \tau)] d\mathbf{X}^{(\text{miss})} \\
&= \mathbb{E}_{r(\mathbf{X}^{(\text{miss})}, \tau)} [u^\top(\mathbf{X}^{(\text{miss})}, \tau) \nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) - \lambda \nabla_{\mathbf{X}^{(\text{miss})}} \cdot u(\mathbf{X}^{(\text{miss})}, \tau)], \tag{B.16}
\end{aligned}$$

where (i) and (ii) are based on integration by parts. \square

Proposition (3.3). *When the velocity field $u(\mathbf{X}^{(\text{miss})}, \tau)$ is constrained by the norm of RKHS, the problem of finding the steepest gradient ascent direction can be formulated as follows:*

$$u(\mathbf{X}^{(\text{miss})}, \tau) = \arg \max_{v(\mathbf{X}^{(\text{miss})}, \tau) \in \mathcal{H}^d} \left\{ \mathbb{E}_{r(\mathbf{X}^{(\text{miss})}, \tau)} [v^\top(\mathbf{X}^{(\text{miss})}, \tau) \nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) - \lambda \nabla_{\mathbf{X}^{(\text{miss})}} \cdot v(\mathbf{X}^{(\text{miss})}, \tau)] \right\} - \frac{1}{2} \|v(\mathbf{X}^{(\text{miss})}, \tau)\|_{\mathcal{H}}^2. \tag{B.17}$$

The corresponding optimal solution is given by:

$$u(\mathbf{X}^{(\text{miss})}, \tau) = \mathbb{E}_{r(\tilde{\mathbf{X}}^{(\text{miss})}, \tau)} \left\{ \begin{array}{l} -\lambda \nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \mathcal{K}(\mathbf{X}^{(\text{miss})}, \tilde{\mathbf{X}}^{(\text{miss})}) \\ + [\nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \log \hat{p}(\tilde{\mathbf{X}}^{(\text{miss})} | \mathbf{X}^{(\text{obs})})]^\top \mathcal{K}(\mathbf{X}^{(\text{miss})}, \tilde{\mathbf{X}}^{(\text{miss})}) \end{array} \right\}, \tag{B.18}$$

where $\mathcal{K}(\mathbf{X}^{(\text{miss})}, \tilde{\mathbf{X}}^{(\text{miss})})$ is kernel function.

Proof. Assume we have a map function $\phi(x)$, the kernel function can be given as follows:

$$\mathcal{K}(x, y) = \langle \phi(x), \phi(y) \rangle_{\mathcal{H}}. \tag{B.19}$$

Based on this, the regularization term that control the magnitude of $v(\mathbf{X}^{(\text{miss})}, \tau)$ can be given by $\frac{1}{2} \|v(\mathbf{X}^{(\text{miss})}, \tau)\|_{\mathcal{H}}$, and the spectral decomposition of kernel function can be given as follows:

$$\mathcal{K}(x, y) = \sum_{i=1}^{\infty} \xi_i \phi_i(x) \phi_i(y), \tag{B.20}$$

where $\phi_i(\cdot)$ indicates the orthonormal basis and ξ_i is the corresponding eigen-value. For any function $v(\mathbf{X}^{(\text{miss})}, \tau) \in \mathcal{H}$, the following decomposition is given:

$$v(\mathbf{X}^{(\text{miss})}, \tau) = \sum_{i=1}^{\infty} v_i \sqrt{\xi_i} \phi_i(\mathbf{X}^{(\text{miss})}, \tau), \tag{B.21}$$

where v_i and $\sum_{i=1}^{\infty} \|v_i\|^2 < \infty$.

The learning objective defined in Eq. (B.17) can be reformulated as follows:

$$\begin{aligned}
& v^*(\mathbf{X}^{(\text{miss})}, \tau) \\
&= \arg \max_{v(\mathbf{X}^{(\text{miss})}, \tau) \in \mathcal{H}^d} \left\{ \mathbb{E}_{r(\mathbf{X}^{(\text{miss})}, \tau)} [v^\top (\mathbf{X}^{(\text{miss})}, \tau) \nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) \right. \\
&\quad \left. - \lambda \nabla_{\mathbf{X}^{(\text{miss})}} \cdot v(\mathbf{X}^{(\text{miss})}, \tau)] \right\} - \frac{1}{2} \|v(\mathbf{X}^{(\text{miss})}, \tau)\|_{\mathcal{H}^d}^2, \\
&\stackrel{(i)}{=} \arg \max_{v(\mathbf{X}^{(\text{miss})}, \tau) \in \mathcal{H}^d} \left\{ \mathbb{E}_{r(\tilde{\mathbf{X}}^{(\text{miss})}, \tau)} \left[\sum_{i=1}^{\infty} \sqrt{\xi_i} \nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \log \hat{p}(\tilde{\mathbf{X}}^{(\text{miss})} | \mathbf{X}^{(\text{obs})})^\top v_i \phi_i(\tilde{\mathbf{X}}^{(\text{miss})}, \tau) \right. \right. \\
&\quad \left. \left. - \lambda \nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \cdot \sum_{i=1}^{\infty} v_i \sqrt{\xi_i} \phi_i(\tilde{\mathbf{X}}^{(\text{miss})}, \tau) \right] \right\} - \frac{1}{2} \sum_{i=1}^{\infty} \|v_i\|^2,
\end{aligned} \tag{B.22}$$

Take the right-hand-side of (i) with-respect-to v_i , and set it to 0, we can get:

$$\sqrt{\xi_i} \left\{ \mathbb{E}_{r(\tilde{\mathbf{X}}^{(\text{miss})}, \tau)} \left[[\nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \log \hat{p}(\tilde{\mathbf{X}}^{(\text{miss})} | \mathbf{X}^{(\text{obs})})]^\top \phi_i(\tilde{\mathbf{X}}^{(\text{miss})}, \tau) - \lambda \nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \phi_i(\tilde{\mathbf{X}}^{(\text{miss})}, \tau) \right] \right\} - v_i = 0. \tag{B.23}$$

On this basis, v_i^* can be given as follows:

$$v_i^* = \sqrt{\xi_i} \left\{ \mathbb{E}_{r(\tilde{\mathbf{X}}^{(\text{miss})}, \tau)} \left[[\nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \log \hat{p}(\tilde{\mathbf{X}}^{(\text{miss})} | \mathbf{X}^{(\text{obs})})]^\top \phi_i(\tilde{\mathbf{X}}^{(\text{miss})}, \tau) - \lambda \nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \phi_i(\tilde{\mathbf{X}}^{(\text{miss})}, \tau) \right] \right\}, \tag{B.24}$$

and hence, $u(\mathbf{X}^{(\text{miss})}, \tau)$ can be given as follows:

$$\begin{aligned}
& u(\mathbf{X}^{(\text{miss})}, \tau) \\
&= \sum_{i=1}^{\infty} \sqrt{\xi_i} v_i^* \phi_i(\mathbf{X}^{(\text{miss})}, \tau) \\
&= \mathbb{E}_{r(\tilde{\mathbf{X}}^{(\text{miss})}, \tau)} \left[\begin{array}{c} -\lambda \nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \mathcal{K}(\mathbf{X}^{(\text{miss})}, \tilde{\mathbf{X}}^{(\text{miss})}) \\ + [\nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \log \hat{p}(\tilde{\mathbf{X}}^{(\text{miss})} | \mathbf{X}^{(\text{obs})})]^\top \mathcal{K}(\mathbf{X}^{(\text{miss})}, \tilde{\mathbf{X}}^{(\text{miss})}) \end{array} \right].
\end{aligned} \tag{B.25}$$

□

Proposition (3.4). Suppose the proposal distribution $r(\mathbf{X}^{(\text{joint})})$ is factorized by $r(\mathbf{X}^{(\text{joint})}) := r(\mathbf{X}^{(\text{miss})})p(\mathbf{X}^{(\text{obs})})$. The cost functional concerned with joint distribution defined by the following equation:

$$\mathcal{F}_{\text{joint-NER}} := \mathbb{E}_{r(\mathbf{X}^{(\text{joint})})} [\log \hat{p}(\mathbf{X}^{(\text{joint})})] - \lambda \mathbb{H}[r(\mathbf{X}^{(\text{joint})})], \tag{B.26}$$

results in the velocity field defined in Eq. (10), and is a lower bound of \mathcal{F}_{NER} where the gap is a constant. (i.e. $\mathcal{F}_{\text{joint-NER}} = \mathcal{F}_{\text{NER}} - \text{const}, \text{const} \geq 0$.)

Proof. Our proof will be divided into two parts namely ‘velocity field derivation’ and ‘upper bound acquirement’.

Velocity Field Derivation:

the following continuity equation should be satisfied during the optimization of $r(\mathbf{X}^{(\text{miss})}, \tau)$:

$$\begin{aligned}
& \frac{\partial r(\mathbf{X}^{(\text{miss})}, \tau)}{\partial \tau} = -\nabla_{\mathbf{X}^{(\text{miss})}} \cdot [r(\mathbf{X}^{(\text{miss})}, \tau) u(\mathbf{X}^{(\text{miss})}, \tau)] \\
& \Rightarrow \frac{\partial r(\mathbf{X}^{(\text{miss})}, \tau)}{\partial \tau} \times p(\mathbf{X}^{(\text{obs})}) = -\nabla_{\mathbf{X}^{(\text{miss})}} \cdot [r(\mathbf{X}^{(\text{miss})}, \tau) u(\mathbf{X}^{(\text{miss})}, \tau)] \times p(\mathbf{X}^{(\text{obs})}) \\
& \stackrel{(i)}{\Rightarrow} \frac{\partial r(\mathbf{X}^{(\text{joint})}, \tau)}{\partial \tau} = -\nabla_{\mathbf{X}^{(\text{miss})}} \cdot [r(\mathbf{X}^{(\text{joint})}, \tau) u(\mathbf{X}^{(\text{joint})}, \tau)],
\end{aligned} \tag{B.27}$$

where (i) is based on the fact that $\mathbf{X}^{(\text{obs})}$ remains unchanged during imputation process. And thus, according to Eq. (B.16), the evolution of $\mathcal{F}_{\text{joint-NER}}$ along time τ , $\frac{d\mathcal{F}_{\text{joint-NER}}}{d\tau}$, can be given as follows

based on the chain rule:

$$\begin{aligned}
& \frac{d\mathcal{F}_{\text{joint-NER}}}{d\tau} \\
&= \int \frac{\partial r(\mathbf{X}^{(\text{joint})}, \tau)}{\partial \tau} \left[\log \hat{p}(\mathbf{X}^{(\text{joint})}) + \lambda \log r(\mathbf{X}^{(\text{joint})}, \tau) + \lambda \right] d\mathbf{X}^{(\text{joint})} \\
&= \int -\{\nabla_{\mathbf{X}^{(\text{miss})}} \cdot [r(\mathbf{X}^{(\text{joint})}, \tau)u(\mathbf{X}^{(\text{joint})}, \tau)]\} [\log \hat{p}(\mathbf{X}^{(\text{joint})}) + \lambda \log r(\mathbf{X}^{(\text{joint})}, \tau) + \lambda] d\mathbf{X}^{(\text{joint})} \\
&\stackrel{(i)}{=} \int [r(\mathbf{X}^{(\text{joint})}, \tau)u(\mathbf{X}^{(\text{joint})}, \tau)]^\top \nabla_{\mathbf{X}^{(\text{miss})}} [\log \hat{p}(\mathbf{X}^{(\text{joint})}) + \lambda \log r(\mathbf{X}^{(\text{joint})}, \tau) + \lambda] d\mathbf{X}^{(\text{joint})} \\
&= \int [r(\mathbf{X}^{(\text{joint})}, \tau)u(\mathbf{X}^{(\text{joint})}, \tau)]^\top \{\nabla_{\mathbf{X}^{(\text{miss})}} [\log \hat{p}(\mathbf{X}^{(\text{joint})}) + \lambda \log r(\mathbf{X}^{(\text{joint})}, \tau)]\} d\mathbf{X}^{(\text{joint})} \\
&= \int [u(\mathbf{X}^{(\text{joint})}, \tau)]^\top [r(\mathbf{X}^{(\text{joint})}, \tau) \nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{joint})}) + \lambda r(\mathbf{X}^{(\text{joint})}, \tau) \nabla_{\mathbf{X}^{(\text{miss})}} \log r(\mathbf{X}^{(\text{joint})}, \tau)] d\mathbf{X}^{(\text{joint})} \\
&= \int [u(\mathbf{X}^{(\text{joint})}, \tau)]^\top [r(\mathbf{X}^{(\text{joint})}, \tau) \nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{joint})}) + \lambda \nabla_{\mathbf{X}^{(\text{miss})}} r(\mathbf{X}^{(\text{joint})}, \tau)] d\mathbf{X}^{(\text{joint})} \\
&\stackrel{(ii)}{=} \int r(\mathbf{X}^{(\text{joint})}, \tau) [u^\top(\mathbf{X}^{(\text{joint})}, \tau) \nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{joint})}) - \lambda \nabla_{\mathbf{X}^{(\text{miss})}} \cdot u(\mathbf{X}^{(\text{joint})}, \tau)] d\mathbf{X}^{(\text{joint})} \\
&= \mathbb{E}_{r(\mathbf{X}^{(\text{joint})}, \tau)} [u^\top(\mathbf{X}^{(\text{joint})}, \tau) \nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{joint})}) - \lambda \nabla_{\mathbf{X}^{(\text{miss})}} \cdot u(\mathbf{X}^{(\text{joint})}, \tau)], \tag{B.28}
\end{aligned}$$

where (i) and (ii) are based on integration by parts.

