

# Missing Data Imputation by Reducing Mutual Information with Rectified Flows

Jiahao Yu<sup>1</sup>, Qizhen Ying<sup>2</sup>, Leyang Wang<sup>3</sup>, Ziyue Jiang<sup>1</sup>, Song Liu<sup>1\*</sup>

<sup>1</sup>*School of Mathematics, University of Bristol, United Kingdom*

<sup>2</sup>*Department of Computer Science, University of Oxford, United Kingdom*

<sup>3</sup>*Department of Computer Science, University College London, United Kingdom*

Correspond to <song.liu@bristol.ac.uk>

## Abstract

This paper introduces a novel iterative method for missing data imputation that sequentially reduces the mutual information between data and their corresponding missing mask. Inspired by GAN-based approaches, which train generators to decrease the predictability of missingness patterns, our method explicitly targets the reduction of mutual information. Specifically, our algorithm iteratively minimizes the KL divergence between the joint distribution of the imputed data and missing mask, and the product of their marginals from the previous iteration. We show that the optimal imputation under this framework corresponds to solving an ODE, whose velocity field minimizes a rectified flow training objective. We further illustrate that some existing imputation techniques can be interpreted as approximate special cases of our mutual-information-reducing framework. Comprehensive experiments on synthetic and real-world datasets validate the efficacy of our proposed approach, demonstrating superior imputation performance. The code of the proposed algorithm can be found at <https://github.com/aneugethubname/MIRI>.

## 1 Introduction

The problem of missing data, referring to absent components (NaNs) in a dataset, can arise in a wide range of scenarios [Rubin, 1976, Little and Rubin, 2019]. For example, in a survey dataset, respondents might skip some questions accidentally or intentionally, resulting in missing values. fMRI data could have missing voxels due to incomplete brain coverage and spatial variations in images acquired across subjects, [Vaden et al., 2012]. Simple “one-shot” methods replace missing entries with summary statistics, such as the sample mean or median of the observed data. In contrast, approaches like Multiple Imputation by Chained Equations (MICE) [van Buuren and Groothuis-Oudshoorn, 2011], MissForest [Stekhoven and Bühlmann, 2012], and HyperImpute [Jarrett et al., 2022] employ an iterative, dimension-wise scheme. At each iteration, they model one feature conditional on all others and draw from that model to replace the feature’s missing entries. They cycle through the features repeatedly until the imputations stabilize. Because the updates are sequential and applied one dimension at a time, these methods are difficult to parallelize and scale poorly to high-dimensional data such as images [Brini et al., 2024]. Their accuracy is also highly sensitive to the choice of per-variable models [van Buuren and Groothuis-Oudshoorn, 2011, Laqueur et al., 2022].

In recent years, research on generative modeling has made huge progress. Since imputation is naturally a data-generating task, generative models have been applied to impute missing data in some recent works. [Uehara et al., 2020] proposes to maximize the likelihood of the imputed data using an EM algorithm, where the likelihood function is modelled using a normalizing flow [Rezende and

Mohamed, 2015, Dinh et al., 2015]. Generative Adversarial Imputation Nets (GAIN) [Yoon et al., 2018] leverages Generative Adversarial Nets (GAN) [Goodfellow et al., 2014] to impute missing components. It trains a generator (a neural network) so that a binary classifier cannot distinguish whether a given component is imputed or not. [Li et al., 2019] proposes MisGAN to train a generator so that a discriminator cannot differentiate between generated data and observed data through missing masks.

Both GANs and normalizing flows come with their own challenges. GANs rely on unstable adversarial training and may suffer from mode collapse [Li et al., 2018, Zhang et al., 2018], while normalizing flows need specially designed architectures to ensure their invertibility, limiting their applications. Currently, diffusion models and *flow-based methods* represent the state of the art in generative modeling techniques [Song et al., 2021, Lipman et al., 2023, Liu et al., 2023]. They first train a velocity field according to some criterion. Then, they obtain samples from the target distribution by solving an ODE or SDE using the trained velocity field as the drift. The main challenge is that the flow-based methods are designed to transport reference samples to match the target distribution. However, it is unclear how to formulate the missing data imputation problem as a distribution transport problem. Nonetheless, some progress along this line has been made. For example, MissDiff Ouyang et al. [2023] learns a conditional score model from partially observed samples, then uses an inpainter [Lugmayr et al., 2022] to impute missing values.

Authors in [Liu et al., 2024] observe that GAIN’s criterion encourages the imputed dataset to be independent of the missingness pattern. In this paper, we continue their path to explore this idea but under a more rigorous framework. Our contributions are summarized as follows: (1) We introduce a novel imputation framework, named Mutual Information Reducing Iterations (MIRI), and prove that this iterative algorithm reduces the mutual information between imputed data and their corresponding missingness mask, when an optimal imputer is employed; (2) We show that an optimal MIRI imputer can be obtained by solving an ordinary differential equation (ODE), where the velocity field is trained using a rectified flow objective. This formulation naturally integrates flow-based generative modeling into the imputation process; (3) We show that several existing imputation methods can be viewed as approximate special cases of the MIRI framework; (4) We demonstrate that our proposed approach achieves promising empirical performance on both tabular and image datasets.

## 2 Background

We now briefly review the missing data imputation problem, the Generative Adversarial Imputation Nets (GAIN) [Yoon et al., 2018] and Rectified Flow [Liu et al., 2023], which are the foundations of our proposed method.

**Notations:** We denote scalars as lowercase letters (e.g.,  $x, y$ ), and vectors as bold lowercase letters (e.g.,  $\mathbf{x}, \mathbf{y}$ ). Random variables are uppercase letters, e.g.,  $X, Y$ . Random vectors are bold uppercase letters  $\mathbf{X}, \mathbf{Y}$ . The superscript  $t \in \{1, \dots, T\}$  denotes the iteration count in the algorithm. The subscript 0, 1 and  $\tau$  denote the continuous time index in an ODE. Given a fixed missing mask  $\mathbf{m}$  and a vector  $\mathbf{x}$  with the same dimension, we write  $\mathbf{x}_{\mathbf{m}}$  as the subvector of  $\mathbf{x}$  whose elements correspond to  $m_i = 1$  and  $\mathbf{x}_{1-\mathbf{m}}$  as the subvector of  $\mathbf{x}$  whose elements corresponds to  $m_i = 0$ . Following the convention, we refer to  $\mathbf{x}_{\mathbf{m}}$  as the “non-missing components” of  $\mathbf{x}$  and  $\mathbf{x}_{1-\mathbf{m}}$  as the “missing components” of  $\mathbf{x}$ .  $A \stackrel{d}{=} B$  means random variables  $A$  and  $B$  are equal in distribution.

### 2.1 Sequential Imputation

Now, we formally define the missing data imputation problem. The true data vector  $\mathbf{X}^* = [X_1^*, \dots, X_d^*]$  is a random vector taking values on  $\mathcal{X} \subset \mathbb{R}^d$ , while the missing mask  $\mathbf{M} = [M_1, \dots, M_d]$  is a random vector in  $\{0, 1\}^d$ . In the missing completely at random (MCAR) Rubin [1976] setting, we assume that the true data is independent of the missing mask, i.e.,  $\mathbf{X}^* \perp \mathbf{M}$ .

The observed data vector  $\tilde{\mathbf{X}} = (\tilde{X}_1, \dots, \tilde{X}_d)$  is defined as  $\tilde{X}_j := \begin{cases} X_j^* & M_j = 1 \\ \text{NaN} & M_j = 0 \end{cases}$ . In words, the variable  $j$  in the observation  $\tilde{\mathbf{X}}$  is missing when  $M_j = 0$  and is observed when  $M_j = 1$ . Equivalently,

$$\tilde{\mathbf{X}} = \mathbf{M} \odot \mathbf{X}^* + (1 - \mathbf{M}) \odot \text{NaN},$$

where  $\odot$  represents element-wise product. In this paper, we want to construct an imputation vector

$$\mathbf{X}(\mathbf{g}) = \mathbf{M} \odot \tilde{\mathbf{X}} + (1 - \mathbf{M}) \odot \mathbf{g}(\mathbf{Z}_0; \tilde{\mathbf{X}}, \mathbf{M}), \quad (1)$$

where  $\mathbf{g} : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is called an *imputer*.  $\mathbf{Z}_0$  is a random seed of the imputer. Different random seeds  $\mathbf{Z}_0$  lead to different imputations.

In many cases, the imputation is a sequential process: Instead of computing  $\mathbf{X}$  in one go, we improve the imputation by iterative updates. The sequential imputation considers the following procedure. For iteration  $t = 1, \dots, T$ , we construct the imputation

$$\mathbf{X}^{(t)}(\mathbf{g}_t) := \mathbf{M} \odot \mathbf{X}^{(t-1)} + (1 - \mathbf{M}) \odot \mathbf{g}_t(\mathbf{Z}_0^{(t)}; \mathbf{X}^{(t-1)}, \mathbf{M}), \quad (2)$$

where  $\mathbf{X}^{(t)}$  indicates the imputed data vector at iteration  $t$ .  $\mathbf{X}^{(0)}$  is initialized with any standard imputation technique (e.g., mean or median imputation). We show a schematic of the sequential imputation algorithm in Figure 8 of Appendix G.

The central task of this iterative imputation is to learn the imputer  $\mathbf{g}_t$  at each iteration  $t$ . In the following sections, we show how a popular missing data imputation scheme can be formulated within this framework and how  $\mathbf{g}_t$  can be trained using a GAN-type objective.

## 2.2 Generative Adversarial Imputation Nets (GAIN) [Yoon et al., 2018]

GAIN arises from a natural intuition: With perfect imputation, imputed entries are indistinguishable from the originally observed ones. Thus, GAIN trains a generator so that a classifier cannot predict the missing mask  $\mathbf{M}$  accurately.

Let  $f(m_j, \mathbf{x}, \mathbf{m}_{-j}) \in [0, 1]$  be a probabilistic classifier which models the conditional probability  $\mathbb{P}_{M_j | \mathbf{X}^{(t-1)}, \mathbf{M}_{-j}}$ , the probability of  $j$ -th component being imputed given the previous imputation  $\mathbf{X}^{(t-1)}$  and the missing mask excluding the  $j$ -th component. For iteration  $t = 1, \dots, T$ , GAIN performs the following steps:

1.  $f_{tj} := \arg \min_f -\mathbb{E} [\log f(M_j, \mathbf{X}^{(t-1)}, \mathbf{M}_{-j})]$ , for all  $j = 1 \dots d$ .
2.  $\mathbf{g}_t := \arg \max_{\mathbf{g}} -\sum_j \mathbb{E} [\log f_{tj}(M_j, \mathbf{X}^{(t)}(\mathbf{g}), \mathbf{M}_{-j})] - \lambda(\text{reconstruction error})$ .
3.  $\mathbf{X}^{(t)} \xleftarrow{\text{impute}} \mathbf{X}^{(t)}(\mathbf{g}_t), t \leftarrow t + 1$ .

$\mathbf{X}^{(t)}(\mathbf{g}_t)$  is the imputer we introduced in (2). At iteration  $t$ , we first train classifiers  $f_{tj}$  by *minimizing* the cross-entropy loss. Subsequently, we train the imputer  $\mathbf{g}_t$ , parameterized as a neural network, to reduce the classifiers' predictive accuracy by *maximizing* the sum of cross-entropy losses across all  $j$ , while simultaneously minimizing the reconstruction error.

In practice, both the minimization and maximization steps are performed by a single gradient update, leading to an adversarial training. Thus, the algorithm is named Generative Adversarial Imputation Nets (GAIN). We observe that GAIN works well in practice. However, the adversarial training is known to be difficult and may suffer from mode collapses [Li et al., 2018, Zhang et al., 2018]. Moreover, balancing the reconstruction error and the cross entropy loss requires careful tuning of the hyperparameter  $\lambda$ .

