MISSSCORE: HIGH-ORDER SCORE ESTIMATION IN THE PRESENCE OF MISSING DATA

Anonymous authors

Paper under double-blind review

ABSTRACT

The first order derivative (score) of data density, typically estimated via denoising score matching, has emerged as an effective tool for modeling data distribution and generating synthetic data. Extending this concept to higher-order scores could uncover more detailed local information of the data distribution, enabling new applications. However, learning these high-order scores usually requires complete data, which is often unavailable in real-world scenarios such as healthcare and finance due to privacy and cost constraints. In this work, we introduce MissScore, a novel score-based framework for learning high-order scores from observations with missing data. We derive objective functions for estimating high-order scores under different missing data mechanisms and propose a new algorithm to handle missing data effectively. Our empirical results demonstrate that MissScore efficiently and accurately approximates high-order scores with missing data, while enhancing sampling speed and data quality, as validated through several downstream tasks, including data generation and causal discovery.

023 024 025

026 027

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

The first-order derivative of the log data density, also known as (Stein) score (Liu et al., 2016), 028 plays an important role in various machine learning applications, including image and tabular data 029 synthesis (Song & Ermon, 2019; 2020; Kim et al., 2022), super-resolution (Li et al., 2022), and inverse problems in medical imaging (Song et al., 2021; Chung & Ye, 2022). Denoising Score 031 Matching (DSM) (Vincent, 2011), an efficient method for estimating the score of the data density from samples, has become widely used in training score-based generative models (Ho et al., 2020; 033 Song & Ermon, 2020). Beyond the first-order score, high-order derivatives of the data density, which 034 we refer to as high-order scores, offer more refined local approximations of the data distribution, such as its curvature, and enable new model capabilities. For instance, they can improve the mixing speed of sampling methods (Dalalyan & Karagulyan, 2019; Sabanis & Zhang, 2019; Meng et al., 2021) and provide insights into quantifying the uncertainty in denoising problem Meng et al. (2021). 037 Additionally, Lu et al. (2022) empirically demonstrated that incorporating high-order score matching improves the likelihood of score-based diffusion ordinary differential equations on both synthetic and real data, while maintaining high-quality generation. Furthermore, high-order scores have been 040 utilized in recovering causal structures by bridging the gap between the score of data density and the 041 underlying graph topology (Rolland et al., 2022; Sanchez et al., 2022; Liu et al., 2024). 042

Despite their promise, learning high-order scores usually requires training the model on complete 043 data (Meng et al., 2021; Lu et al., 2022). However, in many real-world scenarios such as health-044 care, finance, and social networks, data often contain missing values due to privacy constraints or high sampling costs (Rubin, 1976; Shpitser, 2016). A straightforward approach to handling missing 046 data is to impute the missing values and train the model on the imputed dataset. Yet, imputation 047 methods can compromise data quality, potentially leading to biased results and significantly degrad-048 ing performance in downstream tasks (Ouyang et al., 2023). Specifically, imputation methods fail to account for the inherent uncertainty in the missing data, resulting in a distribution over imputed data that does not accurately reflect the true data distribution. Some alternative approaches, such as 051 using generative adversarial networks (GANs) or variational auto-encoders (VAEs) to directly approximate the data generation model from incomplete data (Li et al., 2019; Gain & Shpitser, 2018), 052 require training additional networks, which can be computationally expensive and may also result in model inconsistency. To address this issue, some works propose to explicitly constrain the learning objective within the model, which can enhance performance (Städler & Bühlmann, 2012; Gao et al., 2022). These studies emphasize the need for alternative unbiased approaches that can directly and more efficiently handle missing data.

057 Back to score estimation from incomplete data, Ouyang et al. (2023) adopts a similar way by intro-058 ducing a diffusion-based framework that learns the first-order score directly from incomplete data. 059 In principle, high-order scores could be estimated from a learned first-order score model trained on 060 incomplete data using automatic differentiation. However, this approach becomes computationally 061 impractical for high-dimensional data and large model sizes, particularly when using deep neu-062 ral networks based models (Meng et al., 2021). Furthermore, automatic differentiation introduces 063 additional estimation errors, as small errors does not always lead to a small estimation error for 064 high-order scores. Moreover, methods based on GANs (Goodfellow et al., 2020) or VAEs (Kingma, 2013) do not inherently capture score information, regardless of whether data is missing or com-065 plete. In contrast, score-based models naturally integrate this information (Li et al., 2019; Ho et al., 066 2020; Gain & Shpitser, 2018). These limitations highlight the need of using score-based models 067 when estimating high-order scores in the presence of missing data. 068