Similar to the proof of proposition 3.3, we can restrict the velocity field in RKHS and find the steepest gradient boosting direction as follows according to Eqs. (B.19) to (B.21):

$$\begin{aligned}
& v^*(\mathbf{X}^{(\text{joint})}, \tau) \\
&= \arg \max_{v(\mathbf{X}^{(\text{joint})}, \tau) \in \mathcal{H}^d} \{ \mathbb{E}_{r(\mathbf{X}^{(\text{joint})}, \tau)} [v^\top(\mathbf{X}^{(\text{joint})}, \tau) \nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{joint})}) \\
&\quad - \lambda \nabla_{\mathbf{X}^{(\text{miss})}} \cdot v(\mathbf{X}^{(\text{miss})}, \tau)] \} - \frac{1}{2} \|v(\mathbf{X}^{(\text{joint})}, \tau)\|_{\mathcal{H}^d}, \tag{B.29} \\
&\stackrel{(i)}{=} \arg \max_{v(\mathbf{X}^{(\text{joint})}, \tau) \in \mathcal{H}^d} \{ \mathbb{E}_{r(\tilde{\mathbf{X}}^{(\text{joint})}, \tau)} \left[\sum_{i=1}^{\infty} \sqrt{\xi_i} \nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \log \hat{p}(\tilde{\mathbf{X}}^{(\text{joint})})^\top v_i \phi_i(\tilde{\mathbf{X}}^{(\text{joint})}, \tau) \right. \\
&\quad \left. - \lambda \nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \cdot \sum_{i=1}^{\infty} v_i \sqrt{\xi_i} \phi_i(\tilde{\mathbf{X}}^{(\text{joint})}, \tau) \right] \} - \frac{1}{2} \sum_{i=1}^{\infty} \|v_i\|^2,
\end{aligned}$$

Take the right-hand-side of (i) with-respect-to v_i , and set it to 0, we can get:

$$\sqrt{\xi_i} \{ \mathbb{E}_{r(\tilde{\mathbf{X}}^{(\text{joint})}, \tau)} [[\nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \log \hat{p}(\tilde{\mathbf{X}}^{(\text{joint})})]^\top \phi_i(\tilde{\mathbf{X}}^{(\text{joint})}, \tau) - \lambda \nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \phi_i(\tilde{\mathbf{X}}^{(\text{joint})}, \tau)] \} - v_i = 0. \tag{B.30}$$

On this basis, v_i^* can be given as follows:

$$v_i^* = \sqrt{\xi_i} \{ \mathbb{E}_{r(\tilde{\mathbf{X}}^{(\text{joint})}, \tau)} [[\nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \log \hat{p}(\tilde{\mathbf{X}}^{(\text{joint})})]^\top \phi_i(\tilde{\mathbf{X}}^{(\text{joint})}, \tau) - \lambda \nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \phi_i(\tilde{\mathbf{X}}^{(\text{joint})}, \tau)] \}, \tag{B.31}$$

and hence, $u(\mathbf{X}^{(\text{joint})}, \tau)$ can be given as follows:

$$\begin{aligned}
& u(\mathbf{X}^{(\text{joint})}, \tau) \\
&= \sum_{i=1}^{\infty} \sqrt{\xi_i} v_i^* \phi_i(\mathbf{X}^{(\text{joint})}, \tau) \\
&= \mathbb{E}_{r(\mathbf{X}^{(\text{joint})}, \tau)} \left[\begin{aligned} & -\lambda \nabla_{\tilde{\mathbf{X}}^{(\text{joint})}} \mathcal{K}(\mathbf{X}^{(\text{joint})}, \tilde{\mathbf{X}}^{(\text{joint})}) \\ & + \nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \log \hat{p}(\tilde{\mathbf{X}}^{(\text{joint})}) \mathcal{K}(\mathbf{X}^{(\text{miss})}, \tilde{\mathbf{X}}^{(\text{miss})}) \end{aligned} \right]. \tag{B.32}
\end{aligned}$$

Lower Bound Acquirement:

Before proving the proposition, we should notice that given the unchanged observational data $\mathbf{X}^{(\text{obs})}$, the distribution $p(\mathbf{X}^{(\text{obs})})$ is a constant. On this basis, consider the definition of \mathcal{F}_{NER} (right-hand-side of Eq. (5)), the first term and the second term are denoted by ‘term 1’ and ‘term 2’ for simplicity:

$$\underbrace{\mathbb{E}_{r(\mathbf{X}^{(\text{miss})})}[\log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})})]}_{\text{term 1}} + \lambda \times \left[\underbrace{-\mathbb{H}[r(\mathbf{X}^{(\text{miss})})]}_{\text{term 2}} \right]. \quad (\text{B.33})$$

For term 1, we can obtain the following derivation:

$$\begin{aligned} & \int r(\mathbf{X}^{(\text{miss})}) \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) d\mathbf{X}^{(\text{miss})} \\ & \geq \int r(\mathbf{X}^{(\text{miss})}) \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) d\mathbf{X}^{(\text{miss})} + \underbrace{\int p(\mathbf{X}^{(\text{obs})}) \log p(\mathbf{X}^{(\text{obs})}) d\mathbf{X}^{(\text{obs})}}_{\text{negative entropy (negative constant)}} \\ & = \iint p(\mathbf{X}^{(\text{obs})}) r(\mathbf{X}^{(\text{miss})}) \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) d\mathbf{X}^{(\text{miss})} d\mathbf{X}^{(\text{obs})} \\ & \quad + \underbrace{\int p(\mathbf{X}^{(\text{obs})}) \log p(\mathbf{X}^{(\text{obs})}) d\mathbf{X}^{(\text{obs})}}_{\text{negative constant}} \\ & = \iint p(\mathbf{X}^{(\text{obs})}) r(\mathbf{X}^{(\text{miss})}) \log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) d\mathbf{X}^{(\text{miss})} d\mathbf{X}^{(\text{obs})} \\ & \quad + \underbrace{\iint r(\mathbf{X}^{(\text{miss})}) p(\mathbf{X}^{(\text{obs})}) \log p(\mathbf{X}^{(\text{obs})}) d\mathbf{X}^{(\text{miss})} d\mathbf{X}^{(\text{obs})}}_{\text{negative constant}} \\ & = \iint \underbrace{p(\mathbf{X}^{(\text{obs})}) r(\mathbf{X}^{(\text{miss})})}_{r(\mathbf{X}^{(\text{miss})}, \mathbf{X}^{(\text{obs})})} \underbrace{[\log \hat{p}(\mathbf{X}^{(\text{miss})} | \mathbf{X}^{(\text{obs})}) + \log p(\mathbf{X}^{(\text{obs})})]}_{\log \hat{p}(\mathbf{X}^{(\text{miss})}, \mathbf{X}^{(\text{obs})})} d\mathbf{X}^{(\text{miss})} d\mathbf{X}^{(\text{obs})} \\ & = \mathbb{E}_{r(\mathbf{X}^{(\text{miss})}, \mathbf{X}^{(\text{obs})})}[\log \hat{p}(\mathbf{X}^{(\text{miss})}, \mathbf{X}^{(\text{obs})})]. \end{aligned} \quad (\text{B.34})$$

Similarly, the term 2 can be reformulated as follows:

$$\begin{aligned} & -\mathbb{H}[r(\mathbf{X}^{(\text{miss})})] \\ & \geq -\mathbb{H}[r(\mathbf{X}^{(\text{miss})})] + \underbrace{\int p(\mathbf{X}^{(\text{obs})}) \log p(\mathbf{X}^{(\text{obs})}) d\mathbf{X}^{(\text{obs})}}_{\text{negative entropy (negative constant)}} \\ & = \iint p(\mathbf{X}^{(\text{obs})}) r(\mathbf{X}^{(\text{miss})}) \log r(\mathbf{X}^{(\text{miss})}) d\mathbf{X}^{(\text{miss})} d\mathbf{X}^{(\text{obs})} \\ & \quad + \underbrace{\iint p(\mathbf{X}^{(\text{obs})}) r(\mathbf{X}^{(\text{miss})}) \log p(\mathbf{X}^{(\text{obs})}) d\mathbf{X}^{(\text{miss})} d\mathbf{X}^{(\text{obs})}}_{\text{negative entropy (negative constant)}} \\ & = \iint \underbrace{p(\mathbf{X}^{(\text{obs})}) r(\mathbf{X}^{(\text{miss})})}_{r(\mathbf{X}^{(\text{obs})}, \mathbf{X}^{(\text{miss})})} \underbrace{[\log r(\mathbf{X}^{(\text{miss})}) + \log p(\mathbf{X}^{(\text{obs})})]}_{r(\mathbf{X}^{(\text{obs})}, \mathbf{X}^{(\text{miss})})} d\mathbf{X}^{(\text{miss})} d\mathbf{X}^{(\text{obs})} \\ & = -\mathbb{H}[r(\mathbf{X}^{(\text{obs})}, \mathbf{X}^{(\text{miss})})]. \end{aligned} \quad (\text{B.35})$$

Combine Eqs. (B.34) and (B.35), we can obtain the following relationship:

$$\mathcal{F}_{\text{NER}} - \text{const} = \mathcal{F}_{\text{joint-NER}}, \quad (\text{B.36})$$

and constant const is greater than 0. □

Corollary (3.5). *The following equation can be satisfied:*

$$u(\mathbf{X}^{(\text{joint})}, \tau) = u(\mathbf{X}^{(\text{miss})}, \tau). \quad (\text{B.37})$$

Proof. This corollary can be easily proven by according to Eq. (B.36):

$$\begin{aligned}
\mathcal{F}_{\text{NER}} &= \mathcal{F}_{\text{joint-NER}} + \text{const} \\
\Rightarrow \nabla_{\mathbf{X}^{(\text{miss})}} \frac{\delta \mathcal{F}_{\text{NER}}}{\delta r(\mathbf{X}^{(\text{miss})})} &= \nabla_{\mathbf{X}^{(\text{miss})}} \frac{\delta \mathcal{F}_{\text{joint-NER}} + \text{const}}{\delta r(\mathbf{X}^{(\text{miss})})} \\
\Rightarrow \nabla_{\mathbf{X}^{(\text{miss})}} \frac{\delta \mathcal{F}_{\text{NER}}}{\delta r(\mathbf{X}^{(\text{miss})})} &= \nabla_{\mathbf{X}^{(\text{miss})}} \frac{\delta \mathcal{F}_{\text{joint-NER}}}{\delta r(\mathbf{X}^{(\text{miss})})}.
\end{aligned} \tag{B.38}$$

Plug Eq. (B.38) into Eqs. (A.5) and (A.6), we can see that the density functions for $\mathbf{X}^{(\text{miss})}$ within functional \mathcal{F}_{NER} and $\mathcal{F}_{\text{joint-NER}}$ are identical. \square

Appendix C Detailed Information for KnewImp Implementation

C.1 Forward Euler’s Method for ODE Simulation

Suppose we have the following ODE:

$$\frac{dx_\tau}{d\tau} = f(x_\tau, \tau), \tag{C.1}$$

and the initial value at $\tau = 0$ is given $x_0 = x_{\text{init}}$, the value at time η can be derived as follows:

$$x_\eta = x_0 + \int_0^\eta f(x_\tau, \tau) d\tau. \tag{C.2}$$

To alleviate the intergal term, the forward Euler’s method attempts to convert the integral term to summation term as follows:

$$x_\eta = x_0 + f(x_\tau, \tau) \times (\eta - 0). \tag{C.3}$$

On this basis, the value at time T can be obtained by repeating Eq. (C.3) from $\tau = 0$ to $\tau = T$, which is the forward Euler’s method.

Algorithm 1: Algorithm for Forward Euler’s Method

Input: ODE $f(x_\tau, \tau)$; start point τ_0 ; end point τ_T ; step size η ; initial value x_{τ_0} .

Output: Predicted value x_{τ_T} at τ_T .

- 1 Repeating times j calculation: $j \leftarrow (\tau_T - \tau_0) / \eta$
 - 2 **for** $t = \tau_0 + \eta, \tau_0 + 2\eta, \dots, \tau_0 + j\eta$ **do**
 - 3 | $x_{\tau_T} \leftarrow x(t - \eta) + f(x_{t-\eta}, t - \eta) \times \eta$
 - 4 **end**
-

C.2 Algorithms for KnewImp

As we pointed out in Fig. 1, the KnewImp mainly consists of two parts, namely ‘Impute’ and ‘Estimate’. Based on this, we first give the algorithm for the ‘Impute’ and ‘Estimate’ parts in this subsection.

Algorithm 2: Impute Part Algorithm for KnewImp

Input: Initialized Missing Data $\mathbf{X}^{(\text{imp})}$, Score Function: $\nabla_{\mathbf{X}^{(\text{joint})}} \log \hat{p}(\mathbf{X}^{(\text{joint})})$, and Mask Matrix \mathbf{M} .

Hyperparameters:

Simulation Time: \mathbb{T} , Discretization Step Size η , and Bandwidth of RBF kernel h .