Recently, an observation was made that the generator training of GAIN may be viewed as a process of *breaking the dependency* between  $\mathbf{M}$  and  $\mathbf{X}^{(t)}(\mathbf{g})$  [Liu et al., 2024]. Indeed, if  $\mathbf{X}^{(t)}(\mathbf{g}) \perp\!\!\!\perp \mathbf{M}$ , implying  $\mathbf{X}^{(t)}(\mathbf{g})$  is not predictive of  $\mathbf{M}$  at all, then the cross-entropy loss is maximized. In this paper, we rigorously explore the idea of imputation by reducing dependency, which will guide us in developing a novel iterative imputation algorithm that reduces the dependency between  $\mathbf{M}$  and  $\mathbf{X}^{(t)}(\mathbf{g})$ .

## 2.3 Flow-based Sampling and Rectified Flow

As an alternative to GANs, flow-based generative models have attracted significant interest in recent years. These algorithms first train a “velocity field” over time  $\tau \in [0, 1]$  according to some loss function, and then generate samples by solving the ODE or SDE using the velocity field as the drift. In some cases, the loss function is a simple least-squares loss, leading to a training routine

that is simpler and more stable than GAN. Methods in this category include Score-based generative models [Song and Ermon, 2019, Song et al., 2021], Rectified Flow [Liu et al., 2023], Flow Matching [Lipman et al., 2023], Diffusion Schrödinger Bridge [De Bortoli et al., 2021], among others.

Rectified flow is one of the simplest flow-based generative models. Given two random vectors  $\mathbf{X}_0$  and  $\mathbf{X}_1$ , the velocity field  $\mathbf{v}^*$  is trained by minimizing the following objective:

$$\mathbf{v}^* = \arg \min_{\mathbf{v}} \int_0^1 \mathbb{E} [\|\mathbf{X}_1 - \mathbf{X}_0 - \mathbf{v}(\mathbf{X}_\tau, \tau)\|^2] d\tau, \quad (3)$$

where  $\mathbf{X}_\tau$  is the interpolated random vector between  $\mathbf{X}_0$  and  $\mathbf{X}_1$ , defined as  $\mathbf{X}_\tau = \tau \mathbf{X}_1 + (1 - \tau) \mathbf{X}_0$  for  $\tau \in [0, 1]$ . The sampling ODE is:  $\frac{d\mathbf{Z}_\tau}{d\tau} = \mathbf{v}^*(\mathbf{Z}_\tau, \tau)$ . One can prove that, if  $\mathbf{Z}_0 \stackrel{d}{=} \mathbf{X}_0$ , then  $\mathbf{Z}_\tau \stackrel{d}{=} \mathbf{X}_\tau$  for all  $\tau \in [0, 1]$ . This ‘‘marginal preserving property’’ implies: If one initializes the ODE using a sample  $\mathbf{X}_0$  from the reference distribution, one can obtain a sample from the target distribution by solving the ODE at  $\tau = 1$ .

### 3 Reducing Dependency

In this section, we propose the MIRI imputation framework, obtain the optimal imputer under this framework using rectified flow, and theoretically justify its validity.

#### 3.1 Mutual Information Reducing Iterations (MIRI)

Now we study the problem of training an imputer  $\mathbf{g}$  to *reduce the dependency* between the imputed sample  $\mathbf{X}(\mathbf{g})$  and  $\mathbf{M}$ . The approach adopted by [Liu et al., 2024] is to minimize the mutual information between them, i.e., the optimal  $\mathbf{g}$  is given as:

$$\mathbf{g} \in \arg \min_{\mathbf{g}} D[\mathbb{P}_{\mathbf{X}(\mathbf{g}), \mathbf{M}} | \mathbb{P}_{\mathbf{X}(\mathbf{g})} \otimes \mathbb{P}_{\mathbf{M}}], \quad (4)$$

where  $D[\mathbb{P}|\mathbb{Q}]$  stands for the Kullback-Leibler (KL) divergence of probability distributions  $\mathbb{Q}$  from  $\mathbb{P}$ . However, the mutual information is not directly computable as we do not have access to the true distribution  $\mathbb{P}_{\mathbf{X}(\mathbf{g}), \mathbf{M}}$  and  $\mathbb{P}_{\mathbf{X}(\mathbf{g})} \otimes \mathbb{P}_{\mathbf{M}}$ . A potential solution is to transform this optimization problem into a bi-level min-max adversarial optimization problem: Fixing  $\mathbf{g}$ , it is possible to estimate the mutual information using samples from  $\mathbb{P}_{\mathbf{X}(\mathbf{g}), \mathbf{M}}$  and  $\mathbb{P}_{\mathbf{X}(\mathbf{g})} \otimes \mathbb{P}_{\mathbf{M}}$  using Mutual Information Neural Estimation (MINE) [Belghazi et al., 2018] which *maximizes* the Donsker-Varadhan lowerbound [Donsker and Varadhan, 1975] to approximate the KL divergence. After that, we can *minimize* the estimated mutual information with respect to  $\mathbf{g}$ . Repeat until convergence. If one solves both optimization problems one gradient step at a time, this algorithm is the classic adversarial training that GAN is known for.

However, *contrary* to the suggestion in [Liu et al., 2024], the mutual information in (4) *cannot* be minimized simply by directly applying forward or reverse KL Wasserstein Gradient Flow (WGF). While WGF can minimize a KL divergence  $D[\mathbb{P}|\mathbb{Q}]$  when the distribution to be optimized is either  $\mathbb{P}$  or  $\mathbb{Q}$ , the distribution of interest in our setting ( $\mathbb{P}_{\mathbf{X}(\mathbf{g})}$ ) appears simultaneously in both  $\mathbb{P}$  and  $\mathbb{Q}$ . As a result, neither the forward nor reverse KL formulation of WGF is directly applicable. Therefore, while the WGF procedure proposed in [Liu et al., 2024] is conceptually appealing, it does not correctly minimize the mutual information objective as formulated in (4).

Motivated by the limitation of both adversarial training and WGF, we propose a sequential imputation algorithm that reduces the mutual information between  $\mathbf{X}^{(t)}$  and  $\mathbf{M}$ , called Mutual Information Reducing Iterations (MIRI). Let us denote the joint probability of the pair  $(\mathbf{X}^{(t)}(\mathbf{g}), \mathbf{M})$  as  $\mathbb{P}_{\mathbf{X}^{(t)}(\mathbf{g}), \mathbf{M}}$  and the product measure of  $\mathbf{X}^{(t-1)}$  and  $\mathbf{M}$  as  $\mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}}$ . For  $t = 1 \dots T$ , the algorithm performs the following steps:

1.  $\mathbf{g}_t \in \arg \min_{\mathbf{g}} D[\mathbb{P}_{\mathbf{X}^{(t)}(\mathbf{g}), \mathbf{M}} | \mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}}]$ .
2.  $\mathbf{X}^{(t)} \stackrel{\text{impute}}{\leftarrow} \mathbf{X}^{(t)}(\mathbf{g}_t), t \leftarrow t + 1$ .

The difference between the first step and (4) is that we replace  $\mathbb{P}_{\mathbf{X}(\mathbf{g})} \otimes \mathbb{P}_{\mathbf{M}}$  in (4) with  $\mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}}$ , the product measure of *the previous iteration*. We can also confirm that MIRI does indeed reduce the mutual information:

**Proposition 1.** *The mutual information between  $\mathbf{X}^{(t)}$  and  $\mathbf{M}$  is non-increasing after each iteration.*

See Appendix A for the proof. In Section 5.1, we show that the GAIN algorithm can be viewed as an approximate implementation of the above MIRI algorithm.

Moreover, we show that minimizing the KL divergence in the first step of MIRI has a sufficient and necessary condition:

**Proposition 2.** *Let  $p_{\mathbf{g}}(\mathbf{x}, \mathbf{m})$  be the density of  $(\mathbf{X}^{(t)}(\mathbf{g}_t), \mathbf{M})$  and  $q(\mathbf{x})$  be the density of  $\mathbf{X}^{(t-1)}$ .  $\mathbf{g} \in \arg \min_{\mathbf{g}} D[\mathbb{P}_{\mathbf{X}^{(t)}(\mathbf{g}), \mathbf{M}} | \mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}}]$  if and only if*

$$p_{\mathbf{g}}(\mathbf{x}_{1-\mathbf{m}} | \mathbf{x}_{\mathbf{m}}, \mathbf{m}) = q(\mathbf{x}_{1-\mathbf{m}} | \mathbf{x}_{\mathbf{m}}), \forall \mathbf{x}, \mathbf{m}. \quad (5)$$

The proof can be found in the Appendix D. This result shows that the optimal  $\mathbf{g}$  should sample from the conditional distribution of the missing components given the non-missing components according to the marginal distribution of  $\mathbf{X}^{(t-1)}$ , the imputed data in the previous iteration. (5) inspires us to construct a flow-based generative model with a terminal distribution  $q(\mathbf{x}_{1-\mathbf{m}} | \mathbf{x}_{\mathbf{m}})$ .

As we discuss in Section 5.2 and 5.3, (5) is also the key design principle behind many classic and contemporary data imputation algorithms.

### 3.2 Imputation by Rectified Flow

In this section, we show how to construct an optimal imputer for MIRI using *an imputation ODE* obtained via Rectified Flow training. Let us zoom in on Step 1 of the MIRI algorithm at iteration  $t$ .

Let  $(\mathbf{X}_1, \mathbf{M}_1)$  be a pair of random vectors drawn from the product measure  $\mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}}$  and  $(\mathbf{X}_0, \mathbf{M}_0)$  from the joint distribution  $\mathbb{P}_{\mathbf{X}^{(t-1)}, \mathbf{M}}$ . The pairs are drawn such that  $\mathbf{X}_1 \perp\!\!\!\perp (\mathbf{X}_0, \mathbf{M}_0)$ .  $\mathbf{X}_{\tau}$  is defined as the linear interpolation of  $\mathbf{X}_0$  and  $\mathbf{X}_1$ , following the same construction as in Section 2.3.

Consider the following rectified flow training objective function:

$$\mathbf{v}^* := \arg \min_{\mathbf{v}} \int_0^1 \mathbb{E} [\|\mathbf{X}_1 - \mathbf{X}_0 - \mathbf{v}(\mathbf{X}_{\tau, 1-\mathbf{M}_0}, \mathbf{X}_{0, \mathbf{M}_0}, \mathbf{X}_{1, \mathbf{M}_0}, \mathbf{M}_0, \tau)\|^2] d\tau, \quad (6)$$

where  $\mathbf{X}_{0, \mathbf{M}_0} := \mathbf{M}_0 \odot \mathbf{X}_0$  and others are defined similarly.

We can define an *imputation process* by using the optimal velocity field  $\mathbf{v}^*$ . Given a partially observed vector where the missing entries are padded with zeros, i.e.,  $[\mathbf{0}, \mathbf{x}_{\mathbf{m}}]$ , the imputation process dynamics are governed by the following ODE:

$$\frac{d\mathbf{Z}_{\tau}}{d\tau} = (1 - \mathbf{m}) \odot \mathbf{v}^*(\mathbf{Z}_{\tau}, [\mathbf{0}, \mathbf{x}_{\mathbf{m}}], [\mathbf{0}, \mathbf{x}_{\mathbf{m}}], \mathbf{m}, \tau), \quad (7)$$

with initial condition

$$\mathbf{Z}_0 \sim \mathbb{P}_{\mathbf{X}_{0, 1-\mathbf{m}} | \{\mathbf{X}_{0, \mathbf{m}} = [\mathbf{0}, \mathbf{x}_{\mathbf{m}}], \mathbf{M}_0 = \mathbf{m}\}}. \quad (8)$$

We define the imputer  $\mathbf{g}^*$  as the solution to the imputation ODE (7) evaluated at terminal time  $\tau = 1$ :

$$\mathbf{g}^*(\mathbf{Z}_0, \mathbf{x}_{\mathbf{m}}, \mathbf{m}) := \mathbf{Z}_1. \quad (9)$$

Moreover, we can show that,  $\mathbf{g}^*$  is an *optimal imputer* to the MIRI algorithm that we introduced in Section 3.1.