Contributions. In this work, we propose a novel score-based framework, which we call MissScore, 069 for learning high-order scores from incomplete data. We derive objective functions for estimating high-order scores under different missing mechanisms, leveraging DSM to recover the true score 071 function. While our framework can be applied to scores of any order, we focus on second-order 072 scores (i.e., the Hessian of the log density) in our experiments. Our empirical results demonstrate 073 that the proposed models efficiently and accurately approximate high-order scores with missing 074 data. Moreover, we show that our high-order score-based model enhances both sampling speed 075 and data quality in data generation tasks, even with missing data. The quality of these generated 076 samples is further validated across several downstream tasks. Additionally, we introduce a novel 077 causal discovery method for missing data, which scales effectively with the number of variables and samples, and performs competitively with state-of-the-art methods (Gao et al., 2022; Vo et al., 2024) 079 in causal discovery.

The paper is organized as follows: Section 2 introduces the background on high-order DSM and formulates the objective functions for estimating high-order scores under different missing data mechanisms. In Section 3, we detail the proposed training method and evaluate the accuracy of the learned scores. Sections 4 and 5 showcase experimental results on downstream tasks, illustrating the method's effectiveness. Finally, Section 6 concludes the paper. The Appendix includes related work, proofs, and supplementary experimental details.

2 ESTIMATING HIGH-ORDER SCORES WITH MISSING DATA

In this section, we provide background on high-order denoising score matching and present our method to deal with missing data.

2.1 BACKGROUND ON HIGH-ORDER DENOISING SCORE MATCHING

Consider a data distribution $p_{data}(\mathbf{x})$ and a model distribution $p(\mathbf{x}; \boldsymbol{\theta})$ over \mathbb{R}^d . The score functions of $p_{data}(\mathbf{x})$ and $p(\mathbf{x}; \boldsymbol{\theta})$ are denoted as $\mathbf{s}_1(\mathbf{x}) = \nabla_{\mathbf{x}} \log p_{data}(\mathbf{x})$ and $\mathbf{s}_1(\mathbf{x}; \boldsymbol{\theta}) = \nabla_{\mathbf{x}} \log p(\mathbf{x}; \boldsymbol{\theta})$, respectively. In DSM, instead of directly estimating the score function from the original data, the method works by introducing noise from a predefined noise distribution $q_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})$ into the data. The objective is then to estimate the score of the perturbed data distribution $q_{\sigma}(\tilde{\mathbf{x}}) = \int q_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})p_{data}(\mathbf{x})d\mathbf{x}$. To achieve so, DSM minimizes the following objective function

$$\frac{1}{2} \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \mathbb{E}_{q_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})} \left[\| \mathbf{s}_{1}(\tilde{\mathbf{x}}; \boldsymbol{\theta}) - \nabla_{\tilde{\mathbf{x}}} \log q_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x}) \|_{2}^{2} \right].$$
(1)

102 It has been shown that minimizing Eq. (1) is equivalent to minimizing the score matching loss 103 between $\mathbf{s}_1(\tilde{\mathbf{x}}; \boldsymbol{\theta})$ and $\mathbf{s}_1(\tilde{\mathbf{x}})$ (Vincent, 2011) under certain regularity conditions. When the noise 104 distribution $q_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})$ is Gaussian, i.e., $\mathcal{N}(\tilde{\mathbf{x}}|\mathbf{x}, \sigma^2 \mathbf{I})$, the objective simplifies to

105 106

107

100 101

087

088 089

091 092

$$\mathcal{L}_{\text{DSM}}(\boldsymbol{\theta}) = \frac{1}{2} \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \mathbb{E}_{q_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})} \left[\left\| \mathbf{s}_{1}(\tilde{\mathbf{x}}; \boldsymbol{\theta}) + \frac{1}{\sigma^{2}} (\tilde{\mathbf{x}} - \mathbf{x}) \right\|_{2}^{2} \right].$$
(2)