Output: Imputed Result: $\mathbf{X}^{(\text{imp})}$.

```

1 Set  $\tau = 0$ ,
2 while  $\tau < \mathbb{T}$  do
   /* Velocity Field Acquisition */
3    $\nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{joint})}) \leftarrow \nabla_{\mathbf{X}^{(\text{joint})}} \log \hat{p}(\mathbf{X}^{(\text{joint})}) \odot (\mathbb{1}_{N \times D} - \mathbf{M}) + 0 \times \mathbf{M}$ ,
4    $u(\mathbf{X}^{(\text{joint})}, \tau) \leftarrow \mathbb{E}_{r(\mathbf{X}^{(\text{joint})}, \tau)} \left\{ \begin{array}{l} -\lambda \nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \mathcal{K}(\mathbf{X}^{(\text{joint})}, \tilde{\mathbf{X}}^{(\text{joint})}) \\ + [\nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{joint})})]^\top \mathcal{K}(\mathbf{X}^{(\text{joint})}, \tilde{\mathbf{X}}^{(\text{joint})}) \end{array} \right\}$ ,
   /* ODE Simulation By Forward Euler's Method */
5    $\mathbf{X}^{(\text{imp})} \leftarrow \mathbf{X}^{(\text{imp})} + \eta \times u(\mathbf{X}^{(\text{joint})}, \tau)$ ,
6    $\tau \leftarrow \tau + 1$ .
7 end

```

Algorithm 3: Estimate Part Algorithm for KnewImp

Input: Imputed Data $\mathbf{X}^{(\text{imp})}$, and $\nabla_{\mathbf{X}^{(\text{joint})}} \log \hat{p}(\mathbf{X}^{(\text{joint})})$ parameterized by Neural Network with Parameter θ .

Hyperparameters:

Network Learning Rate lr , Training Epoch \mathcal{E} , and Network Hidden Unit HU_{score} .

Output: Score Function: $\nabla_{\mathbf{X}^{(\text{joint})}} \log \hat{p}(\mathbf{X}^{(\text{joint})})$.

```

1 while  $e \leq \mathcal{E}$  do
   /* Data Noising */
2    $\hat{\mathbf{X}}^{(\text{joint})} \leftarrow \mathbf{X}^{(\text{joint})} + \epsilon, \epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ ,
3    $\nabla_{\hat{\mathbf{X}}^{(\text{joint})}} \log q_\sigma(\hat{\mathbf{X}}^{(\text{joint})} | \mathbf{X}^{(\text{joint})}) \leftarrow -\frac{\hat{\mathbf{X}}^{(\text{joint})} - \mathbf{X}^{(\text{joint})}}{\sigma^2}$ ,
   /* Score Function Training */
4    $\mathcal{L}_{\text{DSM}} \leftarrow \frac{1}{2} \mathbb{E}_{q_\sigma(\hat{\mathbf{X}}^{(\text{joint})} | \mathbf{X}^{(\text{joint})})} [\|\nabla_{\hat{\mathbf{X}}^{(\text{joint})}} \log \hat{p}(\mathbf{X}^{(\text{joint})}) - \nabla_{\hat{\mathbf{X}}^{(\text{joint})}} \log q_\sigma(\hat{\mathbf{X}}^{(\text{joint})} | \mathbf{X}^{(\text{joint})})\|^2]$ ,
5    $\theta \leftarrow \theta - lr \times \nabla_\theta \mathcal{L}_{\text{DSM}}$ .
6 end

```

On this basis, the algorithm for KnewImp is summarized as follows:

Algorithm 4: KnewImp Algorithm for MDI

Input: Missing Data $\mathbf{X}^{(\text{miss})}$, and Mask Matrix M .

Hyperparameters:

Loop Time: \mathcal{T} , Simulation Time: T , Discretization Step Size η , Bandwidth of RBF kernel h ,
Network Learning Rate lr , Training Epoch \mathcal{E} , and Network Hidden Unit HU_{score} .

```
1  $\mathbf{X}^{(\text{imp})} \leftarrow \text{Initialize}(\mathbf{X}^{(\text{imp})})$ 
2 while  $t < \mathcal{T}$  do
   | /* 'Estimate' Part                                     */
3   |  $\nabla_{\mathbf{X}^{(\text{joint})}} \log \hat{p}(\mathbf{X}^{(\text{joint})}) \leftarrow \text{Algorithm 3}$ 
   | /* 'Impute' Part                                     */
4   |  $\mathbf{X}^{(\text{imp})} \leftarrow \text{Algorithm 2}$ 
5 end
```

Appendix D Detailed Information for Experiments

D.1 Background & Simulation of Missing Data

According to reference [31], missing data can be classified into three categories: Missing Completely at Random (MCAR), where the absence of data is completely unrelated to any observed or unobserved variables; Missing at Random (MAR), where the likelihood of missing data depends solely on observed data; and Missing Not at Random (MNAR), where missingness is influenced by unobserved data. In the cases of MCAR and MAR, the patterns of missing data are considered ‘ignorable’ because it is unnecessary to explicitly model the distribution of the missing values. Conversely, MNAR scenarios, where missing data can introduce significant biases that are not easily corrected without imposing domain-specific assumptions, constraints, or parametric forms on the missingness mechanism, present more complex challenges [28, 13]. Therefore, our discussion is primarily focused on numerical tabular data within the MCAR and MAR contexts.

To simulate missing data, we adopt the methodologies outlined in reference [13]:

- **MAR:** Initially, a random subset of features is selected to remain non-missing. The masking of the remaining features is conducted using a logistic model, which employs the non-missing features as predictors. This model is parameterized with randomly selected weights, and the bias is adjusted to achieve the desired missingness rate.
- **MCAR:** For each data point, the masking variable is generated from a Bernoulli distribution with a predetermined fixed mean, ensuring that the probability of missingness is the same across all data points.
- **MNAR:** Although MNAR scenarios are not the primary focus of this manuscript, we include experiments in this context. Missingness is introduced either by additional masking of the MAR-selected features using a Bernoulli process with a fixed mean, or through direct self-masking of values using interval-censoring techniques. In this paper, we mainly consider the former strategy. In other words, the mechanism of MNAR we used in this paper is identical to the previously described MAR mechanism, but the inputs of the logistic model are then masked by an MCAR mechanism.

Based on this, the datasets listed in Table D.1 are adopted in this paper.

Table D.1: Detailed dataset descriptions, where ‘Dimension’ denotes the variate number of each dataset. ‘Numer’ denotes the total number of item.

Abbreviation	Dataset Name	Numer (N)	Dimension (D)
BT	Blood Transfusion	748	4
BCD	Breast Cancer Diagnostic	569	30
CC	Concrete Compression	1030	7
CBV	Connectionist Bench Vowel	990	10
IS	Ionosphere	351	34
PK	Parkinsons	195	23
QB	QSAR Biodegradation	1055	41
WQW	Wine Quality White	4898	11

D.2 Training Protocols of Different Models

In this study, we employ a two-layer Multi-Layer Perceptron (MLP) to model $\nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{miss})})$. Each layer is configured with 256 hidden units (HU_{score}). The activation function is set as ‘Swish’ function [30], and the variance scale σ for DSM is set as 0.1. The network is trained using an Adam optimizer with a learning rate of 1.0×10^{-3} , and the batch size is dynamically set to N . For the ‘impute’ part, we specify a simulation time (T) of 500, a step size of 0.1, and a bandwidth (h) of 0.5. The loop time \mathcal{T} for KnewImp is set as 2. For baseline models, the batch size is uniformly set at 512. Models incorporating neural architectures are optimized with the Adam optimizer at a learning rate of 1.0×10^{-2} , in line with the practices recommended by Kingma and Ba [18]. The MIWAE model features a latent dimension of 16 and 32 hidden units. The settings for the TDM model include 16 hidden units per layer and two layers. For the CSDI_T and MissDiff models, the parameters are set as follows: particle number at 50, diffusion embedding dimension at 128, batch size at 512 (for Sink and TDM, if $N < 512$, the batch size is set as $2^{\lfloor \frac{N}{2} \rfloor}$), and learning rate at 1.0×10^{-3} , and diffusion steps at 100.

To ensure fairness and reproducibility, all experiments are conducted on a workstation equipped with an Intel Xeon E5 processor with four cores, eight Nvidia GTX 1080 GPUs, and 128 GB of RAM. Each experiment is replicated at least five times, utilizing six distinct random seeds to guarantee robustness in the results.

D.3 Evaluation Protocols

Imputation methods are assessed using two metrics: the mean absolute error (MAE), which is a pointwise metric, and the squared Wasserstein distance (abbreviated as Wass), which evaluates empirical distributions. Based on reference [28], consider a dataset $\mathbf{X} \in \mathbb{R}^{N \times D}$ with missing values. For any entry (i, j) identified as missing, let $\mathbf{X}^{(\text{imp})}[i, d]$ represent the corresponding imputation, and $\mathbf{X}^{(\text{true})}[i, d]$ denote the ground truth. Define \mathbf{m}_0 as the total number of missing entries, $\mathbf{m}_0 := \#\{(i, d), \mathbf{M}[i, d] = 0\}$, and \mathbf{m}_1 as the number of data points that have at least one missing value, $\mathbf{m}_1 := \#\{i : \exists d, \mathbf{M}[i, d] = 0\}$. The set \mathbf{M}_1 encompasses indices of data points with any missing values, $\mathbf{M}_1 := \{i : \exists d, \mathbf{M}[i, d] = 0\}$. The metrics used to evaluate the accuracy of the imputation, MAE and Wass, are calculated as follows:

$$\text{MAE} := \frac{1}{\mathbf{m}_0} \sum_{(i,d): \mathbf{M}[i,d]=0} |\mathbf{X}^{(\text{true})}[i, d] - \mathbf{X}^{(\text{imp})}[i, d]|, \quad (\text{D.1})$$

$$\text{Wass} := \mathcal{W}_2^2 \left[\frac{1}{\mathbf{m}_1} \sum_{k=1}^K \Delta_{\mathbf{x}_{\mathbf{M}_1}^{(\text{imp})}}, \frac{1}{\mathbf{m}_1} \sum_{k=1}^K \Delta_{\mathbf{x}_{\mathbf{M}_1}^{(\text{true})}} \right], \quad (\text{D.2})$$

where $\Delta_{\mathbf{x}}$ is the Dirac distribution (measure) concentrated on \mathbf{x} .

Appendix E Additional Empirical Evidence

E.1 Additional Experimental Results with MNAR Scenario

In this subsection, we expand upon the results presented in Table 1 by including the MNAR scenario, as detailed in Table E.1. Additionally, we report on the outcomes of an ablation study and sensitivity analysis in Tables E.3 and E.4 and Figure E.1. These extended results lead to several pertinent observations:

- Across three different missing data scenarios, the models consistently exhibit the poorest performance under the MNAR condition. For instance, in the MNAR scenario, nearly all models show a significant decrease in imputation accuracy and an increase in standard deviation. This supports the assertion made in Appendix D.1 that addressing the MNAR scenario requires the incorporation of relevant domain knowledge to mitigate biases introduced by the pattern of missing data.
- The findings from the ablation study under the MNAR scenario are consistent with those observed in both MAR and MCAR scenarios in Section 4.3. This consistency underscores the importance of including the NER term and adopting the joint distribution modeling approach.
- Similarly, the results from the sensitivity analysis under the MNAR scenario align with those from MAR and MCAR scenarios in Section 4.4. This alignment reinforces our interpretations of model performance across different groups of hyperparameters under MAR and MCAR scenarios.

Table E.1: Performance of MAE and Wass metrics at 30% missing rate, and “*” marks that KnewImp outperforms significantly at p -value < 0.05 over paired samples t -test.