**Theorem 1.**  *$\mathbf{g}^*$  defined in (9) is an optimal imputer in the sense that*

$$\mathbf{g}^* \in \arg \min_{\mathbf{g}} D[\mathbb{P}_{\mathbf{X}^{(t)}(\mathbf{g}), \mathbf{M}} | \mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}}].$$

The main technical challenge of the proof is to show that the ODE (7) can indeed establish a “marginal-preserving” process from  $\mathbb{P}_{\mathbf{X}^{(t-1)} | \mathbf{M}}$  to  $\mathbb{P}_{\mathbf{X}^{(t-1)}}$ . This is different from the classic proof since both our training objective and ODE are different from the regular rectified flow. See Appendix B for the proof. With an imputer  $\mathbf{g}^*$ , we can proceed to the second step of the MIRI algorithm.

---

**Algorithm 1** MIRI with Rectified Flow (Single Imputation)

---

**Require:** Paired, i.i.d. observations and masks  $\{(\tilde{\mathbf{X}}, \mathbf{M})\}$ , Maximum Iterations  $T$ , Maximum SGD steps  $N$ , Batch Size  $B$ , Optimizer SGD\_update.

```
1:  $(\mathbf{X}^{(0)}, \mathbf{M}) \leftarrow \text{initial\_impute}(\tilde{\mathbf{X}}, \mathbf{M})$ .
2: for  $t = 1$  to  $T$  do
3:   Initialize a neural network model  $\mathbf{v}$ .
4:   for  $n = 1$  to  $N$  do
5:     Sample batch  $\{(\mathbf{X}_0^{(j)}, \mathbf{M}_0^{(j)})\}_{j=1}^B \sim \mathbb{P}_{\mathbf{X}^{(t-1)}, \mathbf{M}_0}$ ,
6:     Sample batch  $\{\mathbf{X}_1^{(j)}\}_{j=1}^B \sim \mathbb{P}_{\mathbf{X}^{(t-1)}}$ , and time indices  $\{\tau_j\}_{j=1}^B \sim \text{Uniform}(0, 1)$ 
7:      $\forall j, \mathbf{X}_{\tau_j}^{(j)} \leftarrow (1 - \tau_j)\mathbf{X}_0^{(j)} + \tau_j\mathbf{X}_1^{(j)}, \mathbf{Y}^{(j)} \leftarrow \mathbf{X}_1^{(j)} - \mathbf{X}_0^{(j)}$ 
8:      $\mathbf{v}_t \leftarrow \text{SGD\_update} \left( \nabla_{\mathbf{v}} \sum_j \left\| \mathbf{Y}^{(j)} - \mathbf{v}(\mathbf{X}_{\tau_j}^{(j)}, \mathbf{M}_0^{(j)} \odot \mathbf{X}_0^{(j)}, \mathbf{M}_0^{(j)} \odot \mathbf{X}_1^{(j)}, \mathbf{M}_0^{(j)}, \tau_j) \right\|^2 \right)$ 
9:   end for
10:   $\mathbf{X}^{(t)} \leftarrow \mathbf{M} \odot \mathbf{X}^{(t-1)} + (1 - \mathbf{M}) \odot \text{ODE\_Solver}((1 - \mathbf{M}) \odot \mathbf{v}, \mathbf{X}^{(t-1)}, \mathbf{M})$ 
11: end for
12: return  $\{\mathbf{X}_i^{(T)}\}$ 
```

---

### 3.3 Initial Condition

To run the imputation process, we require at least a sample of  $\mathbf{X}_{0,1-\mathbf{m}} | \{\mathbf{X}_{0,\mathbf{m}} = \mathbf{x}_{\mathbf{m}}, \mathbf{M} = \mathbf{m}\}$ . By the definition of  $\mathbf{X}_0$ , this is equivalent to sampling from  $\mathbf{X}_{1-\mathbf{m}}^{(t-1)} | \{\mathbf{X}_{\mathbf{m}}^{(t-1)} = \mathbf{x}_{\mathbf{m}}, \mathbf{M} = \mathbf{m}\}$ . This poses no difficulty, since at each iteration  $t$ , we maintain the access to the joint sample  $(\mathbf{X}^{(t-1)}, \mathbf{M})$ , from which we can extract the relevant conditional sample by slicing. However, in the first iteration, we need samples from  $\mathbf{X}_{1-\mathbf{m}}^{(0)} | \{\mathbf{X}_{\mathbf{m}}^{(0)} = \mathbf{x}_{\mathbf{m}}, \mathbf{M} = \mathbf{m}\}$  to kickstart the MIRI algorithm.

One can sample from any distribution as those initial samples and in our experiment, we simply draw independent samples from the normal  $N(0, 1)$  or uniform  $U(0, 1)$  to fill out each missing component. Alternatively, one can use another imputation algorithm to suggest initial samples.

## 4 Practical Implementation

Algorithm 1 presents the full MIRI procedure on a finite sample set  $(\tilde{\mathbf{X}}, \mathbf{M})$ . Note that the superscript  $(t)$  represents iteration count while  $(j)$  is the sample index in a batch. We model  $\mathbf{v}_t$  with a neural network trained via stochastic gradient descent. The function `initial_impute` replaces NaNs with initial guesses.

According to the requirement in Section 3.2, we need to ensure  $(\mathbf{X}_0, \mathbf{M}_0) \perp\!\!\!\perp \mathbf{X}_1$ . In practice, we propose to sample a batch  $\{\mathbf{X}_{0,j}, \mathbf{M}_{0,j}\}_{j=1}^B$  from the paired dataset and sample a batch  $\{\mathbf{X}_{1,j}\}_{j=1}^B$  from a *shuffled set* of  $\{\mathbf{X}_{0,j}\}$  to weaken their dependency. Although, in this case,  $\mathbf{X}_1$  is still not independent of  $\mathbf{X}_0$  and  $\mathbf{M}_0$ , we observe that this strategy works well in practice.

The imputation ODE is solved using the Euler method (see Appendix H). However, other methods may also be used. The Euler method with 100 steps generally provides good performance.

## 5 MIRI and Other Imputation Algorithms

In this section, we show that how MIRI is related to other existing imputation algorithms.

### 5.1 GAIN

Let  $q(\mathbf{x}, \mathbf{m})$  be the density function of the joint distribution of  $(\mathbf{X}^{(t-1)}, \mathbf{M})$ . Suppose the discriminator is well-trained in the first step in the GAIN algorithm, i.e.,  $f(m_j, \mathbf{x}, \mathbf{m}_{-j}) \approx q(m_j | \mathbf{x}, \mathbf{m}_{-j})$ ,

and the generator training in the second step (without the reconstruction error) is

$$\mathbf{g}_t = \arg \min_{\mathbf{g}} \mathbb{E} \left[ \sum_j \log q(M_j | \mathbf{X}_{\mathbf{g}}^{(t)}, \mathbf{M}_{-j}) \right] = \arg \min_{\mathbf{g}} \mathbb{E} \left[ \log \prod_j q(M_j | \mathbf{X}_{\mathbf{g}}^{(t)}, \mathbf{M}_{-j}) \right]$$

Using the pseudo-likelihood approximation [Besag, 1975],  $\prod_j q(M_j | \mathbf{X}_{\mathbf{g}}^{(t)}, \mathbf{M}_{-j}) \approx q(\mathbf{M} | \mathbf{X}_{\mathbf{g}}^{(t)})$ , thus:

$$\mathbf{g}_t \approx \arg \min_{\mathbf{g}} \mathbb{E} \left[ \log q(\mathbf{M} | \mathbf{X}_{\mathbf{g}}^{(t)}) \right] = \arg \min_{\mathbf{g}} \mathbb{E} \left[ \log \frac{q(\mathbf{M}, \mathbf{X}_{\mathbf{g}}^{(t)})}{q(\mathbf{X}_{\mathbf{g}}^{(t)})p(\mathbf{M})} \right]. \quad (10)$$

By Gibbs' inequality,

$$\mathbb{E} \left[ \log \frac{q(\mathbf{M}, \mathbf{X}_{\mathbf{g}}^{(t)})}{q(\mathbf{X}_{\mathbf{g}}^{(t)})p(\mathbf{M})} \right] \leq \mathbb{E} \left[ \log \frac{p_{\mathbf{g}}(\mathbf{M}, \mathbf{X}_{\mathbf{g}}^{(t)})}{q(\mathbf{X}_{\mathbf{g}}^{(t)})p(\mathbf{M})} \right] = D \left[ \mathbb{P}_{\mathbf{X}^{(t)}(\mathbf{g}), \mathbf{M}} | \mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}} \right],$$

thus, we can see that GAIN learns its imputer  $\mathbf{g}_t$  by approximately minimizing a lower bound of  $D \left[ \mathbb{P}_{\mathbf{X}^{(t)}(\mathbf{g}), \mathbf{M}} | \mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}} \right]$ .

There are two approximations made in the above argument: the pseudo-likelihood approximation and the Gibbs' inequality. First, since GAIN only performs one gradient step update to the imputer  $\mathbf{g}_t$  at each iteration, the gap between  $D \left[ \mathbb{P}_{\mathbf{X}^{(t)}(\mathbf{g}), \mathbf{M}} | \mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}} \right]$  and its lower bound may not be large. Second, the pseudo-likelihood approximation can be a good approximation, for example, when the correlation among  $M_j$  is not too strong. The discriminator training of GAIN resembles a dimension-wise strategy that has been widely used in Ising model estimation [Ravikumar et al., 2010, Meinshausen and Bühlmann, 2006]. Note that the optimization in (10) is fundamentally intractable due to the high-dimensional density  $q(\mathbf{m} | \mathbf{x}^{(t)}(\mathbf{g}))$ .  $\mathbf{m}$  is a binary variable defined on  $\{0, 1\}^d$ , thus, the density not have a tractable normalizing constant. This further restricts us to the dimension-wise pseudo-likelihood approach in GAIN algorithm.

## 5.2 Round-Robin Approaches: MICE, HyperImpute, etc.

Recall that Proposition 2 states that minimizing the KL divergence  $D \left[ \mathbb{P}_{\mathbf{X}^{(t)}(\mathbf{g}), \mathbf{M}} | \mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}} \right]$  is equivalent to choosing a  $\mathbf{g}$  such that  $p_{\mathbf{g}}(\mathbf{x}_{1-\mathbf{m}} | \mathbf{x}_{\mathbf{m}}, \mathbf{m}) = q(\mathbf{x}_{1-\mathbf{m}} | \mathbf{x}_{\mathbf{m}})$ , where  $p_{\mathbf{g}}(\mathbf{x}, \mathbf{m})$  is the density of  $(\mathbf{X}^{(t)}(\mathbf{g}), \mathbf{M})$  and  $q(\mathbf{x})$  is the density of  $\mathbf{X}^{(t-1)}$ . One can simply build an imputer that samples from  $q(\mathbf{x}_{1-\mathbf{m}} | \mathbf{x}_{\mathbf{m}})$  to impute the missing components. Then, by construction,  $p_{\mathbf{g}}(\mathbf{x}_{1-\mathbf{m}} | \mathbf{x}_{\mathbf{m}}, \mathbf{m}) = q(\mathbf{x}_{1-\mathbf{m}} | \mathbf{x}_{\mathbf{m}})$ .

This algorithm has two issues: First,  $\mathbf{m}$  changes for every sample, so we need to learn a different imputer for every sample. Second, sampling from the conditional probability distribution  $q(\mathbf{x}_{1-\mathbf{m}} | \mathbf{x}_{\mathbf{m}})$  is not trivial as the dimension of  $\mathbf{x}_{1-\mathbf{m}}$  can still be high.