The learned score function implicitly learns how to "denoise" the perturbed data \tilde{x} , guiding it back 109 toward the true data distribution through the optimization of Eq. (2). By focusing on estimating the 110 score of the noise-perturbed distribution $q_{\sigma}(\tilde{\mathbf{x}})$ instead of the original data distribution $p_{\text{data}}(\mathbf{x})$, DSM 111 offers a more efficient approach to score estimation compared to other techniques (Hyvärinen & 112 Dayan, 2005; Song et al., 2020). When σ approaches zero, the perturbed distribution $q_{\sigma}(\tilde{\mathbf{x}})$ closely approximates $p_{\text{data}}(\mathbf{x})$, meaning that the score estimated by DSM for $q_{\sigma}(\tilde{\mathbf{x}})$ is nearly identical to that 113 of $p_{\text{data}}(\mathbf{x})$. Meng et al. (2021) provide a derivation of DSM using Tweedie's formula (Efron, 2011), 114 and they generalize this approach to incorporate high-order moments of x based on \tilde{x} , allowing them 115 to develop an objective function for learning high-order scores. 116

117 118 119 120 **Theorem 1** (Meng et al., 2021) $\mathbb{E}[\otimes^n \mathbf{x}|\tilde{\mathbf{x}}] = f_n(\tilde{\mathbf{x}}, \mathbf{s}_1, ..., \mathbf{s}_n)$, where $\otimes^n \mathbf{x} \in \mathbb{R}^{D^n}$ denotes *n*-fold tensor multiplications, $f_n(\tilde{\mathbf{x}}, \mathbf{s}_1(\tilde{\mathbf{x}}), ..., \mathbf{s}_n(\tilde{\mathbf{x}}))$ is a polynomial of $\tilde{\mathbf{x}}$, $\mathbf{s}_1(\tilde{\mathbf{x}})$, ..., $\mathbf{s}_n(\tilde{\mathbf{x}})$, and $\mathbf{s}_k(\tilde{\mathbf{x}})$ represents the k-th order score of $q_\sigma(\tilde{\mathbf{x}}) = \int p_{\text{data}}(\mathbf{x})q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})d\mathbf{x}$.

Theorem 1 shows that there exists an equality between high-order moments of the posterior distribution of x given \tilde{x} and high-order scores with respect to \tilde{x} . Leveraging Theorem 1 and the least squares estimation of $\mathbb{E}[\otimes^k x | \tilde{x}]$, the objectives for approximating the *k*-th order scores $s_k(\tilde{x})$ can be constructed as follows.

Theorem 2 (Meng et al., 2021) Given score functions $\mathbf{s}_1(\tilde{\mathbf{x}}), \ldots, \mathbf{s}_{k-1}(\tilde{\mathbf{x}})$, a k-th order score model $\mathbf{s}_k(\tilde{\mathbf{x}}; \boldsymbol{\theta})$ can be obtained by optimizing the following objective:

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \mathbb{E}_{q_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})} \left[\left\| \otimes^k \mathbf{x} - f_k(\tilde{\mathbf{x}}, \mathbf{s}_1(\tilde{\mathbf{x}}), \dots, \mathbf{s}_{k-1}(\tilde{\mathbf{x}}), \mathbf{s}_k(\tilde{\mathbf{x}}; \boldsymbol{\theta})) \right\|_2^2 \right].$$

where f_k is a polynomial of $\{\tilde{\mathbf{x}}, \mathbf{s}_1(\tilde{\mathbf{x}}), \dots, \mathbf{s}_k(\tilde{\mathbf{x}})\}$ such that

$$f_k(\tilde{\mathbf{x}}, \mathbf{s}_1(\tilde{\mathbf{x}}), \dots, \mathbf{s}_k(\tilde{\mathbf{x}})) = \begin{cases} \tilde{\mathbf{x}} + \sigma^2 \mathbf{s}_1(\tilde{\mathbf{x}}), & \text{if } k = 1, \\ \sigma^2 \frac{\partial}{\partial \tilde{\mathbf{x}}} f_{k-1}(\tilde{\mathbf{x}}, \mathbf{s}_1, \dots, \mathbf{s}_{k-1}) & \\ + \sigma^2 f_{k-1}(\tilde{\mathbf{x}}, \mathbf{s}_1, \dots, \mathbf{s}_{k-1}) \otimes \left(\mathbf{s}_1(\tilde{\mathbf{x}}) + \frac{\tilde{\mathbf{x}}}{\sigma^2}\right), & \text{if } k \ge 2. \end{cases}$$
(3)

133 134 135

136 137

138

125

126

127 128 129

130 131

We have $\mathbf{s}_k(\mathbf{\tilde{x}}; \boldsymbol{\theta}^*) = \mathbf{s}_k(\mathbf{\tilde{x}})$ for almost all $\mathbf{\tilde{x}}$.