Scenario	Model	BT		BCD		CC		CBV		IS		PK		QB		WQW	
		MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass
MAR	CSDL_T	0.93*	3.44*	0.92*	18.2*	0.85*	2.82*	0.81*	3.86*	0.70*	16.9*	0.99*	15.9*	0.65*	20.1*	0.77*	4.13*
	MissDiff	0.85*	2.20*	0.91*	16.5*	0.87*	1.59*	0.83*	3.87*	0.72*	13.3*	0.92*	17.1*	0.63*	26.3*	0.75*	6.88*
	GAIN	0.75*	0.65*	0.54*	1.64*	0.75*	0.67*	0.68*	0.68*	0.56*	1.88*	0.59*	1.90*	0.65*	5.05*	0.68*	0.87*
	MIRACLE	0.62*	0.38	0.55*	1.92*	0.43	0.25	0.55*	0.46*	3.39*	35.1*	4.14*	34.1*	0.46	2.87*	0.51*	0.56
	MIWAE	0.64	0.53	0.52*	1.54*	0.76*	0.64*	0.82*	0.92*	0.50*	1.87*	0.65*	1.98*	0.55*	5.05*	0.62*	0.75*
	Sink	0.87*	0.92*	0.92*	3.84*	0.88*	0.83*	0.84*	0.98*	0.75*	2.43*	0.94*	3.61*	0.65*	4.71*	0.76*	1.04*
	TDM	0.83*	0.89*	0.83*	3.47*	0.81*	0.73*	0.76*	0.85*	0.62*	1.96*	0.86*	3.36*	0.59*	4.46*	0.73*	0.99*
	KnewImp	0.52	0.38	0.34	0.82	0.35	0.25	0.31	0.20	0.39	1.31	0.44	1.21	0.45	3.50	0.46	0.55
MCAR	CSDL_T	0.73*	1.93*	0.73*	15.5*	0.85*	2.71*	0.83*	3.79*	0.76*	15.2*	0.72*	12.4*	0.57*	19.9*	0.78*	4.11*
	MissDiff	0.72*	1.62*	0.73*	14.4*	0.84*	1.23*	0.82*	3.31*	0.75*	13.0*	0.71*	14.1*	0.56*	19.7*	0.76*	4.95*
	GAIN	0.72*	0.39*	0.38*	1.41*	0.78*	0.73*	0.72*	0.99*	0.57*	3.72*	0.46*	1.70	0.42*	3.62	0.73*	1.14*
	MIRACLE	0.52	0.15*	0.44*	1.94*	0.53*	0.35	0.61*	0.72*	2.99*	52.9*	3.38*	42.8*	0.35	2.71*	0.56*	0.75
	MIWAE	0.58*	0.24	0.50*	2.55*	0.76*	0.69*	0.83*	1.24*	0.64*	4.95*	0.51*	2.05*	0.48*	5.87*	0.67*	0.95*
	Sink	0.73*	0.48*	0.75*	4.39*	0.84*	0.85*	0.82*	1.27*	0.75*	4.94*	0.74*	3.36*	0.61*	5.92*	0.76*	1.25*
	TDM	0.68*	0.42*	0.63*	3.57*	0.77*	0.75*	0.77*	1.15*	0.66*	4.20*	0.64*	2.89*	0.52*	5.34*	0.74*	1.20*
	KnewImp	0.48	0.18	0.25	0.80	0.47	0.34	0.42	0.44	0.44	3.05	0.32	1.01	0.34	3.66	0.53	0.76
MNAR	CSDL_T	0.83*	2.29*	0.82*	15.7*	0.85*	2.78*	0.83*	3.83*	0.74*	15.5*	0.84*	12.2*	0.62*	19.8*	0.78*	4.09*
	MissDiff	0.78*	1.43*	0.81*	14.9*	0.84*	1.27*	0.83*	3.53*	0.72*	13.3*	0.81*	16.0*	0.61*	21.6*	0.76*	4.70*
	GAIN	0.77*	0.57*	0.62*	3.94*	0.78*	0.79*	0.78*	1.15*	0.71*	4.85*	0.70*	4.20*	0.76*	10.5*	0.75*	1.23*
	MIRACLE	0.63	0.35	0.60*	4.26*	0.52*	0.35	0.63*	0.77*	3.10*	55.6*	3.49*	44.8*	0.52*	5.61	0.58*	0.80
	MIWAE	0.66*	0.42	0.56*	3.31*	0.74*	0.68*	0.85*	1.30*	0.59*	4.33*	0.60*	3.06*	0.53*	7.21*	0.67*	0.97*
	Sink	0.79*	0.68*	0.83*	5.90*	0.83*	0.89*	0.84*	1.36*	0.75*	4.86*	0.84*	5.02*	0.64*	7.23*	0.77*	1.33*
	TDM	0.76*	0.64*	0.74*	5.18*	0.76*	0.77*	0.79*	1.24*	0.64*	4.02*	0.76*	4.54*	0.57*	6.45	0.74*	1.23*
	KnewImp	0.60	0.35	0.32	1.46	0.44	0.34	0.46	0.52	0.40	2.68	0.39	1.56	0.42	5.57	0.55	0.81

Table E.2: Standard deviation of MAE and Wass metrics at 30% missing rate.

Scenario	Model	BT		BCD		CC		CBV		IS		PK		QB		WQW	
		MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass
MAR	CSDL_T	4.8E-2	2.9E-1	5.9E-2	2.6E+0	2.7E-2	1.4E-1	2.0E-2	9.2E-2	2.1E-2	2.0E+0	3.9E-2	3.1E+0	2.0E-2	8.3E-1	1.8E-2	7.7E-2
	MissDiff	4.0E-2	4.9E-1	3.1E-2	2.6E+0	3.4E-2	2.6E-1	1.9E-2	1.2E+0	5.4E-2	5.6E-1	3.0E-2	2.4E+0	1.8E-2	4.9E+0	1.6E-2	1.3E+0
	GAIN	1.0E-1	1.6E-1	4.3E-2	2.3E+1	3.4E-2	8.9E-2	1.9E-2	4.0E-2	5.8E-2	3.4E-1	5.4E-2	3.7E-1	6.9E-2	8.4E-1	3.9E-2	5.6E-2
	MIRACLE	2.0E-2	6.3E-2	5.1E-2	4.0E-1	1.1E-2	2.0E-2	1.8E-2	2.1E-2	6.7E-1	1.2E+1	4.2E-1	5.6E+0	1.6E-2	1.9E-1	1.3E-2	2.8E-2
	MIWAE	6.5E-2	1.5E-1	5.2E-2	2.5E-1	6.1E-2	1.2E-1	2.4E-2	4.6E-2	5.3E-2	1.8E-1	2.7E-2	1.8E-1	3.8E-2	2.8E-1	1.9E-2	2.7E-2
	Sink	4.6E-2	1.2E-1	3.2E-2	1.6E-1	2.6E-2	7.4E-2	2.4E-2	5.5E-2	5.0E-2	2.0E-1	1.7E-2	9.3E-2	1.7E-2	7.6E-2	2.2E-2	4.4E-2
	TDM	4.5E-2	1.2E-1	1.9E-2	8.2E-2	3.4E-2	8.6E-2	2.6E-2	5.2E-2	6.2E-2	2.1E-1	2.7E-2	1.7E-1	1.1E-2	8.1E-2	2.1E-2	4.7E-2
KnewImp	2.0E-2	4.0E-2	2.7E-2	1.1E-1	5.6E-2	6.4E-2	1.6E-2	2.2E-2	1.9E-2	1.1E-1	1.1E-2	8.8E-2	1.9E-2	2.7E-1	1.6E-2	3.3E-2	
MCAR	CSDL_T	1.0E-2	1.5E-1	8.7E-3	5.7E-1	8.7E-3	8.2E-2	4.6E-3	4.6E-2	4.6E-3	3.6E-1	1.1E-2	8.7E-1	3.7E-3	3.1E-1	1.2E-2	5.0E-02
	MissDiff	6.4E-3	3.3E-1	8.2E-3	8.3E-1	3.5E-3	2.3E-1	5.9E-3	8.4E-1	7.1E-3	1.8E-1	4.6E-3	2.5E+0	6.2E-3	2.4E+0	4.1E-3	6.5E-1
	GAIN	6.1E-2	1.0E-1	7.9E-3	2.6E-2	2.4E-2	5.2E-2	1.4E-2	3.4E-2	1.8E-2	1.8E-1	4.5E-2	3.8E-1	3.7E-3	1.8E-1	2.0E-2	5.5E-2
	MIRACLE	2.6E-2	8.4E-3	1.6E-2	1.9E-1	1.7E-2	1.5E-2	5.2E-3	1.4E-2	4.3E-2	1.2E+0	4.6E-2	1.1E+0	1.0E-2	1.7E-1	1.1E-3	5.5E-3
	MIWAE	3.1E-2	3.9E-2	4.8E-3	4.9E-2	7.6E-3	1.3E-2	1.6E-2	4.3E-2	9.4E-3	1.3E-1	1.0E-2	7.8E-2	9.1E-3	2.7E-1	4.1E-3	9.5E-3
	Sink	7.3E-3	3.4E-2	4.6E-3	2.5E-2	7.0E-3	6.6E-3	4.5E-3	4.5E-3	4.2E-3	1.4E-1	1.0E-2	5.9E-2	3.3E-3	1.9E-1	3.9E-3	1.2E-2
	TDM	4.9E-3	2.8E-2	8.7E-3	3.1E-2	9.8E-3	6.6E-3	6.9E-3	7.9E-3	1.0E-3	1.9E-3	3.3E-3	3.6E-2	9.3E-3	1.5E-1	5.1E-3	1.3E-2
KnewImp	3.3E-3	3.7E-3	1.9E-3	4.6E-2	1.1E-2	1.8E-2	4.1E-3	1.8E-2	5.7E-3	1.1E-1	6.4E-3	3.7E-2	4.8E-3	1.7E-1	2.2E-3	1.1E-2	
MNAR	CSDL_T	2.9E-2	2.2E-1	8.7E-3	7.8E-1	2.2E-2	1.3E-1	7.4E-3	7.4E-2	1.0E-2	5.9E-1	2.2E-02	1.8E+0	2.6E-3	4.6E-1	2.6E-3	4.6E-1
	MissDiff	3.7E-2	3.7E-1	2.4E-3	9.7E-1	5.9E-3	2.4E-1	5.5E-3	8.2E-1	1.4E-2	3.3E-1	1.0E-2	2.1E+0	8.7E-3	3.3E+0	4.0E-3	5.2E-1
	GAIN	4.9E-2	1.2E-1	6.2E-2	6.9E-1	5.3E-2	8.6E-2	4.1E-2	9.3E-2	5.5E-3	4.8E-2	2.5E-2	4.7E-1	5.0E-2	1.2E+0	4.0E-2	1.0E-1
	MIRACLE	6.6E-2	9.5E-2	1.9E-2	4.7E-1	1.3E-2	1.3E-2	4.0E-3	1.7E-2	9.9E-2	3.5E+0	6.9E-2	1.6E+0	1.7E-2	1.7E-1	7.5E-3	1.2E-2
	MIWAE	3.3E-2	6.4E-2	8.3E-3	3.7E-2	2.4E-2	3.5E-2	3.0E-2	8.7E-2	6.6E-3	7.2E-2	2.3E-2	3.2E-1	1.2E-2	1.5E-1	9.1E-3	2.0E-2
	Sink	1.9E-2	6.2E-2	1.4E-2	1.4E-1	1.0E-2	6.3E-3	1.3E-2	3.9E-2	7.2E-3	5.1E-2	1.8E-2	3.5E-1	6.9E-3	1.6E-1	4.6E-3	2.9E-2
	TDM	2.2E-2	6.8E-2	1.4E-2	1.2E-1	9.4E-3	8.3E-3	1.5E-2	3.8E-2	2.0E-2	7.7E-2	1.8E-2	3.7E-1	3.9E-3	1.8E-1	7.1E-3	2.1E-2
KnewImp	2.5E-2	1.0E-1	3.9E-3	1.3E-1	1.9E-2	2.9E-2	8.4E-3	1.2E-2	9.0E-3	1.3E-1	8.5E-3	5.0E-2	7.1E-3	6.8E-1	5.8E-3	1.6E-2	

Table E.3: Ablation Study Results with missing rate at 30%, and “**” marks that KnewImp outperforms significantly at p -value < 0.05 over paired samples t -test. Best results are **bolded**.

Missing	NER	Joint	BT		BCD		CC		CBV		IS		PK		QB		WQW	
			MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass
MAR	✗	✗	0.96*	3.82*	1.05*	20.2*	1.04*	5.47*	0.86*	5.81*	0.67*	20.2*	1.06*	15.6*	0.72*	22.5*	0.79*	6.49*
	✗	✓	0.54	0.42	0.34	0.82	0.61*	0.40*	0.58*	0.47*	0.43*	1.34	0.46*	1.25*	0.47*	3.56*	0.75*	0.64*
	✓	✗	0.96*	3.83*	1.05*	20.3*	1.04*	5.49*	0.86*	5.83*	0.67*	20.2*	1.06*	15.7*	0.72*	22.5*	0.79*	6.51*
	✓	✓	0.52	0.38	0.34	0.82	0.35	0.25	0.31	0.20	0.39	1.31	0.44	1.21	0.45	3.50	0.46	0.55
MCAR	✗	✗	0.72*	2.11*	0.74*	16.7*	0.85*	3.72*	0.83*	5.22*	0.74*	18.4*	0.71*	12.7*	0.58*	20.1*	0.76*	5.57*
	✗	✓	0.52*	0.17*	0.25	0.79	0.62*	0.46*	0.61*	0.71*	0.46	3.05	0.34	1.09	0.36*	3.74*	0.58*	0.82*
	✓	✗	0.72*	2.12*	0.73*	16.8*	0.86*	3.73*	0.83*	5.24*	0.74*	18.4*	0.71*	12.8*	0.58*	20.1*	0.76*	5.60*
	✓	✓	0.48	0.18	0.25	0.80	0.47	0.34	0.42	0.44	0.44	3.05	0.32	1.01	0.34	3.66	0.53	0.76
MNAR	✗	✗	0.81*	2.47*	0.89*	18.2*	0.87*	3.85*	0.85*	5.26*	0.69*	17.6*	0.87*	13.0*	0.64*	20.6*	0.77*	5.71*
	✗	✓	0.62	0.37	0.32	1.47	0.61*	0.47*	0.64*	0.79*	0.44	2.79	0.43*	1.88*	0.44*	5.65	0.60*	0.87
	✓	✗	0.82*	2.57*	0.89*	18.3*	0.87*	3.86*	0.85*	5.28*	0.69*	17.7*	0.88*	13.5*	0.64*	20.7*	0.77*	5.73*
	✓	✓	0.60	0.35	0.32	1.46	0.44	0.34	0.46	0.52	0.40	2.68	0.39	1.56	0.42	5.57	0.55	0.81

Table E.4: Standard deviation of Ablation Study Results with missing rate at 30%.