In principle, we can learn the *joint probability*  $q$  then sample from the conditional distributions  $q(\mathbf{x}_{1-\mathbf{m}} | \mathbf{x}_{\mathbf{m}})$  for any fixed  $\mathbf{m}$  (see the next section). However, learning a high dimensional joint probability density function  $q$  is not trivial. Thus, we adopt the pseudo-likelihood approximation again: factorize  $q(\mathbf{x})$  over dimension-wise conditional probabilities and learn  $q(x_j | \mathbf{x}_{-j})$  for every  $j$ . This dimension-wise conditional probability learning scheme not only makes the learning of the joint distribution easier but also solves the second issue as we can perform *Gibbs sampling*, one missing component at a time. This pseudo-likelihood approximation + Gibbs sampling strategy gives rise to the round-robin algorithms such as MICE and HyperImpute.

Let  $f_j$  be the density model of  $q(x_j | \mathbf{x}_{-j})$ , then the "MIRI MICE" performs the following algorithm. For  $t = 1, \dots, T$ :

1.  $\forall j$ , learn  $f_j$  to approximate  $q(x_j | \mathbf{x}_{-j})$ .
2. Let  $\mathbf{g}_t$  be a Gibbs sampler for  $\mathbb{P}_{\mathbf{X}_{1-\mathbf{m}}^{(t-1)} | \mathbf{X}_{\mathbf{m}}^{(t-1)} = \mathbf{x}_{\mathbf{m}}}$  using  $\{f_j\}$ .
3.  $\mathbf{X}^{(t)} \stackrel{\text{impute}}{\leftarrow} \mathbf{X}^{(t)}(\mathbf{g}_t), t \leftarrow t + 1$ .

In practice, MICE interlaces first step and the later two steps, meaning a Gibbs sampling step and imputation is immediately performed after learning each  $f_j$ . However, the sequential nature of Gibbs sampling prevent us from parallelizing this procedure for high-dimensional samples.

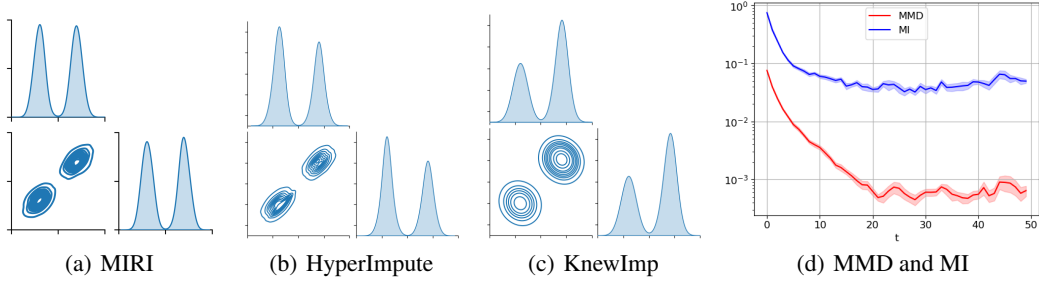


Figure 1: Pairwise density plots of the imputed data and the true data and MMD/MI of MIRI.

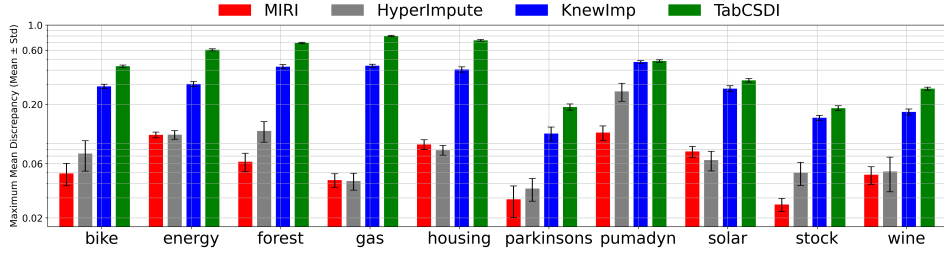


Figure 2: MMD on UCI datasets with 60% data missing. The lower the better.

### 5.3 Contemporary Approaches: Learn Jointly, Sample/Inpaint Conditionally

Traditionally, learning a high-dimensional joint distribution  $q(\mathbf{x})$  is intractable due to its normalising constant. Recently, it was realized that one can learn a joint model  $q(\mathbf{x})$  without dealing with the normalising constant. [Uehara et al., 2020] first estimates the joint model  $q$  via Score Matching [Hyvärinen, 2005] or Noise Contrastive Estimation [Gutmann and Hyvärinen, 2010], then sample from  $q(\mathbf{x}_{1-m}|\mathbf{x}_m)$  via importance sampling. This method works well when  $\mathbf{x}_{1-m}$  is in a relatively low-dimensional space. [Chen et al., 2024] noticed that given a joint model  $q$ , this conditional sampling problem could be solved using Stein Variational Gradient Descent (SVGD) [Liu and Wang, 2016]. However, SVGD utilizes properties of RKHS function, and require kernel computations. Therefore, does not scale to high-dimensional datasets.

In recent years, it is also realized that one can “inpaint” samples given a diffusion model trained on the joint samples [Lugmayr et al., 2022]. This idea gives rise to DiffPuter [Zhang et al., 2025], an imputer that trains a joint diffusion model to sample from  $q$ . Then use an inpainter to impute missing values given the observed vector  $\mathbf{x}_m$ . However, an inpainter is not a standard backward process of diffusion dynamics, thus in general, is not guaranteed to generate samples from  $q(\mathbf{x}_{1-m}|\mathbf{x}_m)$ .

## 6 Experiments

Details on experimental settings and implementations can be found in Appendix E. Additional results can be found in Appendix F. In experiments, we match model architectures across all methods, if possible. For example, if our method uses CNN, we use the same architectures for other neural-network-based methods as well.

### 6.1 Toy Example: Assessing Imputation Quality Using MMD and Mutual Information

In this section, we showcase the performance of the proposed MIRI imputer on a toy example that is a mixture of 2-dimensional Gaussian:  $\mathbf{X}^* \sim \mathcal{N}([-2, -2], 0.5^2 \mathbf{I}) + \mathcal{N}([2, 2], 0.5^2 \mathbf{I})$ . The missing mask is generated as  $\mathbf{M} \sim \prod_{j=1}^d \text{Bernoulli}_{M_j}(0.7)$ , i.e., 30% of the data is missing completely at random. We draw 6000 i.i.d. samples from  $\mathbb{P}_{\mathbf{X}^*}$  and  $\mathbb{P}_{\mathbf{M}}$  to create missing observations. The imputed samples are visualized in the pairwise density plots in Figure 1. We choose HyperImpute [Jarrett et al., 2022] and KnewImp [Chen et al., 2024] as comparisons. It can be seen that MIRI



Table 1: **Quantitative results on CIFAR-10.** Methods are evaluated at three levels of missingness (20%, 40%, 60%) using FID, PSNR, and SSIM. The best results are highlighted in bold.

Method	20% Missingness			40% Missingness			60% Missingness		
	FID ( $\downarrow$ )	PSNR ( $\uparrow$ )	SSIM ( $\uparrow$ )	FID ( $\downarrow$ )	PSNR ( $\uparrow$ )	SSIM ( $\uparrow$ )	FID ( $\downarrow$ )	PSNR ( $\uparrow$ )	SSIM ( $\uparrow$ )
GAIN [Yoon et al., 2018]	164.11	21.21	0.7803	281.62	16.20	0.5576	285.53	11.99	0.2933
KnewImp [Chen et al., 2024]	153.09	18.84	0.6463	193.68	15.81	0.4740	264.40	14.04	0.3317
MissDiff [Ouyang et al., 2023]	90.51	22.29	0.7702	129.84	19.65	0.6648	197.91	16.78	0.4989
HyperImpute [Jarrett et al., 2022]	8.92	<b>34.09</b>	<b>0.9750</b>	65.01	23.22	0.7931	130.36	20.17	0.6533
<b>MIRI (Ours*)</b>	<b>6.01</b>	32.29	0.9736	<b>27.53</b>	<b>27.14</b>	<b>0.9126</b>	<b>68.58</b>	<b>23.22</b>	<b>0.8063</b>

not only captures the bimodal structure of the true data but also preserves the proportion of the two modes. On the other hand, HyperImpute and KnewImp fail to recover the correct proportion of the two modes. We also show performance metrics (MAE, RMSE, MMD) with various other imputation algorithms (such as MissDiff, GAIN, TabCSDI [Zheng and Charoenphakdee, 2022]) in Table 3 in the Appendix F.1. Finally, in Figure 1(d), we show that both MIRI’s MMD and MI decrease over iterations, proving that MI between the imputed data and the missing mask is a good criterion for assessing how well the imputed data recovers  $\mathbb{P}_{\mathbf{X}^*}$ .

## 6.2 Real-World Tabular Imputation: UCI Regression Benchmarks

We evaluate MIRI on 10 standard UCI regression benchmarks [Kelly et al., 2017]. We select these datasets for their diversity in sample size and dimensionality (see Table 2 in Appendix E.1). We simulate missingness by independently masking 20%, 40%, and 60% of entries under an MCAR mechanism. Our method is compared with HyperImpute, KnewImp and TabCSDI. Figure 2 reports the average MMD over 10 trials (error bars denote one standard deviation). As shown, even in the 60% missing data cases, MIRI outperforms or performs comparably to all competitors. Results for the 20% and 40% missingness scenarios, as well as additional comparisons, are provided in Appendix F.2.

## 6.3 Real-World Image Imputation: CIFAR-10 and CelebA Benchmarks

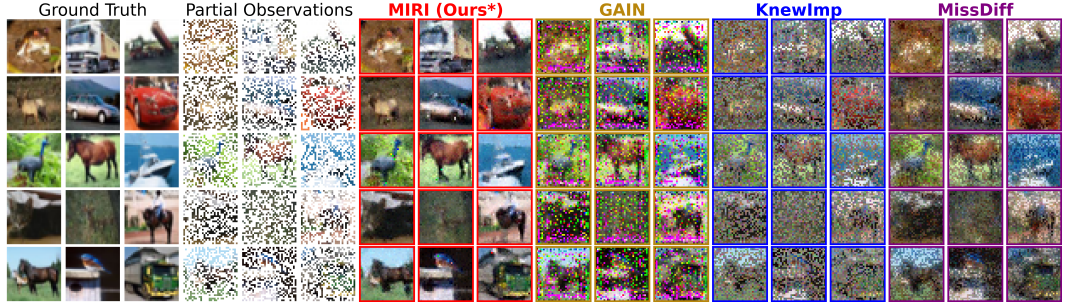
We first assess MIRI on the high-dimensional CIFAR-10 dataset [Krizhevsky et al., 2009] ( $32 \times 32$  RGB images) by randomly removing 60% of pixels (all RGB channels at once), and training on just 5 000 samples ( $< 10\%$  of the full set). Figure 3(a) shows the imputation reconstructions at 60% missingness of the generative model-based approaches GAIN, KnewImp, and MissDiff, alongside our method. Table 1 reports FID, PSNR and SSIM at 20%, 40% and 60% missingness. Across this range, MIRI achieves 65.4 – 93.4% lower FID, 38.1 – 44.9% higher PSNR and 26.4 – 61.6% higher SSIM than the strongest flow-based alternative. These margins remain significant even at 60% missingness, demonstrating robustness to severe sparsity. We have also included results for HyperImpute in Table 1, but in the main text, we only show visual reconstructions for generative model-based approaches. Visual reconstructions for HyperImpute, along with additional results, are provided in Appendix F.3.

Additionally, we evaluate the performance of MIRI on a higher-resolution dataset on CelebA [Liu et al., 2015] ( $64 \times 64$  RGB images), where we adopt a complementary missingness mechanism that *independently* masks each RGB channel. The imputation reconstruction of images using 5 000 samples displayed in Figure 3(b) substantiates MIRI’s strong performance.