2.2 HIGH-ORDER DENOISING SCORE MATCHING WITH MISSING DATA

Consider an observation $\mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$, which is sampled from an unknown data 139 distribution $p_{\text{data}}(\mathbf{x})$. Associated with \mathbf{x} , there is a binary mask $\mathbf{m} = (m_1, m_2, \ldots, m_d) \in \{0, 1\}^d$ 140 where $m_i = 1$ indicates that x_i is missing, and $m_i = 0$ indicates that x_i is observed. The observed 141 data can be express as $\mathbf{x}_{obs} = \mathbf{x} \odot (1 - \mathbf{m}) + na \odot \mathbf{m}$, where \odot is the element-wise multiplication 142 and na indicates the missing value. To perturb the original observation \mathbf{x} , we adopt a Gaussian 143 mechanism, generating $\tilde{\mathbf{x}}$ from the conditional distribution $\tilde{\mathbf{x}} | \mathbf{x} \sim \mathcal{N}(\mathbf{x}, \sigma^2 \mathbf{I}_d)$, where σ is a pre-144 specified constant. Therefore, the conditional density function of $\tilde{\mathbf{x}}$ given \mathbf{x} is defined as $q_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x}) :=$ 145 $(2\pi\sigma^2)^{-\frac{d}{2}}\exp\{-\frac{(\tilde{\mathbf{x}}-\mathbf{x})^{\top}(\tilde{\mathbf{x}}-\mathbf{x})}{2\sigma^2}\}.$ 146

The missing mechanisms can be categorized based on the relationships between the mask m and
the complete data x as follows Rubin (1976) and the data generation process for each mechanism is
detailed in Appendix D:

- Missing Completely at Random: mask m is independent with the completed data x.
- Missing at Random: mask m only depends on the observed value x_{obs}.
- Missing Not at Random: m depends on the observed value \mathbf{x}_{obs} and the missing value.

In the following section, we derive the objective functions for training the high-order DSM model, in the presence of missing data under the M(C)AR mechanisms.

157 158

150

151

152

153

2.2.1 MISSING COMPLETELY AT RANDOM (MCAR)

In this section, we start by exploring the first- and second-order scores, then extend the approach to scores of any order under the MCAR assumption. The following theorem presents our first theoretical result, showing that DSM with a missing mask can recover the oracle score, which is the gradient of $\log p_{\text{data}}(\mathbf{x})$ with respect to \mathbf{x} . **Theorem 3** If the missing mechanism of \mathbf{x} is MCAR, with the missing probability of every element lying between 0 and 1, i.e., $p(m_i = 1) \in [0, 1)$ for all $i \in \{1, 2, ..., d\}$. We denote the objective $\mathcal{J}_{\text{DSM}}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x},\mathbf{m}} \mathbb{E}_{\tilde{\mathbf{x}}|\mathbf{x},\mathbf{m}} \left[\left\| \left\{ \mathbf{s}_1(\tilde{\mathbf{x}}; \boldsymbol{\theta}) + \frac{1}{\sigma^2}(\tilde{\mathbf{x}} - \mathbf{x}) \right\} \odot (\mathbf{1} - \mathbf{m}) \right\|_2^2 \right]$. Then we have,

$$\arg\min_{\boldsymbol{\theta}} \mathcal{J}_{\text{DSM}}(\boldsymbol{\theta}) = \arg\min_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\tilde{\mathbf{x}} \mid \mathbf{x}} \left[\left\| \{ \mathbf{s}_1(\tilde{\mathbf{x}}; \boldsymbol{\theta}) - \mathbf{s}_1(\tilde{\mathbf{x}}) \} \odot \sqrt{\mathbb{P}(\mathbf{m} = 0)} \right\|_2^2 \right]$$

Furthermore, if there exist unique θ^* such that $\mathbf{s}_1(\tilde{\mathbf{x}}) = \mathbf{s}_1(\tilde{\mathbf{x}}; \theta^*)$, then $\theta^* = \underset{\theta}{\operatorname{arg min}} \mathcal{J}_{\text{DSM}}(\theta)$.

171 Since $\mathbb{P}(\mathbf{m} = 0) > 0$, we can conclude that the global optimal of DSM with missing mask coincides 172 with the orcale score. The proof is provided in Appendix C.1. For the second-order score model, 173 building upon Theorem 2, the second-order DSM with missing mask can learn the oracle second-174 order score, which is the Hessian of $\log p_{\text{data}}(\mathbf{x})$ with respect to \mathbf{x} .