Missing	NER	Joint	BT		BCD		CC		CBV		IS		PK		QB		WQW	
			MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass	MAE	Wass
MAR	✗	✗	6.1E-2	4.1E-1	4.8E-2	6.8E-1	1.1E-1	5.6E-1	4.6E-2	4.5E-1	4.4E-2	3.8E+0	1.4E-1	3.9E+0	7.2E-2	2.4E+0	5.0E-2	4.1E-1
MAR	✗	✓	6.9E-2	1.2E-1	2.7E-2	1.1E-1	2.8E-2	4.7E-2	2.5E-2	3.8E-2	2.7E-2	1.3E-1	1.0E-2	9.0E-2	1.8E-2	2.6E-1	2.8E-2	5.6E-2
MAR	✓	✗	6.1E-2	4.1E-1	4.8E-2	6.8E-1	1.1E-1	5.6E-1	4.6E-2	4.5E-1	4.4E-2	3.8E+0	1.4E-1	3.9E+0	7.2E-2	2.4E+0	5.0E-2	4.1E-1
MAR	✓	✓	2.0E-2	4.0E-2	2.7E-2	1.1E-1	5.6E-2	6.4E-2	1.6E-2	2.2E-2	1.9E-2	1.1E-1	1.1E-2	8.8E-2	1.9E-2	2.7E-1	1.6E-2	3.3E-2
MCAR	✗	✗	9.9E-3	5.8E-2	1.0E-2	2.3E-1	3.6E-3	6.0E-2	3.3E-3	2.0E-2	9.5E-3	4.2E-1	1.5E-2	2.3E-1	1.1E-2	1.0E+0	4.0E-3	1.4E-2
MCAR	✗	✓	5.0E-3	5.9E-3	2.8E-3	4.6E-2	1.2E-2	1.5E-2	9.4E-3	1.7E-2	6.4E-3	1.2E-1	1.0E-2	1.2E-1	7.3E-3	1.9E-1	9.4E-04	3.9E-3
MCAR	✓	✗	9.9E-3	5.7E-2	1.0E-2	2.3E-1	3.6E-3	6.0E-2	3.2E-3	2.1E-2	9.5E-3	4.2E-1	1.5E-2	2.3E-1	1.0E-2	1.0E+0	4.0E-3	1.4E-2
MCAR	✓	✓	3.3E-3	3.7E-3	1.9E-3	4.6E-2	1.1E-2	1.8E-2	4.1E-3	1.8E-2	5.7E-3	1.1E-1	6.4E-3	3.7E-2	4.8E-3	1.7E-1	2.2E-3	1.1E-2
MNAR	✗	✗	4.2E-2	1.5E-1	2.3E-2	8.5E-1	3.2E-2	1.8E-1	1.2E-2	5.3E-2	6.9E-3	1.2E-1	3.2E-2	9.6E-1	1.6E-2	8.3E-1	1.4E-2	1.0E-1
MNAR	✗	✓	4.0E-2	1.4E-1	3.4E-3	1.3E-1	1.8E-2	2.1E-2	4.8E-3	1.7E-2	1.0E-2	1.3E-1	1.1E-2	1.8E-1	8.0E-3	7.0E-1	7.7E-3	1.4E-2
MNAR	✓	✗	4.8E-2	1.4E-1	2.4E-2	8.5E-1	3.3E-2	1.8E-1	1.2E-2	5.3E-2	6.9E-3	1.2E-1	1.9E-2	2.5E-1	1.6E-2	8.3E-1	1.4E-2	1.0E-1
MNAR	✓	✓	2.5E-2	1.0E-1	3.9E-3	1.3E-1	1.9E-2	2.9E-2	8.4E-3	1.2E-2	9.0E-3	1.3E-1	8.5E-3	5.0E-2	7.1E-3	6.8E-1	5.8E-3	1.6E-2

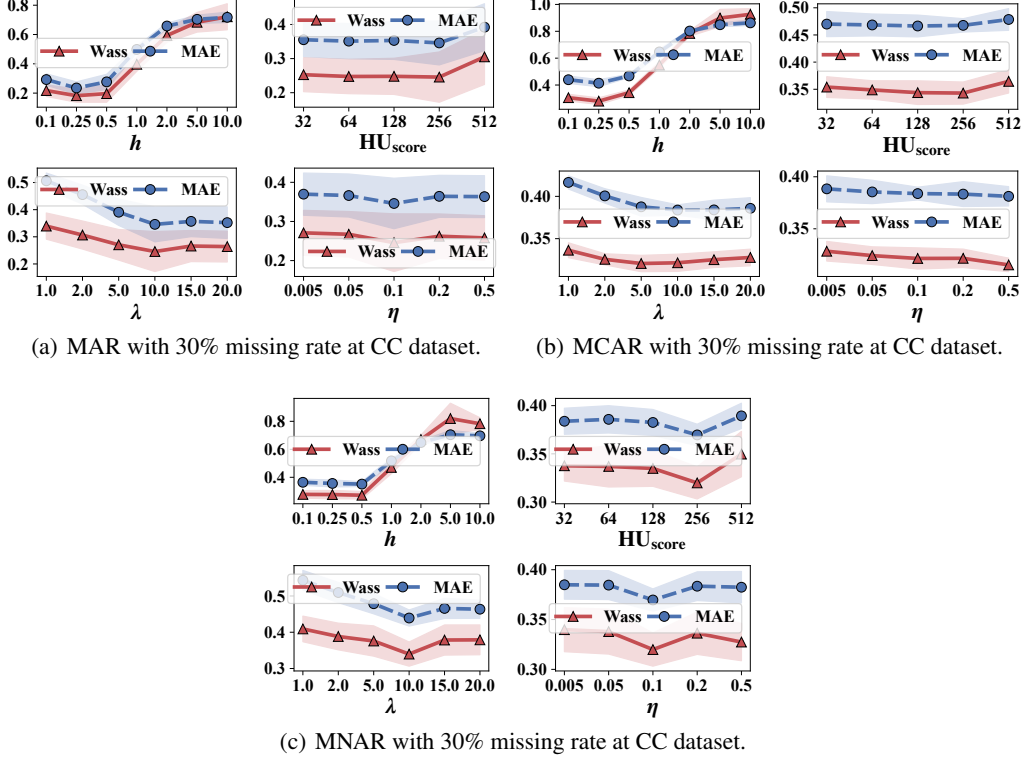


Figure E.1: Parameter sensitivity of KnewImp on bandwidth for kernel function (h), hidden unit of score network HU_{score} , NER weight λ , and discretization step η for Eq. (10) on CC dataset. Mean values and one standard deviations from mean are represented by scatters and shaded area, respectively.

E.2 Time Complexity Analysis

In this subsection, we present an analysis of the complexity of time for our KnewImp approach. The complexity analysis is based on the algorithms described in Algorithms 2 to 4. We begin by estimating the time complexity of the score function $\nabla_{\mathbf{X}^{(\text{joint})}} \log \hat{p}(\mathbf{X}^{(\text{joint})})$. Assuming the number of layers in the score network is L and each layer has an equal number of hidden units denoted as HU_{score} , the time complexity for the imputation algorithm defined in Algorithms 2 and 4 is detailed as follows:

1. Impute Part:

- **Score function computation:** The time complexity for computing the score function is expressed as:

$$\mathcal{O} \left[2 \times N \times \left(D \times \text{HU}_{\text{score}} + (L - 1) \times \text{HU}_{\text{score}}^2 \right) \right], \quad (\text{E.1})$$

where the factor 2 accounts for the backward propagation needed during the score function computation.

- **Kernel function and its gradient:** Employing the RBF kernel $\mathcal{K}(\mathbf{X}, \tilde{\mathbf{X}}) := \exp\left(-\frac{\|\mathbf{X} - \tilde{\mathbf{X}}\|^2}{2h^2}\right)$, the gradient with respect to $\tilde{\mathbf{X}}$ is analytically determined as:

$$[\nabla_{\tilde{\mathbf{X}}} \mathcal{K}(\mathbf{X}, \tilde{\mathbf{X}})][:, j] = -\frac{1}{h^2} \left\{ [\mathcal{K}(\mathbf{X}, \tilde{\mathbf{X}}) \times \tilde{\mathbf{X}}][:, j] + \tilde{\mathbf{X}}[:, j] \odot \sum_{j=1}^D \mathcal{K}(\mathbf{X}, \tilde{\mathbf{X}})[:, j] \right\}. \quad (\text{E.2})$$

The time complexities for calculating the kernel function and its gradient are specified in Eqs. (E.3) and (E.4):

$$\mathcal{O} \left[N^2 \times D + N^2 \right], \quad (\text{E.3})$$

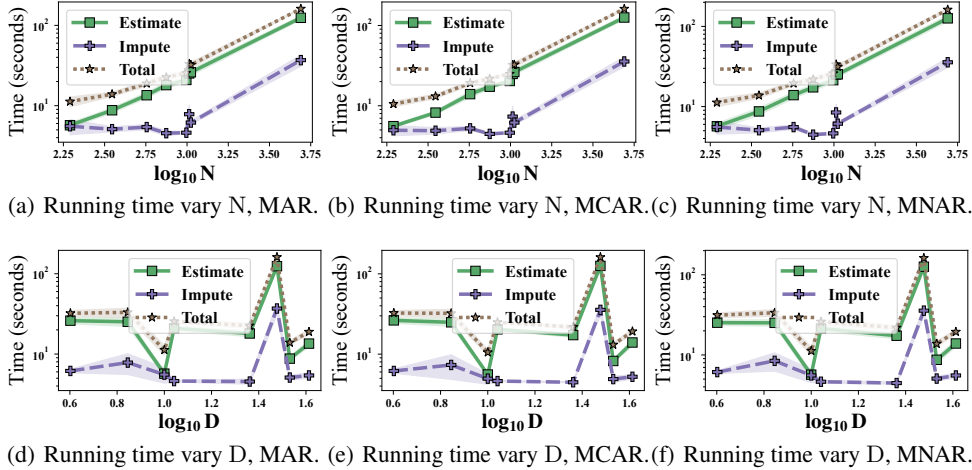


Figure E.2: Average computation time. The scatters and shaded areas indicate the mean and one standard deviation from the mean, respectively.

$$\mathcal{O} [N^2 \times D + N^2 + N \times D]. \quad (\text{E.4})$$

2. **Estimate Part:** Building on the previous item, the time complexity for the estimation algorithm defined in Algorithm 3 is given as:

$$\mathcal{O} [4 \times N \times (D \times \text{HU}_{\text{score}} + (L - 1) \times \text{HU}_{\text{score}}^2)], \quad (\text{E.5})$$

where the factor of 4 comprises three distinct components: backward propagation (1), forward propagation (1), and the acquisition of the sample-wise score function (2). Note that the network parameter size is substantially smaller than the number of data points, thereby making the forward computation of the score function the primary factor in time complexity.

Based on the analysis outlined above, we explore how computational complexity varies with different dataset sizes N and the number of features D , as shown in Figs. E.2 (a) and (b), respectively. From these figures, it is evident that computational time increases with the dataset size N . However, changes in the number of features D do not significantly affect the computation time. This observation underscores that the primary determinant of computational complexity in our context is the dataset size, aligning with our theoretical analysis, which indicates a quadratic relationship between time complexity and the size of the dataset N for the ‘Impute’ part, and $N \gg D$ for the ‘Estimate’ part.

Moreover, the data reveals that the total computational time is predominantly governed by ‘Estimate’ part of our KnewImp approach. This suggests that the training of the score function represents a critical bottleneck in the efficiency of the KnewImp algorithm. Therefore, accelerating the KnewImp algorithm crucially hinges on reducing the computational demands of the ‘Estimate’ part.

E.3 Convergence Analysis

In this subsection, we explore the convergence of the Impute part as defined in Algorithm 2 within our KnewImp approach. Prior to delving into this discussion, it is essential to establish a clear definition of convergence:

Definition E.1. A sequence $\{\mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_T\}$ is said to be convergent if there exists a real number \mathcal{G} such that for any given positive number ε ($\varepsilon > 0$), there exists a positive integer N , such that for all indices n greater than N , the corresponding terms $\mathcal{F}_n, n \geq N$ satisfy the inequality $|\mathcal{F}_n - \mathcal{G}| < \varepsilon$.

Based on Definition E.1, if a sequence is either monotonically increasing or monotonically decreasing and bounded (either bounded above or bounded below), then it is guaranteed to converge according to the celebrated monotone convergence theorem (Section 3.14 in reference [32]). Based on this, we first prove the following proposition for the convergence in ‘Impute’ part of our approach:

Proposition E.1. The convergence of the ‘Imputer’ part can be guaranteed, given that the step size η is small enough.