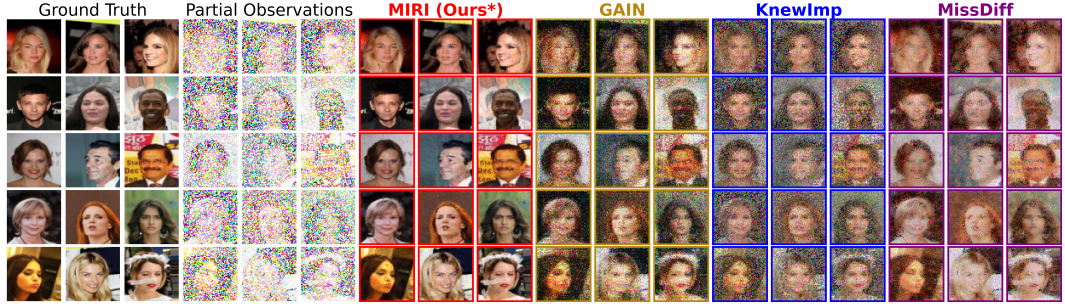
## 7 Limitations and Future Work

Our method is computationally heavy, as for each MIRI iteration, we need to run a full rectified flow algorithm, including training the velocity field and solving the ODE. One interesting future work is to incorporate recent distillation methods to accelerate sample generation.

Our Theorem 1 is a population-level statement, meaning it assumes the population optimal  $\mathbf{v}^*$  is obtained from the training objective (6) and that the imputation ODE (7) is solved exactly. It does not consider errors due to finite sample approximation, optimisation and inaccurate ODE solver.



(a) 15 uncured 32×32 CIFAR-10 images and their imputations. Pixels are removed from *all RGB channels*.



(b) 15 uncured 64×64 CelebA images and their imputations. Pixels are removed from *each RGB channel independently*.

Figure 3: Comparison of imputed samples on CIFAR-10 and CelebA with 60% of missingness.

How to extend our theoretical results to the finite-sample optimal  $\hat{v}$  and the approximate imputation ODE solver (like the Euler method) is another avenue for future works.

## Acknowledgements

Funding for this work was provided by the University of Bristol Undergraduate Summer Research Bursary 2024. Jiahao Yu was funded by the Alumni Foundation, and Leyang Wang by the Heilbronn Institute for Mathematical Research. Partial computational resources were provided by the Advanced Computing Research Centre, University of Bristol (<http://www.bristol.ac.uk/acrc/>).

## References

- M. I. Belghazi, A. Baratin, S. Rajeshwar, S. Ozair, Y. Bengio, A. Courville, and D. Hjelm. Mutual information neural estimation. In *Proceedings of International Conference on Machine Learning*, 2018.
- J. Besag. Statistical analysis of non-lattice data. *The Statistician*, pages 179–195, 1975.
- Alberto Brini, , and Edwin R. van den Heuvel. Missing data imputation with high-dimensional data. *The American Statistician*, 78(2):240–252, 2024. ISSN 0003-1305. doi: 10.1080/00031305.2023.2259962. URL <https://doi.org/10.1080/00031305.2023.2259962>.
- Zhichao Chen, Haoxuan Li, Fangyikang Wang, Odin Zhang, Hu Xu, Xiaoyu Jiang, Zhihuan Song, and Hao Wang. Rethinking the diffusion models for missing data imputation: A gradient flow perspective. In *Advances in Neural Information Processing Systems*, 2024.
- V. De Bortoli, J. Thornton, J. Heng, and A. Doucet. Diffusion Schrödinger bridge with applications to score-based generative modeling. In *Advances in Neural Information Processing Systems*, 2021.

- Laurent Dinh, David Krueger, and Yoshua Bengio. NICE: Non-linear independent components estimation. In *3rd International Conference on Learning Representations, ICLR 2015, Workshop Track Proceedings*, 2015. URL <http://arxiv.org/abs/1410.8516>.
- M. D. Donsker and S. R. S. Varadhan. Asymptotic evaluation of certain markov process expectations for large time, i. *Communications on Pure and Applied Mathematics*, 28(1):1–47, 1975.
- I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, 2014.
- Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of International Conference on Artificial Intelligence and Statistics*, 2010.
- A. Hyvärinen. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6:695–709, 2005.
- D. Jarrett, B. Cebere, T. Liu, A. Curth, and M. van der Schaar. HyperImpute: Generalized iterative imputation with automatic model selection. In *Proceedings of International Conference on Machine Learning*, 2022.
- M. Kelly, R. Longjohn, and K. Nottingham. UCI machine learning repository. Available at: <https://archive.ics.uci.edu/ml>, 2017. (Accessed: 16 May 2025).
- Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. CIFAR-10 and CIFAR-100 datasets. Available at: <https://www.cs.toronto.edu/~kriz/cifar.html>, 2009. (Accessed: 16 May 2025).
- H. S. Laqueur, A. B. Shev, and R. M. C. Kagawa. SuperMICE: An ensemble machine learning approach to multiple imputation by chained equations. *American Journal of Epidemiology*, 191(3):516–525, March 2022.
- J. Li, A. Madry, J. Peebles, and L. Schmidt. On the limitations of first-order approximation in GAN dynamics. In *Proceedings of International Conference on Machine Learning*, 2018.
- S. C.-X. Li, B. Jiang, and B. Marlin. MisGan: Learning from incomplete data with generative adversarial networks. In *International Conference on Learning Representations*, 2019.
- Y. Lipman, R. T. Q. Chen, H. Ben-Hamu, M. Nickel, and M. Le. Flow matching for generative modeling. In *International Conference on Learning Representations*, 2023.
- R. J. A. Little and D. B. Rubin. *Statistical analysis with missing data*. John Wiley & Sons, 2019.
- Q. Liu and D. Wang. Stein variational gradient descent: A general purpose bayesian inference algorithm. In *Advances in Neural Information Processing Systems*, 2016.
- Song Liu, Jiahao Yu, Jack Simons, Mingxuan Yi, and Mark Beaumont. Minimizing  $f$ -divergences by interpolating velocity fields. In *Proceedings of International Conference on Machine Learning*, 2024.
- X. Liu, C. Gong, and Q. Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *International Conference on Learning Representations*, 2023.
- Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2015.
- Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. RePaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- N. Meinshausen and P. Bühlmann. High-dimensional graphs and variable selection with the Lasso. *Annals of Statistics*, 34(3):1436–1462, June 2006.

- 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.
- P. Ravikumar, M. J. Wainwright, and J. D. Lafferty. High-dimensional Ising model selection using  $\ell_1$ -regularized logistic regression. *Annals of Statistics*, 38(3):1287–1319, 2010.
- D. J. Rezende and S. Mohamed. Variational inference with normalizing flows. In *Proceedings of International Conference on Machine Learning*, 2015.
- Donald B Rubin. Inference and missing data. *Biometrika*, 63(3):581–592, 1976.
- Y. Song and S. Ermon. Generative modeling by estimating gradients of the data distribution. In *Advances in Neural Information Processing Systems*, 2019.
- Y. Song, J. Sohl-Dickstein, D.P. Kingma, Kumar A., S. Ermon, and Poole B. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021.
- D. J. Stekhoven and P. Bühlmann. MissForest—non-parametric missing value imputation for mixed-type data. *Bioinformatics*, 28(1):112–118, 2012.
- Masatoshi Uehara, Takeru Matsuda, and Jae Kwang Kim. Imputation estimators for unnormalized models with missing data. In *Proceedings of International Conference on Artificial Intelligence and Statistics*, 2020.
- K. I. Vaden, M. Gebregziabher, S. E. Kuchinsky, and M. A. Eckert. Multiple imputation of missing fMRI data in whole brain analysis. *Neuroimage*, 60(3):1843–1855, 2012.
- S. van Buuren and K. Groothuis-Oudshoorn. MICE: Multivariate imputation by chained equations in R. *Journal of Statistical Software*, 45(3):1–67, 2011.
- J. Yoon, J. Jordon, and M. Schaar. GAIN: Missing data imputation using generative adversarial nets. In *Proceedings of International Conference on Machine Learning*, 2018.
- Hengrui Zhang, Liancheng Fang, Qitian Wu, and Philip S. Yu. DiffPuter: An EM-driven diffusion model for missing data imputation. In *International Conference on Learning Representations*, 2025.
- Zhaoyu Zhang, Mengyan Li, and Jun Yu. On the convergence and mode collapse of GAN. In *SIG-GRAPH Asia 2018 Technical Briefs*, SA ’18, pages 1–4, New York, NY, USA, 2018. Association for Computing Machinery. ISBN 978-1-4503-6062-3. doi: 10.1145/3283254.3283282. URL <https://doi.org/10.1145/3283254.3283282>.
- Shuhan Zheng and Nontawat Charoenphakdee. Diffusion models for missing value imputation in tabular data. In *NeurIPS Table Representation Learning (TRL) Workshop*, 2022.

## A Proof of Proposition 1

*Proof.* Consider the difference between the mutual information at iteration  $t$  and  $t - 1$ :

$$\begin{aligned} & D \left[ \mathbb{P}_{\mathbf{X}, \mathbf{M}}^{(t)} | \mathbb{P}_{\mathbf{X}}^{(t)} \otimes \mathbb{P}_{\mathbf{M}} \right] - D \left[ \mathbb{P}_{\mathbf{X}^{(t-1)}, \mathbf{M}} | \mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}} \right] \\ = & \underbrace{D \left[ \mathbb{P}_{\mathbf{X}, \mathbf{M}}^{(t)} | \mathbb{P}_{\mathbf{X}}^{(t)} \otimes \mathbb{P}_{\mathbf{M}} \right] - D \left[ \mathbb{P}_{\mathbf{X}, \mathbf{M}}^{(t)} | \mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}} \right]}_A + \underbrace{D \left[ \mathbb{P}_{\mathbf{X}, \mathbf{M}}^{(t)} | \mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}} \right] - D \left[ \mathbb{P}_{\mathbf{X}^{(t-1)}, \mathbf{M}} | \mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}} \right]}_B. \end{aligned} \quad (11)$$

Since  $D \left[ \mathbb{P}_{\mathbf{X}^{(t)}(\mathbf{g}), \mathbf{M}} | \mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}} \right]$  is minimized at  $\mathbf{g}_t$  and  $\mathbf{X}^{(t)} = \mathbf{X}^{(t)}(\mathbf{g}_t)$  if  $\mathbf{g}_t$  is optimal,

$$D \left[ \mathbb{P}_{\mathbf{X}, \mathbf{M}}^{(t)} | \mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}} \right] \leq D \left[ \mathbb{P}_{\mathbf{X}^{(t)}(\mathbf{g}), \mathbf{M}} | \mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}} \right], \forall \mathbf{g}.$$

Moreover,  $D \left[ \mathbb{P}_{\mathbf{X}^{(t-1)}, \mathbf{M}} | \mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}} \right] = D \left[ \mathbb{P}_{\mathbf{X}^{(t)}(\mathbf{g}), \mathbf{M}} | \mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}} \right] |_{\mathbf{g}=\text{Identity}}$ , so  $B$  is non-positive. To show  $A$  is non-positive, we write out the KL-divergences in  $A$ . Let  $p(\mathbf{x}, \mathbf{m})$  be the density of  $(\mathbf{X}^{(t)}, \mathbf{M})$  and  $q(\mathbf{x})$  be the density of  $\mathbf{X}^{(t-1)}$ .

$$\begin{aligned} & D \left[ \mathbb{P}_{\mathbf{X}, \mathbf{M}}^{(t)} | \mathbb{P}_{\mathbf{X}}^{(t)} \otimes \mathbb{P}_{\mathbf{M}} \right] - D \left[ \mathbb{P}_{\mathbf{X}, \mathbf{M}}^{(t)} | \mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}} \right] \\ = & \mathbb{E} \left[ \log \frac{p(\mathbf{X}^{(t)}, \mathbf{M})}{p(\mathbf{X}^{(t)})p(\mathbf{M})} \right] - \mathbb{E} \left[ \log \frac{p(\mathbf{X}^{(t)}, \mathbf{M})}{q(\mathbf{X}^{(t)})p(\mathbf{M})} \right] \\ = & -\mathbb{E} \left[ \log p(\mathbf{X}^{(t)}) \right] + \mathbb{E} \left[ \log q(\mathbf{X}^{(t)}) \right] \leq 0, \end{aligned}$$

where the last inequality is due to the Gibbs inequality.