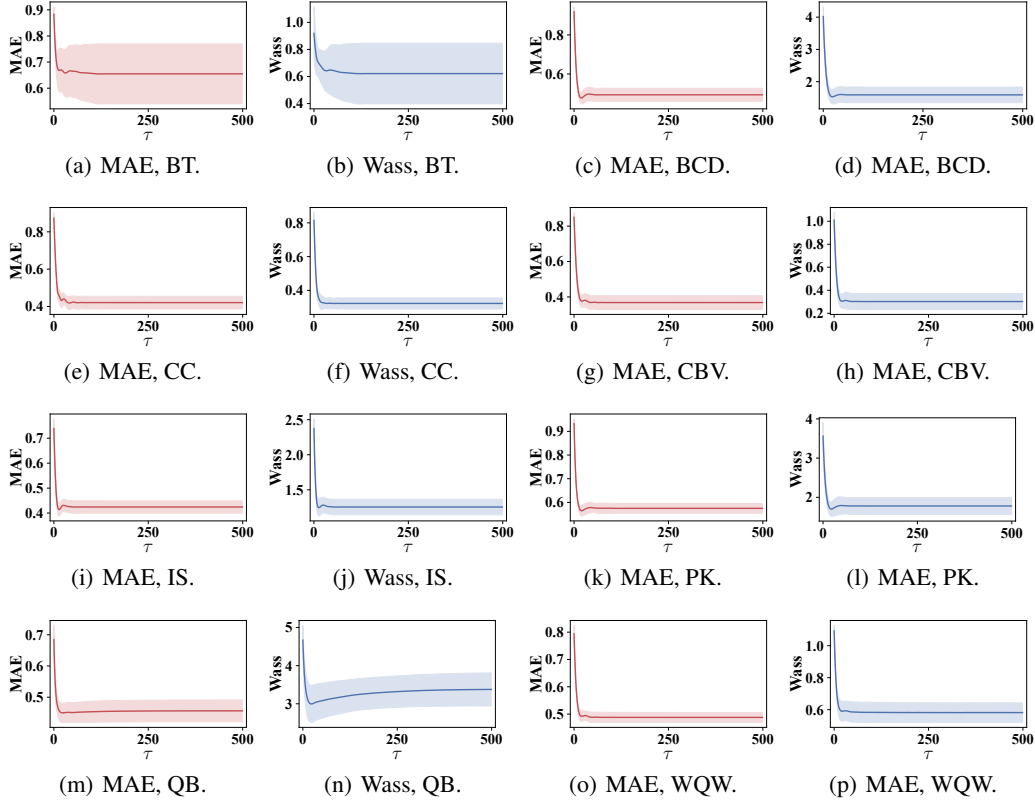


Figure E.3: Evolution of evaluation metrics along iteration time τ under MAR scenario at 30% missing rate. The shaded area indicates the ± 1.0 standard deviation uncertainty interval.

Proof. First, let's reformulate the velocity field as follows:

$$\begin{aligned}
& u(\mathbf{X}^{(\text{joint})}, \tau) \\
&= \mathbb{E}_{r(\tilde{\mathbf{X}}^{(\text{joint})}, \tau)} \left\{ \begin{aligned} & -\lambda \nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \mathcal{K}(\mathbf{X}^{(\text{joint})}, \tilde{\mathbf{X}}^{(\text{joint})}) \\ & + [\nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \log \hat{p}(\tilde{\mathbf{X}}^{(\text{joint})})]^\top \mathcal{K}(\mathbf{X}^{(\text{joint})}, \tilde{\mathbf{X}}^{(\text{joint})}) \end{aligned} \right\} \\
&\stackrel{(i)}{=} \mathbb{E}_{r(\tilde{\mathbf{X}}^{(\text{joint})}, \tau)} \left\{ \begin{aligned} & \lambda [\nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \log r(\tilde{\mathbf{X}}^{(\text{joint})})]^\top \mathcal{K}(\mathbf{X}^{(\text{joint})}, \tilde{\mathbf{X}}^{(\text{joint})}) \\ & + [\nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \log \hat{p}(\tilde{\mathbf{X}}^{(\text{joint})})]^\top \mathcal{K}(\mathbf{X}^{(\text{joint})}, \tilde{\mathbf{X}}^{(\text{joint})}) \end{aligned} \right\} \quad (\text{E.6}) \\
&= \int r(\tilde{\mathbf{X}}^{(\text{joint})}) \left\{ \begin{aligned} & \lambda \nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \log r(\tilde{\mathbf{X}}^{(\text{joint})}) \\ & + \nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \log \hat{p}(\tilde{\mathbf{X}}^{(\text{joint})}) \end{aligned} \right\}^\top \mathcal{K}(\mathbf{X}^{(\text{joint})}, \tilde{\mathbf{X}}^{(\text{joint})}) d\tilde{\mathbf{X}}^{(\text{joint})} \\
&= \int \left\{ \begin{aligned} & \lambda \nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \log r(\tilde{\mathbf{X}}^{(\text{joint})}) \\ & + \nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \log \hat{p}(\tilde{\mathbf{X}}^{(\text{joint})}) \end{aligned} \right\}^\top \mathcal{K}(\mathbf{X}^{(\text{joint})}, \tilde{\mathbf{X}}^{(\text{joint})}) dr(\tilde{\mathbf{X}}^{(\text{joint})}),
\end{aligned}$$

where (i) is based on integration by parts.

Based on this reformulation, the inner product can be given as follows:

$$\begin{aligned}
& \frac{d\mathcal{F}_{\text{joint-NER}}}{d\tau} \\
&= \int \left\langle \nabla_{\mathbf{X}^{(\text{miss})}} \frac{\delta \mathcal{F}_{\text{joint-NER}}}{\delta r(\mathbf{X}^{(\text{joint})})}, u(\mathbf{X}^{(\text{joint})}, \tau) \right\rangle dr(\mathbf{X}^{(\text{miss})}) \\
&= \iint \left\{ \begin{array}{l} \lambda \nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \log r(\tilde{\mathbf{X}}^{(\text{joint})}) \\ + \nabla_{\tilde{\mathbf{X}}^{(\text{miss})}} \log \hat{p}(\tilde{\mathbf{X}}^{(\text{joint})}) \end{array} \right\}^\top \mathcal{K}(\mathbf{X}^{(\text{joint})}, \tilde{\mathbf{X}}^{(\text{joint})}) \times \\
& \qquad \qquad \qquad \left\{ \begin{array}{l} \lambda \nabla_{\mathbf{X}^{(\text{miss})}} \log r(\mathbf{X}^{(\text{joint})}) \\ + \nabla_{\mathbf{X}^{(\text{miss})}} \log \hat{p}(\mathbf{X}^{(\text{joint})}) \end{array} \right\} dr(\tilde{\mathbf{X}}^{(\text{joint})}) dr(\mathbf{X}^{(\text{joint})}) \\
& \stackrel{(i)}{\geq} 0,
\end{aligned} \tag{E.7}$$

where the (i) is predicated on the requirement that the kernel function, $\mathcal{K}(\cdot, \cdot)$, is semi-positive definite; according to the above-mentioned derivation, we can conclude that the evolution of $\mathcal{F}_{\text{joint-NER}}$ is monotonic increasing along τ . Furthermore, $\mathcal{F}_{\text{joint-NER}}$ satisfies the following inequality:

$$\begin{aligned}
& \mathcal{F}_{\text{joint-NER}} \\
& \leq \mathcal{F}_{\text{joint-NER}} - (\lambda + 1) \mathbb{E}_{r(\mathbf{X}^{(\text{joint})})} [\log r(\mathbf{X}^{(\text{joint})})] \\
& = - \mathbb{D}_{\text{KL}} [r(\mathbf{X}^{(\text{joint})}) || \hat{p}(\mathbf{X}^{(\text{joint})})] \\
& \leq 0,
\end{aligned} \tag{E.8}$$

which indicates that $\mathcal{F}_{\text{joint-NER}}$ is upper-bounded by 0.

According to Eqs. (E.7) and (E.8), the cost functional $\mathcal{F}_{\text{joint-NER}}$, driven by the velocity field $u(\mathbf{X}^{(\text{joint})}, \tau)$ along τ , converges. Building on this, employing a smaller step size η results in the iteration curve of $\mathcal{F}_{\text{joint-NER}}$ more closely approximating the ODE defined in Eq. (E.7). Consequently, a smaller η leads to a sequence where $\mathcal{F}_{\text{joint-NER}}$ monotonically increases, aligning with the theoretical expectations of the ODE behavior. \square

Unfortunately, directly obtaining $\mathcal{F}_{\text{joint-NER}}$ is intractable. Nevertheless, we can still observe the changes in Wass and MAE across iteration time τ to demonstrate the convergence of the 'Impute' part. To this end, we present the convergence trends along τ in Figures E.3 to E.5. These figures illustrate that both MAE and Wass generally decrease as the iteration epochs increase and eventually stabilize after $\tau = 250$. This observed behavior supports our theoretical findings regarding the convergence of the 'Impute' part.

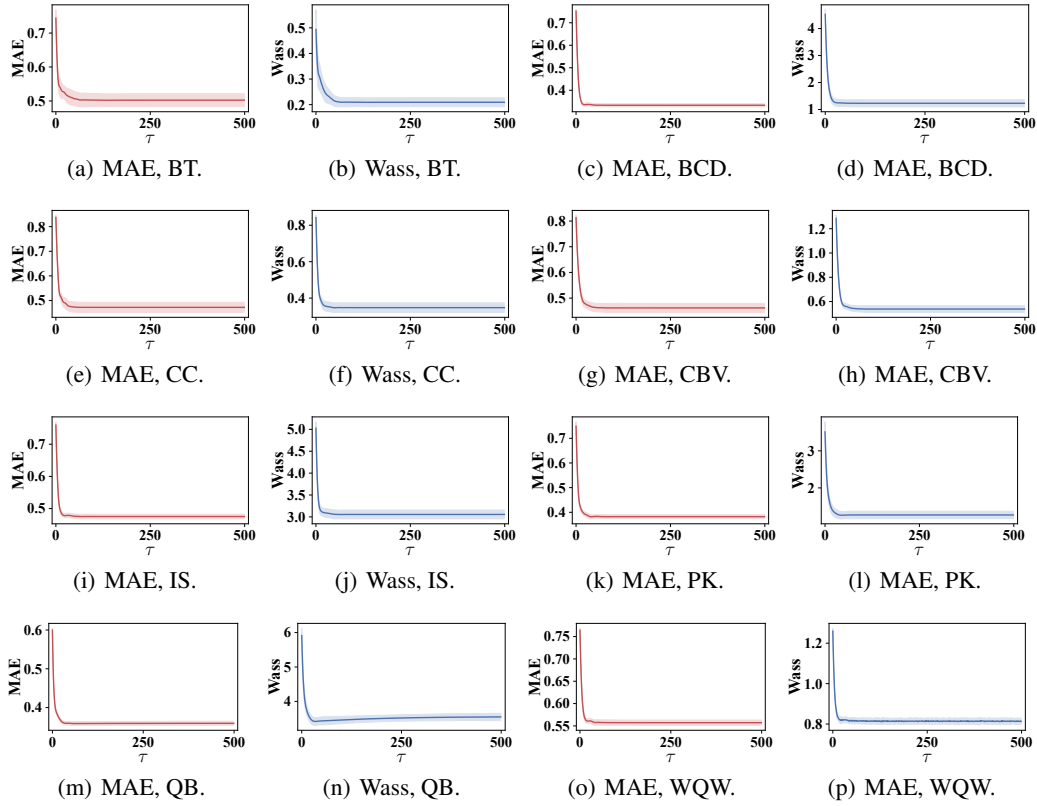


Figure E.4: Evolution of evaluation metrics along iteration time τ under MCAR scenario at 30% missing rate. The shaded area indicates the ± 1.0 standard deviation uncertainty interval.

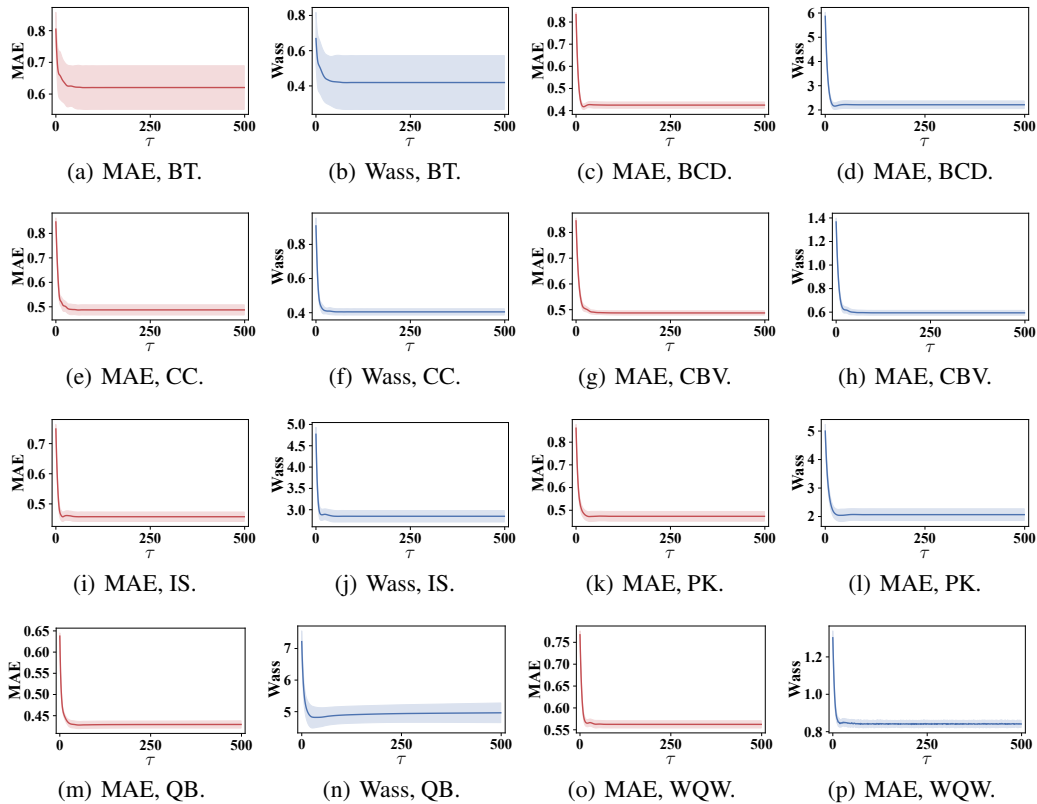


Figure E.5: Evolution of evaluation metrics along iteration time τ under MNAR scenario at 30% missing rate. The shaded area indicates the ± 1.0 standard deviation uncertainty interval.

E.4 Baseline Comparison Vary Different Missing Rates and Scenarios

In this subsection, we present an extended analysis of model performance across varying missing data rates, as detailed in Figures [E.6](#) to [E.11](#) . The results demonstrate that our KnewImp approach performs competitively across a broad spectrum of missing data regimes.

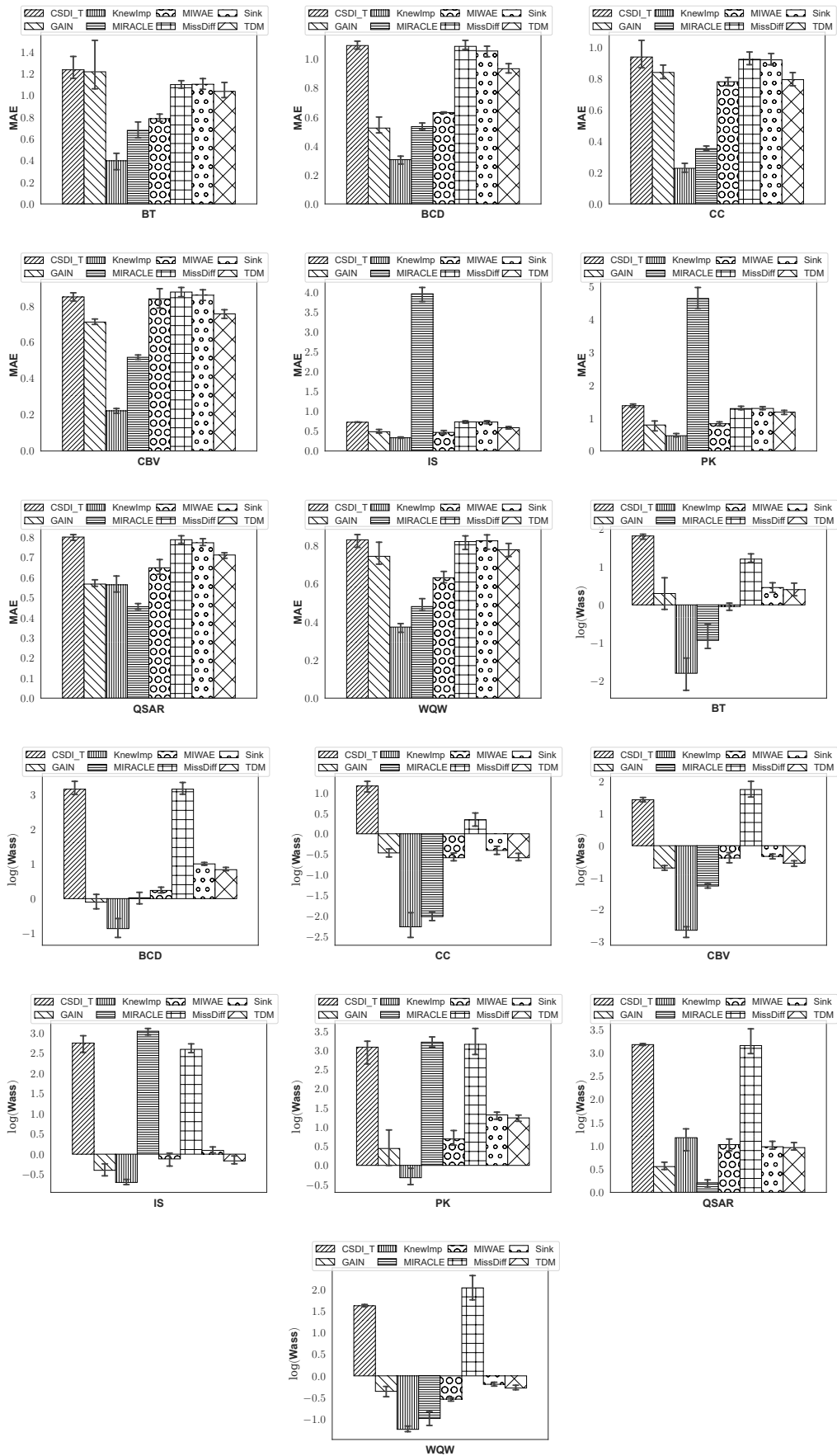


Figure E.6: Imputation accuracy comparison for MAR scenario at 10% missing rate. The error bars indicate the 100% confidence intervals.

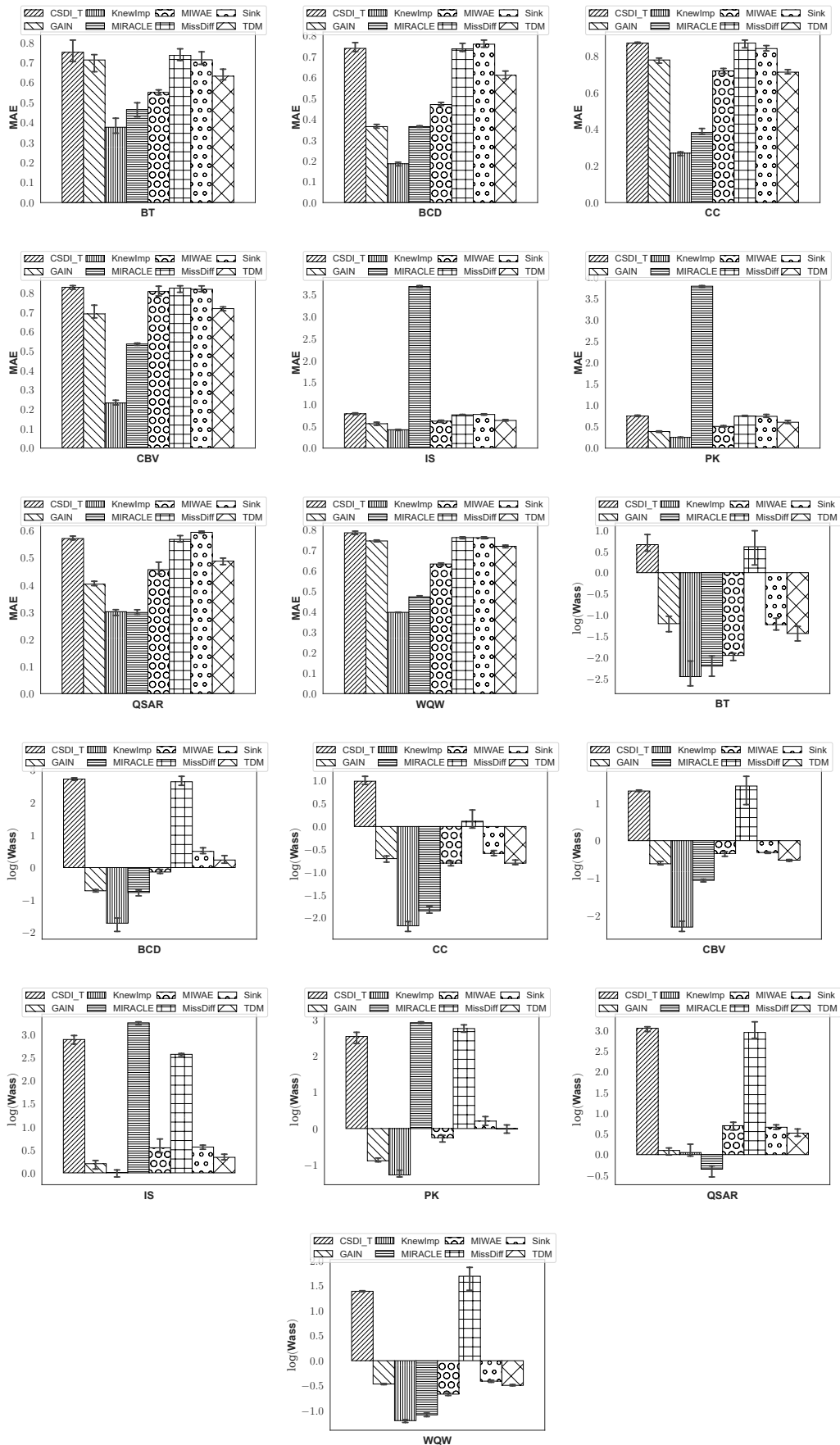


Figure E.7: Imputation accuracy comparison for MCAR scenario at 10% missing rate. The error bars indicate the 100% confidence intervals.

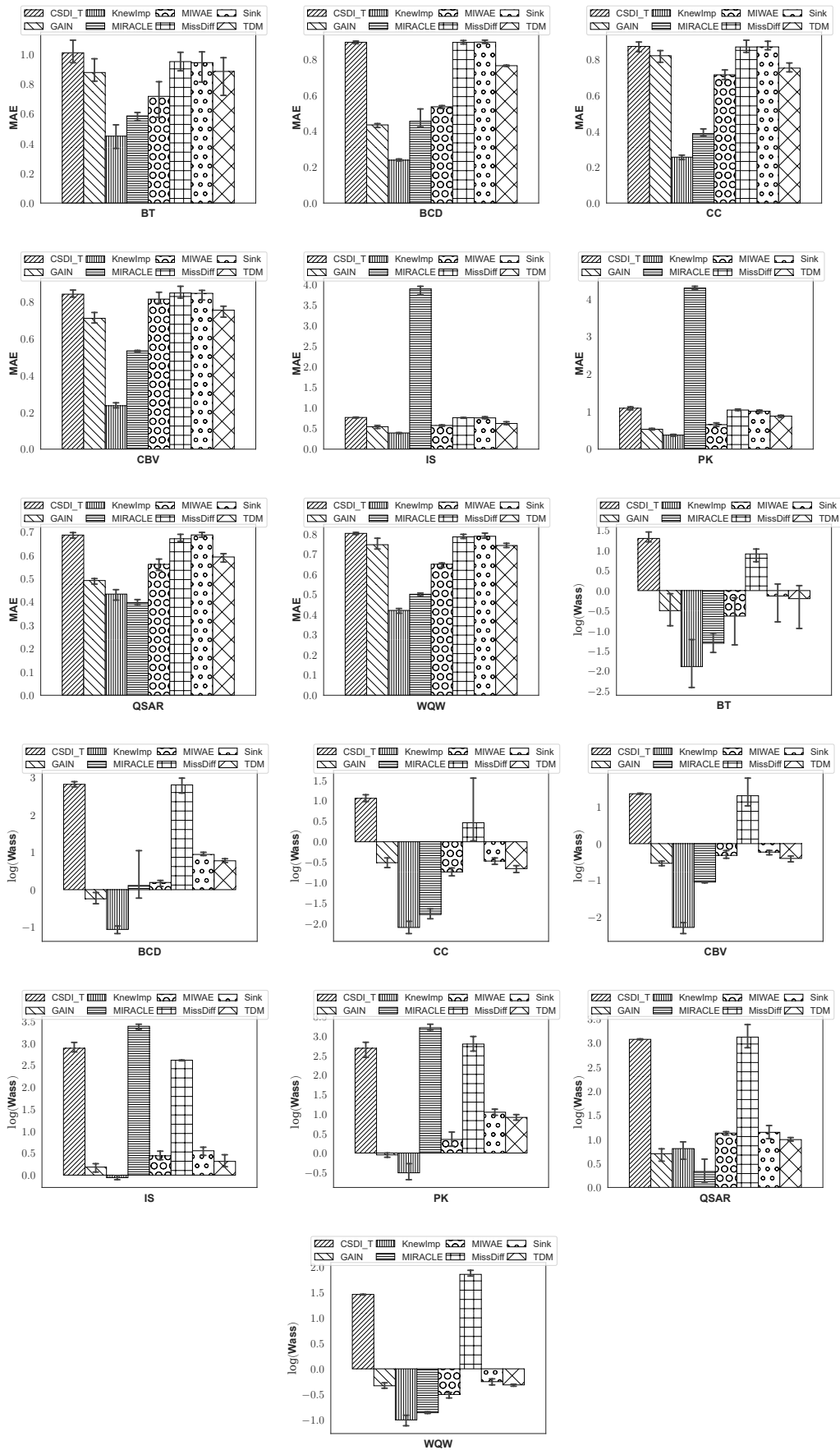


Figure E.8: Imputation accuracy comparison for MNAR scenario at 10% missing rate. The error bars indicate the 100% confidence intervals.