Thus,  $A + B \leq 0$  as desired.  $\square$

## B Proof of Theorem 1

*Proof.* First, we introduce the marginal preserving property for the imputation ODE.

**Lemma 1.** *The imputation ODE with the initial condition given in (8), i.e.,*

$$\mathbf{Z}_0 \stackrel{d}{=} \mathbf{X}_{0,1-\mathbf{m}} | \{ \mathbf{X}_{0,\mathbf{m}} = \mathbf{x}_{\mathbf{m}}, \mathbf{M}_0 = \mathbf{m} \}$$

*has a solution at the terminal time  $\tau = 1$ ,  $\mathbf{Z}_1 \stackrel{d}{=} \mathbf{X}_{1,1-\mathbf{m}} | \{ \mathbf{X}_{1,\mathbf{m}} = \mathbf{x}_{\mathbf{m}} \}$ , for every fixed  $\mathbf{x}, \mathbf{m}$ .*

See the Appendix C for the proof.

Now we derive the main result. Lemma 1 shows that under the initial condition,

$$\mathbf{Z}_1 \stackrel{d}{=} \mathbf{X}_{1,1-\mathbf{m}} | \mathbf{X}_{1,\mathbf{m}} = \mathbf{x}_{\mathbf{m}}. \quad (12)$$

$\mathbf{Z}_1$  is our imputer  $\mathbf{g}_t$  at iteration  $t$ . According to (7), the imputer is constructed as  $\mathbf{g}_t(\cdot; \mathbf{X}_{\mathbf{m}}^{(t-1)}, \mathbf{M})$ , thus

$$\mathbf{Z}_1 \stackrel{d}{=} \mathbf{X}_{\mathbf{g}^*, \mathbf{m}}^{(t)} | \left\{ \mathbf{X}_{1-\mathbf{m}}^{(t-1)} = \mathbf{x}_{1-\mathbf{m}}, \mathbf{M}_0 = \mathbf{m} \right\}. \quad (13)$$

Combining (13) and (12), we obtain the following equality

$$\mathbb{P}_{\mathbf{X}_{\mathbf{g}^*, 1-\mathbf{m}}^{(t)} | \mathbf{X}_{\mathbf{m}}^{(t-1)} = \mathbf{x}_{\mathbf{m}}, \mathbf{M} = \mathbf{m}} = \mathbb{P}_{\mathbf{X}_{1-\mathbf{m}}^{(t-1)} | \mathbf{X}_{\mathbf{m}}^{(t-1)} = \mathbf{x}_{\mathbf{m}}}.$$

Our imputation ODE does not change components where  $m_j = 1$ , so  $\mathbf{X}_{\mathbf{m}}^{(t-1)} = \mathbf{X}_{\mathbf{m}}^{(t)}$ . Thus,

$$\mathbb{P}_{\mathbf{X}_{\mathbf{g}^*, 1-\mathbf{m}}^{(t)} | \mathbf{X}_{\mathbf{m}}^{(t-1)} = \mathbf{x}_{\mathbf{m}}, \mathbf{M} = \mathbf{m}} = \mathbb{P}_{\mathbf{X}_{\mathbf{g}^*, 1-\mathbf{m}}^{(t)} | \mathbf{X}_{\mathbf{m}}^{(t)} = \mathbf{x}_{\mathbf{m}}, \mathbf{M} = \mathbf{m}}.$$

The above two equalities imply that  $p_{\mathbf{g}^*}(\cdot | \mathbf{x}_{\mathbf{m}}, \mathbf{m}) = q(\cdot | \mathbf{x}_{\mathbf{m}})$ . According to Proposition 2, we have shown that  $\mathbf{g}^*$  is an optimal imputer that minimizes  $D \left[ \mathbb{P}_{\mathbf{X}^{(t)}(\mathbf{g}), \mathbf{M}} | \mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}} \right]$ .  $\square$

## C Proof of Lemma 1

**Definition 1.** For fixed values  $\mathbf{y}, \mathbf{u}, \mathbf{v}, \mathbf{m}$ , the event  $\Xi$  is defined as

$$\Xi(\mathbf{u}, \mathbf{w}, \mathbf{m}) = \{\mathbf{X}_{0,\mathbf{m}} = \mathbf{u}, \mathbf{X}_{1,\mathbf{m}} = \mathbf{w}, \mathbf{M}_0 = \mathbf{m}\}.$$

and the optimal solution of (6) is

$$\mathbf{v}_{\mathbf{m}}^*(\mathbf{y}, \mathbf{u}, \mathbf{w}, \mathbf{m}, \tau) := \mathbb{E}[\dot{\mathbf{X}}_{\tau,1-\mathbf{m}} \mid \mathbf{X}_{\tau,1-\mathbf{m}} = \mathbf{y}, \Xi(\mathbf{u}, \mathbf{w}, \mathbf{m})].$$

We denote the density of random variable  $\mathbf{X}_{\tau,1-\mathbf{m}} \mid \Xi(\mathbf{u}, \mathbf{w}, \mathbf{m})$  as  $p_\tau(\mathbf{y} \mid \mathbf{u}, \mathbf{w}, \mathbf{m})$ .

**Lemma 2.**  $p_0(\mathbf{y} \mid \mathbf{u}, \mathbf{w}, \mathbf{m}) = p_0(\mathbf{y} \mid \mathbf{u}, \mathbf{m})$  and  $p_1(\mathbf{y} \mid \mathbf{u}, \mathbf{w}, \mathbf{m}) = p_1(\mathbf{y} \mid \mathbf{w})$ .

*Proof.* We can see that due to the independence  $\mathbf{X}_1 \perp \mathbf{X}_0, \mathbf{M}_0$ ,

$$\mathbf{X}_{0,1-\mathbf{m}} \mid \mathbf{X}_{0,\mathbf{m}}, \mathbf{X}_{1,\mathbf{m}}, \mathbf{M}_0 = \mathbf{X}_{0,1-\mathbf{m}} \mid \mathbf{X}_{0,\mathbf{m}}, \mathbf{M}_0,$$

hence the first equality holds. What's more,

$$\mathbf{X}_{1,1-\mathbf{m}} \mid \mathbf{X}_{0,\mathbf{m}}, \mathbf{X}_{1,\mathbf{m}}, \mathbf{M}_0 = \mathbf{X}_{1,1-\mathbf{m}} \mid \mathbf{X}_{1,\mathbf{m}},$$

hence the second equality holds.  $\square$

**Lemma 3.** The density function  $p_\tau(\mathbf{y} \mid \mathbf{u}, \mathbf{w}, \mathbf{m})$  satisfies the continuity equation

$$\partial_\tau p_\tau(\mathbf{y} \mid \mathbf{u}, \mathbf{w}, \mathbf{m}) + \nabla_{\mathbf{y}} \cdot (p_\tau(\mathbf{y} \mid \mathbf{u}, \mathbf{w}, \mathbf{m}) \mathbf{v}^*(\mathbf{y}, \mathbf{u}, \mathbf{w}, \mathbf{m}, \tau)) = 0 \quad (14)$$

*Proof.* Consider the following expectation:

$$\mathbb{E}[h(\mathbf{X}_{\tau,1-\mathbf{m}}) \mid \Xi(\mathbf{u}, \mathbf{w}, \mathbf{m})] = \int_{\mathbb{R}^{d_{\mathbf{m}}}} h(\mathbf{y}) p_\tau(\mathbf{y} \mid \mathbf{u}, \mathbf{w}, \mathbf{m}) d\mathbf{y} \quad (15)$$

whose time derivative is

$$\partial_\tau \mathbb{E}[h(\mathbf{X}_{\tau,1-\mathbf{m}}) \mid \Xi(\mathbf{u}, \mathbf{w}, \mathbf{m})] = \int h(\mathbf{y}) \partial_\tau p_\tau(\mathbf{y} \mid \mathbf{u}, \mathbf{w}, \mathbf{m}) d\mathbf{y}. \quad (16)$$

On the other hand, we can see that  $\mathbf{X}_\tau$  is a function of  $\mathbf{X}_0$  and  $\mathbf{X}_1$ , thus using the law of unconscious statistician:

$$\begin{aligned} \partial_\tau \mathbb{E}[h(\mathbf{X}_{\tau,1-\mathbf{m}}) \mid \Xi(\mathbf{u}, \mathbf{w}, \mathbf{m})] &= \partial_\tau \mathbb{E}_{\mathbf{X}_0, \mathbf{X}_1} [h(\mathbf{X}_{\tau,1-\mathbf{m}}) \mid \Xi(\mathbf{u}, \mathbf{w}, \mathbf{m})] \\ &= \mathbb{E}_{\mathbf{X}_0, \mathbf{X}_1} [\partial_\tau h(\mathbf{X}_{\tau,1-\mathbf{m}}) \mid \Xi(\mathbf{u}, \mathbf{w}, \mathbf{m})] = \mathbb{E}[\partial_\tau h(\mathbf{X}_{\tau,1-\mathbf{m}}) \mid \Xi(\mathbf{u}, \mathbf{w}, \mathbf{m})]. \end{aligned}$$

Expanding the previous expression using the chain rule,

$$\partial_\tau \mathbb{E}[h(\mathbf{X}_{\tau,1-\mathbf{m}}) \mid \Xi(\mathbf{u}, \mathbf{w}, \mathbf{m})] = \mathbb{E}[\partial_\tau h(\mathbf{X}_{\tau,1-\mathbf{m}}) \mid \Xi(\mathbf{u}, \mathbf{w}, \mathbf{m})] \quad (17)$$

$$= \mathbb{E}[\nabla h(\mathbf{X}_{\tau,1-\mathbf{m}})^\top \dot{\mathbf{X}}_{\tau,1-\mathbf{m}} \mid \Xi(\mathbf{u}, \mathbf{w}, \mathbf{m})] \quad (18)$$

$$= \mathbb{E}[\mathbb{E}[\nabla h(\mathbf{X}_{\tau,1-\mathbf{m}})^\top \dot{\mathbf{X}}_{\tau,1-\mathbf{m}} \mid \mathbf{X}_{\tau,1-\mathbf{m}}; \Xi(\mathbf{u}, \mathbf{w}, \mathbf{m})] \mid \Xi(\mathbf{u}, \mathbf{w}, \mathbf{m})] \quad (19)$$

$$= \mathbb{E}[\nabla h(\mathbf{X}_{\tau,1-\mathbf{m}})^\top \mathbb{E}[\dot{\mathbf{X}}_{\tau,1-\mathbf{m}} \mid \mathbf{X}_{\tau,1-\mathbf{m}}; \Xi(\mathbf{u}, \mathbf{w}, \mathbf{m})] \mid \Xi(\mathbf{u}, \mathbf{w}, \mathbf{m})] \quad (20)$$

$$= \int \nabla h(\mathbf{y})^\top \mathbf{v}^*(\mathbf{y}, \mathbf{u}, \mathbf{w}, \mathbf{m}, \tau) p_\tau(\mathbf{y} \mid \mathbf{u}, \mathbf{w}, \mathbf{m}) d\mathbf{y} \quad (21)$$

$$= - \int h(\mathbf{y}) \nabla \cdot (\mathbf{v}^*(\mathbf{y}, \mathbf{u}, \mathbf{w}, \mathbf{m}, \tau) p_\tau(\mathbf{y} \mid \mathbf{u}, \mathbf{w}, \mathbf{m})) d\mathbf{y} \quad (22)$$

where the last equality by performing integration by parts on the right-hand side (with respect to  $\mathbf{y}$ ).