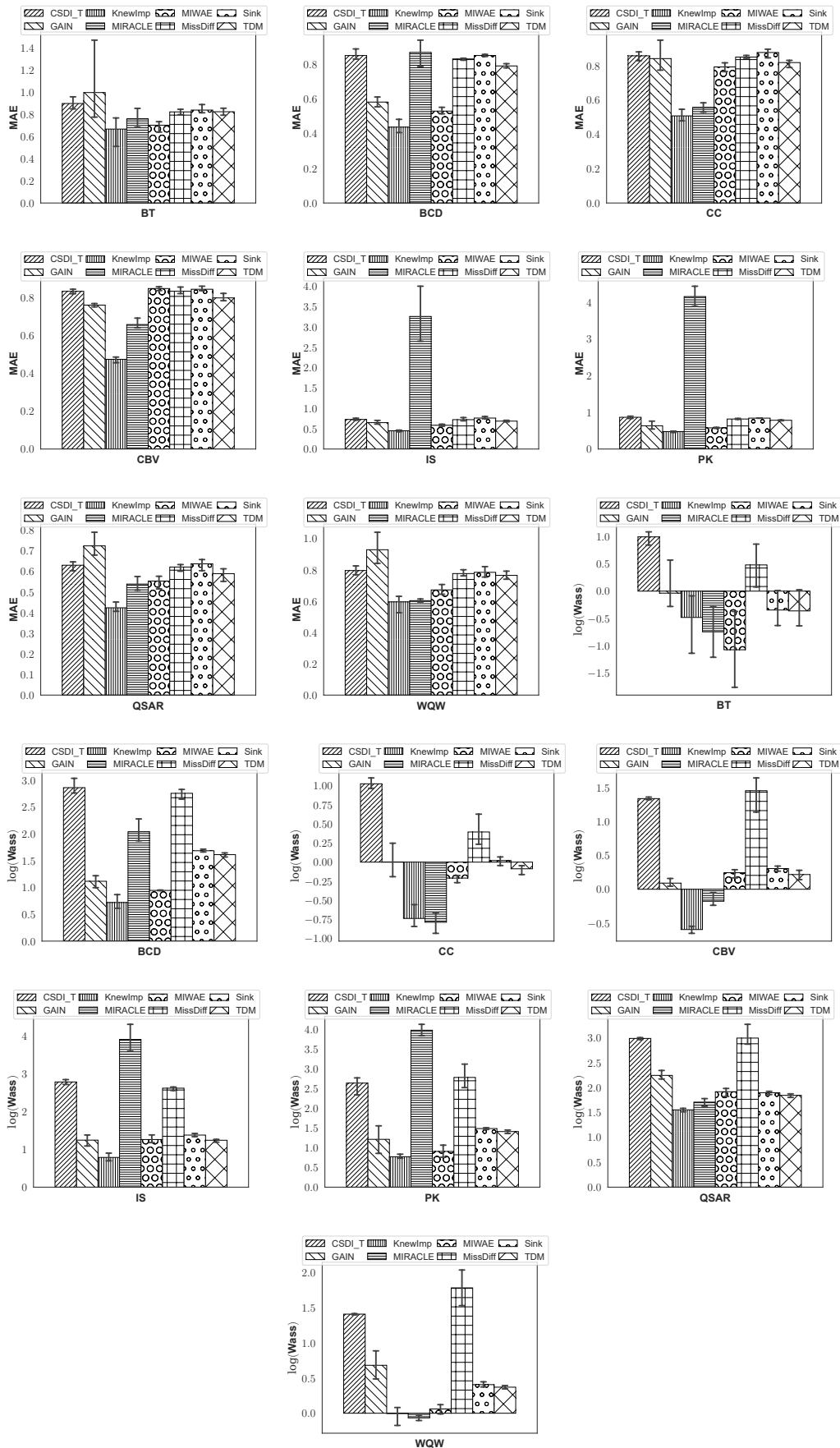


Figure E.9: Imputation accuracy comparison for MAR scenario at 50% missing rate. The error bars indicate the 100% confidence intervals.

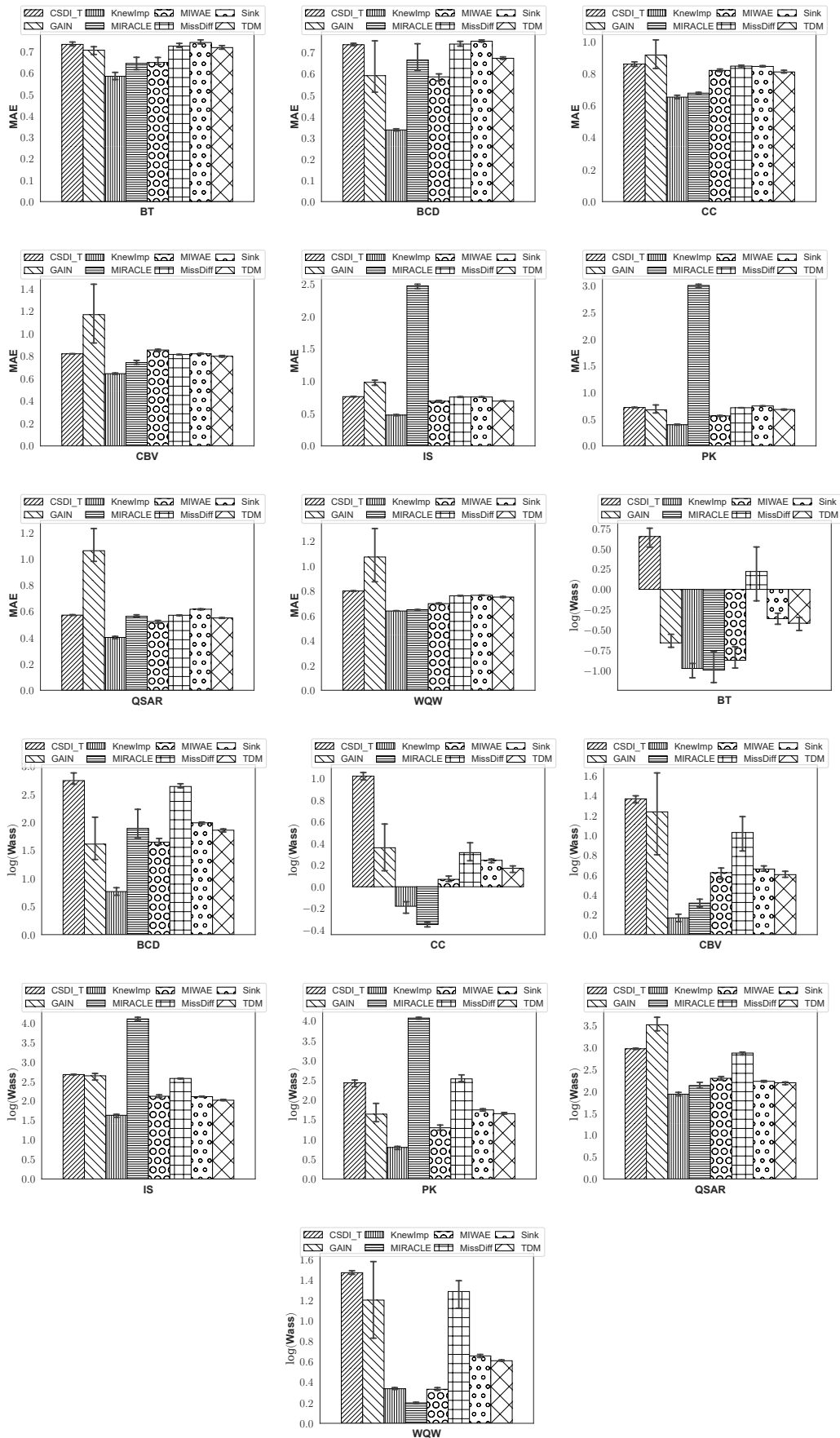


Figure E.10: Imputation accuracy comparison for MCAR scenario at 50% missing rate. The error bars indicate the 100% confidence intervals.

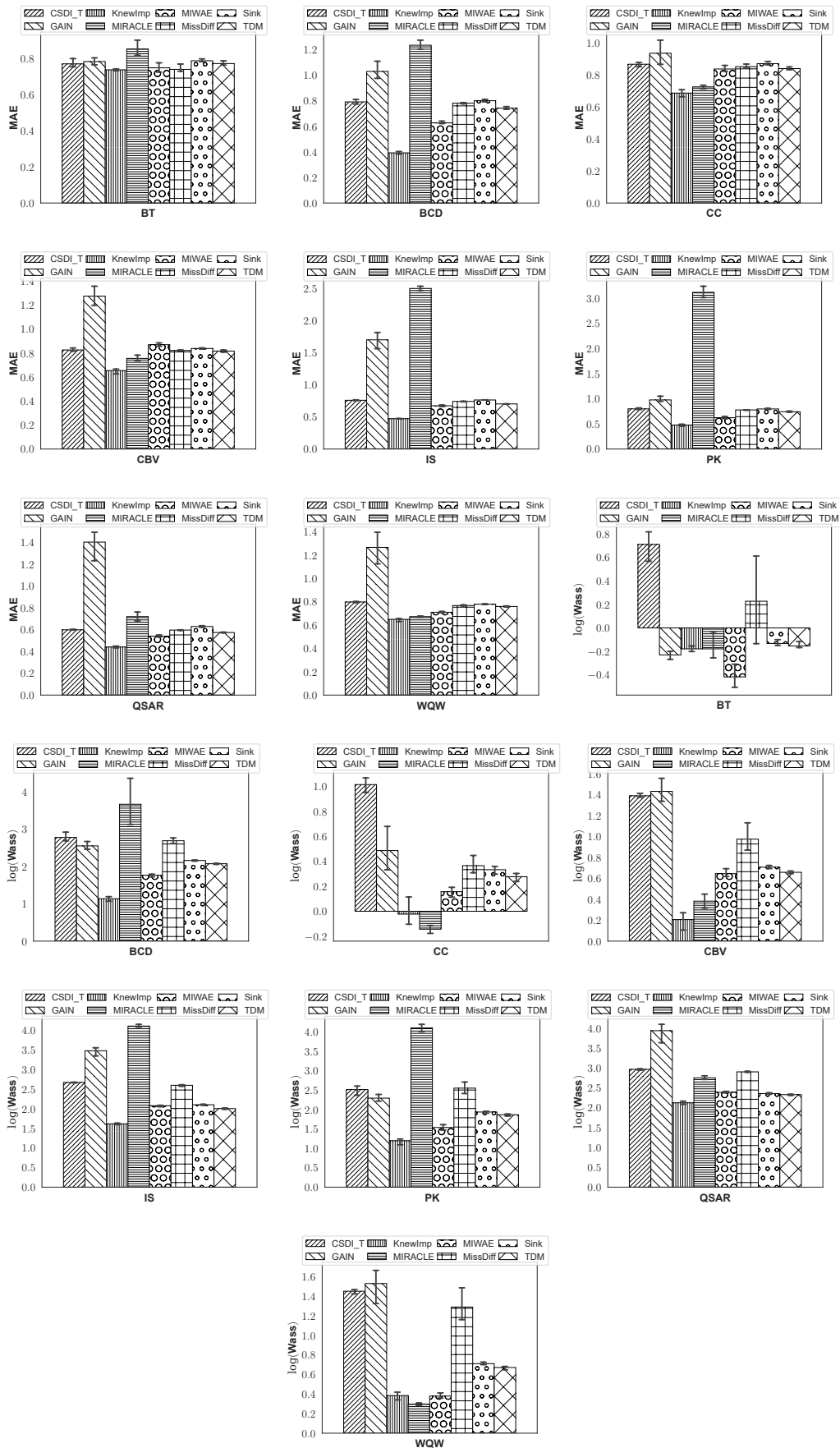


Figure E.11: Imputation accuracy comparison for MNAR scenario at 50% missing rate. The error bars indicate the 100% confidence intervals.

Appendix F Limitations & Future Directions and Broader Impact

F.1 Limitations & Future Directions

The limitations and future research directions of this work can be summarized as follows:

- **Utilization of Kernel Function:** During the derivation of the velocity field, we employ RKHS to ensure implementation easiness. However, this regularization term may impose restrictions on the velocity field’s direction, potentially limiting imputation accuracy in high-dimensional settings. Additionally, the computational complexity tends to scale quadratically with dataset size increases. Exploring alternative regularization terms [8] to replace RKHS presents a promising direction for future research.
- **Training of Score Function:** As discussed in Section E.2, the runtime of KnewImp is predominantly governed by the ‘Estimate’ part. Investigating techniques to reduce the training costs of this part, such as employing sliced score matching [35], represents an intriguing area for future exploration.
- **WGF Framework:** The WGF framework currently operates as a first-order system where each sample is equally weighted. A critical advancement would be the incorporation of second-order systems, such as Hamiltonian dynamics [41, 45], and other gradient flows like, Fisher-Rao gradient flow [53] that assign variable weights to samples. This adaptation aims to decrease computational times inherently.

F.2 Broader Impact Statement

MDI and DMs are pivotal areas within machine learning, each boasting a wide array of real-world applications. While numerous applications exist, this paper does not single out any specific ones; instead, it focuses on addressing fundamental challenges in these fields. This study significantly advances the application of DMs for MDI by tackling prevalent issues such as inaccurate imputation and challenging training processes. We believe that the insights garnered here can be applied to related domains, such as probabilistic time-series forecasting and image inpainting, where accuracy is often more critical than diversity in results. A common challenge across these domains is the nuanced need for precision over variety, which can lead to overlooked opportunities in model application and development. Our proposed method provides a fresh perspective on these tasks through an optimization lens. It evaluates the appropriateness of directly applying existing diffusion models to these tasks and, where necessary, proposes the derivation of novel algorithms. This approach not only enhances the understanding of the underlying mechanisms but also paves the way for more targeted and effective solutions in the future.