Now we equate the two representations ((16) and (22)) for  $\partial_\tau \mathbb{E}[h(\mathbf{X}_{\tau,1-\mathbf{m}}) \mid \Xi(\mathbf{u}, \mathbf{w}, \mathbf{m})]$ :

$$\int h(\mathbf{y}) \partial_\tau p_\tau(\mathbf{y} \mid \mathbf{u}, \mathbf{w}, \mathbf{m}) d\mathbf{y} = - \int h(\mathbf{y}) \nabla \cdot (\mathbf{v}^*(\mathbf{y}, \mathbf{u}, \mathbf{w}, \mathbf{m}, \tau) p_\tau(\mathbf{y} \mid \mathbf{u}, \mathbf{w}, \mathbf{m})) d\mathbf{y}. \quad (23)$$

Since this equality holds for every smooth, compactly supported test function  $h(\mathbf{y})$ , the fundamental lemma of the calculus of variations implies that, in the sense of distributions,

$$\partial_\tau p_\tau(\mathbf{y}|\mathbf{u}, \mathbf{w}, \mathbf{m}) + \nabla \cdot (\mathbf{v}^*(\mathbf{y}, \mathbf{u}, \mathbf{w}, \mathbf{m}, \tau) p_\tau(\mathbf{y}|\mathbf{u}, \mathbf{w}, \mathbf{m})) = 0. \quad (24)$$

We can see that  $\mathbf{v}^*$  satisfies the continuity equation for a time-varying density function  $p_\tau$ .  $\square$

We can see that with the initial condition  $\mathbf{Z}_0 = \mathbf{X}_{0,1-\mathbf{m}} \mid \Xi(\mathbf{u}, \mathbf{w}, \mathbf{m})$ , Lemma 3 says  $\mathbf{Z}_1 = \mathbf{X}_{1,1-\mathbf{m}} \mid \Xi(\mathbf{u}, \mathbf{w}, \mathbf{m})$ . From Lemma 2, we can see that  $\mathbf{Z}_0 = \mathbf{X}_{0,1-\mathbf{m}} \mid \mathbf{X}_{0,\mathbf{m}} = \mathbf{u}, \mathbf{M}_0 = \mathbf{m}$  and  $\mathbf{Z}_1 = \mathbf{X}_{1,1-\mathbf{m}} \mid \mathbf{X}_{1,\mathbf{m}} = \mathbf{w}$ . Letting  $\mathbf{u} = \mathbf{w} = \mathbf{x}_\mathbf{m}$ , we obtain the desired results in Lemma 1.

## D Proof of Proposition 2

*Proof.* Since the sequential imputer (2) never changes  $\mathbf{X}_\mathbf{m}^{(t-1)}$  nor  $\mathbf{M}$ , we have

$$p_\mathbf{g}(\mathbf{x}, \mathbf{m}) = p_\mathbf{g}(\mathbf{x}_{1-\mathbf{m}}|\mathbf{x}_\mathbf{m}, \mathbf{m})p(\mathbf{x}_\mathbf{m}, \mathbf{m}), \quad (25)$$

where the marginal density  $p(\mathbf{x}_{1-\mathbf{m}}, \mathbf{m})$  does not depend on  $\mathbf{g}$ . Now we show that the KL divergence only depends on  $\mathbf{g}$  through the KL divergence  $D[p_\mathbf{g}(\cdot|\mathbf{x}_{1-\mathbf{m}}, \mathbf{m})|q(\cdot|\mathbf{x}_{1-\mathbf{m}})]$ .

$$\begin{aligned} & D[\mathbb{P}_{\mathbf{X}^{(t)}(\mathbf{g}), \mathbf{M}} | \mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}}] \\ &= \mathbb{E}_{(\mathbf{x}, \mathbf{m}) \sim (\mathbf{X}_\mathbf{g}^{(t)}, \mathbf{M})} \left[ \log \frac{p_\mathbf{g}(\mathbf{x}, \mathbf{m})}{q(\mathbf{x})p(\mathbf{m})} \right] \\ &= \mathbb{E}_{(\mathbf{x}, \mathbf{m}) \sim (\mathbf{X}_\mathbf{g}^{(t)}, \mathbf{M})} \left[ \log \frac{p_\mathbf{g}(\mathbf{x}_{1-\mathbf{m}}, \mathbf{x}_\mathbf{m}, \mathbf{m})}{q(\mathbf{x}_{1-\mathbf{m}}, \mathbf{x}_\mathbf{m})p(\mathbf{m})} \right] \\ &= \mathbb{E}_{(\mathbf{x}, \mathbf{m}) \sim (\mathbf{X}_\mathbf{g}^{(t)}, \mathbf{M})} \left[ \log \frac{p_\mathbf{g}(\mathbf{x}_{1-\mathbf{m}}|\mathbf{x}_\mathbf{m}, \mathbf{m})}{q(\mathbf{x}_{1-\mathbf{m}}|\mathbf{x}_\mathbf{m}, \mathbf{m})} \right] + \mathbb{E}_{(\mathbf{x}, \mathbf{m}) \sim (\mathbf{X}_\mathbf{g}^{(t)}, \mathbf{M})} \left[ \log \frac{p(\mathbf{x}_\mathbf{m}, \mathbf{m})}{q(\mathbf{x}_\mathbf{m}, \mathbf{m})p(\mathbf{m})} \right] \\ &= \sum_{\mathbf{m} \in \{0,1\}^d} p(\mathbf{m}) \mathbb{E}_{\mathbf{x} \sim \mathbf{X}_\mathbf{g}^{(t)} | \mathbf{M}=\mathbf{m}} \left[ \log \frac{p_\mathbf{g}(\mathbf{x}_{1-\mathbf{m}}|\mathbf{x}_\mathbf{m}, \mathbf{m})}{q(\mathbf{x}_{1-\mathbf{m}}|\mathbf{x}_\mathbf{m}, \mathbf{m})} \right] + \text{const.} \\ &= \sum_{\mathbf{m} \in \{0,1\}^d} p(\mathbf{m}) \underbrace{\int p(\mathbf{x}_\mathbf{m}|\mathbf{m}) D[p_\mathbf{g}(\mathbf{x}_{1-\mathbf{m}}|\mathbf{x}_\mathbf{m}, \mathbf{m})|q(\mathbf{x}_{1-\mathbf{m}}|\mathbf{x}_\mathbf{m})] d\mathbf{x}_\mathbf{m}}_{A_\mathbf{g}} + \text{const.}, \end{aligned}$$

where const. does not depend on  $\mathbf{g}$ . The KL divergence  $D[p_\mathbf{g}(\mathbf{x}_{1-\mathbf{m}}|\mathbf{x}_\mathbf{m}, \mathbf{m})|q(\mathbf{x}_{1-\mathbf{m}}|\mathbf{x}_\mathbf{m})]$  is zero if and only if  $p_\mathbf{g}(\mathbf{x}_{1-\mathbf{m}}|\mathbf{x}_\mathbf{m}, \mathbf{m}) = q(\mathbf{x}_{1-\mathbf{m}}|\mathbf{x}_\mathbf{m})$ . Assuming  $p(\mathbf{x}_\mathbf{m}|\mathbf{m})$  is positive everywhere,  $A_\mathbf{g}$  is minimized if and only if  $D[p_\mathbf{g}(\mathbf{x}_{1-\mathbf{m}}|\mathbf{x}_\mathbf{m}, \mathbf{m})|q(\mathbf{x}_{1-\mathbf{m}}|\mathbf{x}_\mathbf{m})]$  is zero. Thus, if  $p(\mathbf{m})$  is strictly positive,  $D[\mathbb{P}_{\mathbf{X}^{(t)}(\mathbf{g}), \mathbf{M}} | \mathbb{P}_{\mathbf{X}^{(t-1)}} \otimes \mathbb{P}_{\mathbf{M}}]$  is minimized if and only if  $p_\mathbf{g}(\mathbf{x}_{1-\mathbf{m}}|\mathbf{x}_\mathbf{m}, \mathbf{m}) = q(\mathbf{x}_{1-\mathbf{m}}|\mathbf{x}_\mathbf{m})$ .  $\square$

## E Experimental Setup

This appendix provides all information necessary for reproducibility: datasets, evaluation protocol, and computational resources. The hyperparameter settings and its sensitivity studies can be found in our supplemental material. Hyperparameters were selected through a sensitivity analysis. The synthetic datasets used in this sensitivity study consist of 1000 samples and 20 features. A missing rate of 20% was applied to each dataset. The data types include Gaussian, Uniform Correlated, and Mixed (Gaussian and Uniform). The results indicate that MIRI is *not sensitive* to the hyperparameters selections. See our supplementary material for more details.

### E.1 Datasets

**Synthetic Data.** We generate  $N = 6000$  samples in  $\mathbb{R}^2$  by drawing two equally-sized clusters of  $n = 3000$  points each from isotropic Gaussians. One cluster is centered at  $(-2, -2)$  and the other at  $(2, 2)$ , both with standard deviation  $\sigma = 0.5$ .

Table 2: UCI datasets used in our study.

Dataset	# Samples	# Features
wine	1 599	11
energy	768	8
parkinsons	5 875	20
stock	536	11
pumadyn32nm	8 192	32
housing	506	13
forest	517	12
bike	17 379	17
solar	1 066	10
gas	2 565	128

**UCI Regression Benchmarks.** Table 2 lists the ten UCI datasets used. For each, we report sample size and feature dimensionality.

**CIFAR-10.** We randomly sample 5 000  $32 \times 32$  images and apply pixel-level MCAR masks at varying rates. All three colour channels of each missing pixel are masked.

**CelebA.** We randomly sample 5 000  $64 \times 64$  images and apply channel-level MCAR masks at varying rates. All three colour channels of each missing pixel are masked independently.

## E.2 Evaluation Protocol

For each dataset, we generate ten independent MCAR masks and rerun the full training and imputation pipeline. Reported results are mean  $\pm$  standard deviation over these runs.

## E.3 Hyperparameter Selection

Values were chosen via sensitivity analysis. Complete settings and search ranges are provided in our supplemental material. MIRI exhibits *low sensitivity* to these choices.

## E.4 Computational Resources

**Tabular Experiments.** NVIDIA P100 GPU (16 GB), Intel Xeon E5-2680 v4 CPU (8 cores, 2.4 GHz), 24 GB RAM.

**Image Experiments.** NVIDIA RTX 3090 GPU (24 GB), Intel Xeon Gold 6330 CPU (14 cores, 2.0 GHz), 90 GB RAM.

## E.5 Baseline Implementations

We evaluate **HyperImpute** and **MICE** using the official HyperImpute repository with default settings.<sup>2</sup>

We adapt the official implementations of **KnewImp**<sup>3</sup>, **TabCSDI**<sup>4</sup> and **GAIN**<sup>5</sup> from their respective GitHub repositories. **MissDiff** is reimplemented based on the algorithm described in the original publication [Ouyang et al., 2023]. All baselines use default hyperparameters unless otherwise stated.

<sup>2</sup><https://github.com/vanderschaarlab/hyperimpute>.

<sup>3</sup><https://github.com/JustusvLiebig/NewImp>.

<sup>4</sup><https://github.com/pfnet-research/TabCSDI>.

<sup>5</sup><https://github.com/jsyoon0823/GAIN>.



Table 3: **Performance metrics at 30% missingness.** Metrics computed over 10 runs; values denote mean  $\pm$  standard deviation. The best values are highlighted in bold.

Method	RMSE ( $\downarrow$ )	MAE ( $\downarrow$ )	MMD ( $\downarrow$ )
<b>GAIN</b>	$1.128 \pm 0.004$	$0.600 \pm 0.002$	$0.342 \pm 0.015$
<b>TabCSDI</b>	$1.128 \pm 0.004$	$0.600 \pm 0.002$	$0.337 \pm 0.010$
<b>KnewImp</b>	$1.130 \pm 0.004$	$0.600 \pm 0.002$	$0.335 \pm 0.009$
<b>MissDiff</b>	$0.951 \pm 0.008$	$0.437 \pm 0.004$	$0.189 \pm 0.009$
<b>HyperImpute</b>	<b><math>0.862 \pm 0.020</math></b>	<b><math>0.278 \pm 0.006</math></b>	$0.094 \pm 0.018$
<b>MIRI (Ours*)</b>	$0.938 \pm 0.022$	$0.325 \pm 0.009$	<b><math>0.036 \pm 0.007</math></b>

## F Additional Experimental Results

### F.1 Performance on Synthetic Data

We evaluate each imputation method using three criteria, computed only on entries originally masked under the MCAR mechanism: (1) Root Mean Square Error (RMSE); (2) Mean Absolute Error (MAE); (3) Maximum Mean Discrepancy (MMD). Table 3 reports the mean  $\pm$  standard deviation over ten independent MCAR masks (30% missing).

Overall, HyperImpute is optimal when minimizing per-entry error, achieving roughly 8% lower RMSE and 14% lower MAE than MIRI. In contrast, MIRI reduces MMD by more than 60% relative to HyperImpute, thereby better preserving the underlying data distribution despite a modest increase in point-wise error. MissDiff offers a balanced trade-off, whereas GAIN, TabCSDI, and KnewImp underperform on both fronts.

### F.2 Additional UCI Regression Experiments

Figure 4 shows the MMD results on ten UCI regression benchmarks under three MCAR levels: 20% (top), 40% (middle), and 60% (bottom). Lower values indicate better distributional matching.

**20% Missingness.** MIRI attains an average MMD of  $\approx 0.05$ , compared to HyperImpute ( $\approx 0.07$ ), MissDiff ( $\approx 0.04$ ), and MICE ( $\approx 0.10$ ). KnewImp, TabCSDI and GAIN suffer much higher MMD— $\approx 0.38$ ,  $\approx 0.54$ , and  $\approx 0.60$ , respectively. The largest relative gains of MIRI over HyperImpute occur on high-dimensional datasets (gas, pumadyn32nm), where MIRI reduces MMD by over  $\approx 30\%$ .

**40% Missingness.** MIRI’s MMD increases slightly to  $\approx 0.04$ , while HyperImpute remains at  $\approx 0.06$  and MissDiff at  $\approx 0.03$ . MICE degrades to  $\approx 0.07$ . KnewImp, TabCSDI and GAIN remain elevated at  $\approx 0.27$ ,  $\approx 0.55$ , and  $\approx 0.55$ , respectively.

**60% Missingness.** Under severe missingness, MIRI’s MMD is  $\approx 0.05$ , still below HyperImpute ( $\approx 0.07$ ) and MissDiff ( $\approx 0.04$ ). MICE further deteriorates to  $\approx 0.10$ , while KnewImp, TabCSDI and GAIN exhibit persistently high MMD ( $\approx 0.21$ ,  $\approx 0.55$ , and  $\approx 0.52$ , respectively).

These results confirm that MIRI consistently delivers strong distributional fidelity, remains robust under high missing rates, and scales effectively to high-dimensional datasets.

### F.3 Additional CIFAR-10 Experiments

We evaluate the imputation quality of MIRI, GAIN, KnewImp, and HyperImpute on 15 randomly selected CIFAR-10 images corrupted under pixel-level MCAR with missing rates of 20%, 40% and 60% (see Figures 5, 6 and 7). Across all missing-data regimes, MIRI consistently produces accurate, visually coherent reconstructions, in agreement with the quantitative metrics in Table 1. Although HyperImpute achieves competitive performance at 20% missingness, MIRI outperforms it at both 40% and 60% missingness.

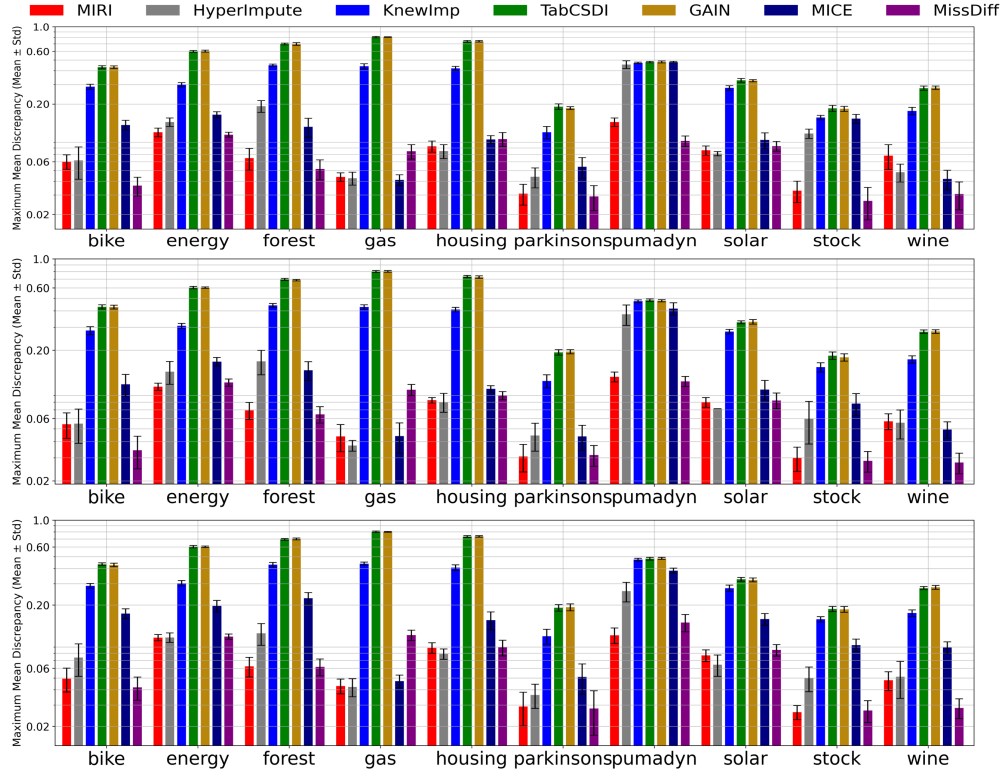


Figure 4: MMD on 10 UCI datasets (Above: 20% missingness, Middle: 40% missingness, Below: 60 % missingness). The lower the better.

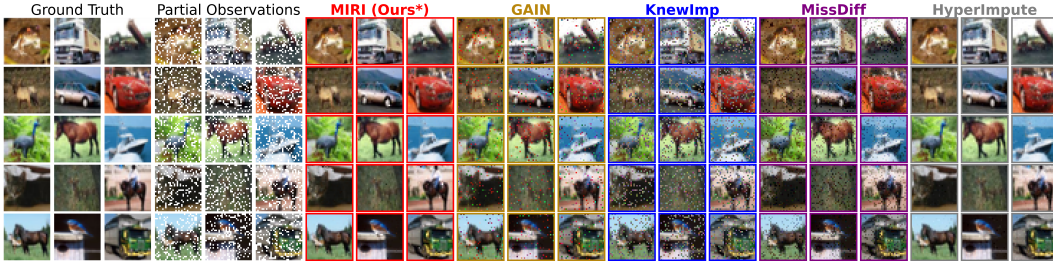


Figure 5: 15 uncurated 32x32 CIFAR-10 images and their imputations. 20% of pixels randomly removed from *all RGB channels*.

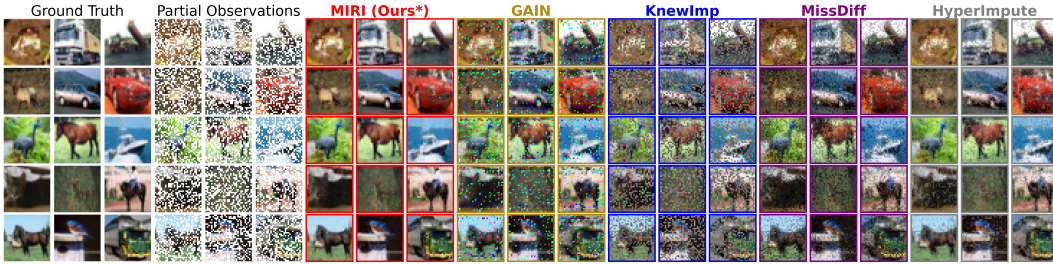


Figure 6: 15 uncurated 32x32 CIFAR-10 images and their imputations. 40% of pixels randomly removed from *all RGB channels*.

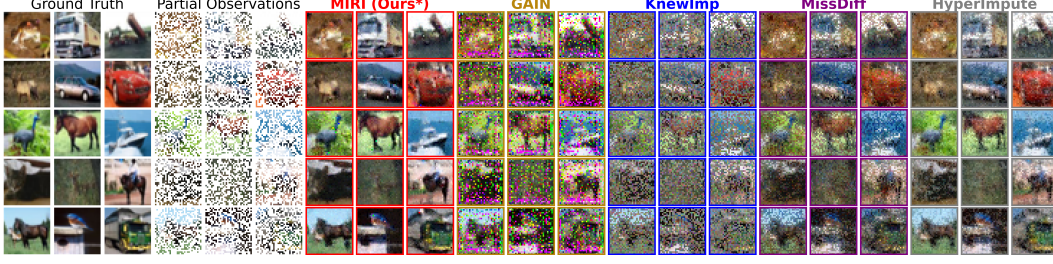


Figure 7: 15 uncurated  $32 \times 32$  CIFAR-10 images and their imputations. 60% of pixels randomly removed from *all* RGB channels.

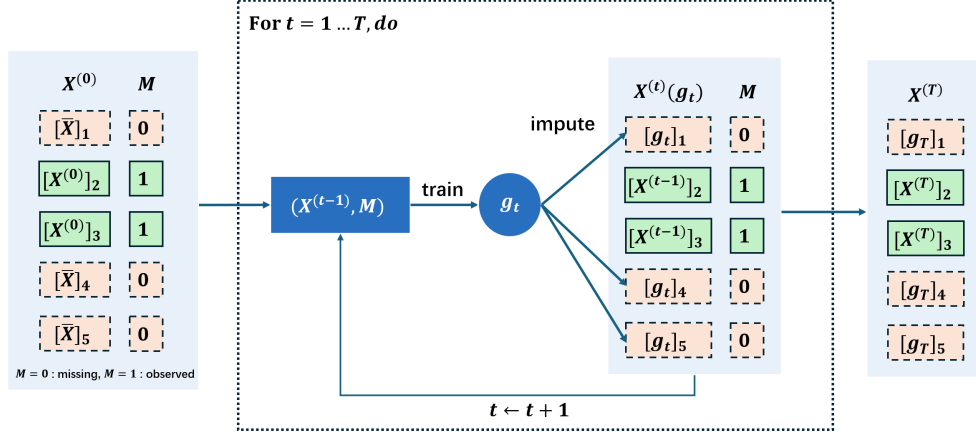


Figure 8: A basic schematic of the sequential imputation algorithm for  $t = 1 \dots T$ .  $[\mathbf{X}]_i$  denotes the  $i$ -th components of  $\mathbf{X}$ . In this example,  $\mathbf{X}^{(0)}$  are initialized using the mean ( $\bar{\mathbf{X}}$ ) imputation.

## G Illustration of Sequential Imputation

Figure 8 illustrates the sequential imputation process of MIRI.

## H Euler ODE Solver

In this section, we provide the Euler solver (Algorithm H) used in our experiments.

---

**Algorithm 2** ODE Solver with Euler Method

---

**Require:** Velocity field  $\mathbf{v}$ , Initial condition  $\{(\mathbf{X}_i, \mathbf{M}_i)\}$ , Number of Euler steps  $N$ .

```
1: for  $k = 1$  to  $N$  do  
2:    $\tau = \frac{k}{N}$   
3:    $\forall i, \mathbf{X}_i \leftarrow \mathbf{X}_i + \frac{1}{N} \cdot \mathbf{v}(\mathbf{X}_i, \mathbf{M}_i, \tau)$   
4: end for  
5: return  $\{\mathbf{X}_i\}$ 
```

---