Theorem 4 Suppose the first-order score $\mathbf{s}_1(\tilde{\mathbf{x}})$ is given, and the missing mechanism of \mathbf{x} is MCAR, with the missing probability of every element lying between 0 and 1, i.e., $p(m_i = 1) \in [0, 1)$ for all $i \in \{1, 2, ..., d\}$. We denote the objective

$$\mathcal{J}_{D_2SM}(\boldsymbol{\theta}) = \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x},\mathbf{m}} \left[\left\| \left\{ \mathbf{s}_2(\tilde{\mathbf{x}};\boldsymbol{\theta}) + \mathbf{s}_1(\tilde{\mathbf{x}})\mathbf{s}_1^\top(\tilde{\mathbf{x}}) + \frac{\mathbf{I} - \mathbf{z}\mathbf{z}^\top}{\sigma^2} \right\} \odot \left\{ (\mathbf{1} - \mathbf{m})(\mathbf{1} - \mathbf{m})^\top \right\} \right\|_2^2 \right],$$

where $\mathbf{z} = \frac{\tilde{\mathbf{x}} - \mathbf{x}}{\sigma}$. Then we have,

$$\arg\min_{\boldsymbol{\theta}} \mathcal{J}_{D_2SM}(\boldsymbol{\theta}) = \arg\min_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\tilde{\mathbf{x}}|\mathbf{x}} \left[\left\| \{ \mathbf{s}_2(\tilde{\mathbf{x}}; \boldsymbol{\theta}) - \mathbf{s}_2(\tilde{\mathbf{x}}) \} \odot \sqrt{\mathbb{P}\left[(\mathbf{1} - \mathbf{m})(\mathbf{1} - \mathbf{m})^\top = 0 \right]} \right\|_2^2 \right]$$

Furthermore, if there exist θ^* such that $\mathbf{s}_2(\tilde{\mathbf{x}}) = \mathbf{s}_2(\tilde{\mathbf{x}}; \theta^*)$, then $\theta^* = \underset{\theta}{\operatorname{arg min}} \mathcal{J}_{D_2SM}(\theta)$.

186 187

191

214 215

185

179 180 181

182 183

166

167 168

169 170

Since $\mathbb{P}\left[(\mathbf{1} - \mathbf{m})(\mathbf{1} - \mathbf{m})^{\top} = 0\right] > 0$, we can conclude that the global optimal of second-order DSM with missing mask aligns with the oracle second order score. The proof is provided in Appendix C.2.

192 2.2.2 MISSING AT RANDOM (MAR)

In this section, we investigate the first- and second-order scores under the MAR assumption. MAR allows missingness to depend on the observed data but not on the unobserved (missing) data. This is more realistic than the more restrictive MCAR assumption, where the probability of missingness is independent of both observed and unobserved data. However, the probability of missing data depends on the observed data, which can introduce bias in estimates if the missing data mechanism is ignored. Thus, we cannot learn the oracle score models through objective functions proposed in Section 2.2.1.

200 Inverse Probability Weighting (IPW) is essential in the MAR framework because it provides a way to 201 correct for the bias introduced by the missing data mechanism (Wooldridge, 2007; Seaman & White, 202 2013). IPW achieves this by reweighting the observed data points, compensating for the fact that 203 some observations are more likely to be missing than others. Each observed data point is weighted 204 by the inverse of its probability of being observed (i.e., the probability that it was not missing). This means that observations with a higher chance of being missing receive a higher weight, while those 205 with a lower chance of being missing receive a lower weight. By reweighting in this way, IPW 206 ensures that the estimates reflect what would have been observed if there were no missing data, thus 207 mitigating the bias introduced by the MAR mechanism. 208

In the following theorem, we present our theoretical result that DSM with missing mask through IPW can learn the oracle score, i.e., the gradient of $\log p_{data}(\mathbf{x})$ w.r.t. \mathbf{x} .

Theorem 5 If the missing mechanism of **x** is MAR, with the missing probability of every element lying between 0 and 1, i.e., $p(m_i = 1) \in [0, 1)$ for all $i \in \{1, 2, ..., d\}$. We denote the objective

$$\mathcal{J}_{DSM}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x},\mathbf{m}} \mathbb{E}_{\tilde{\mathbf{x}}|\mathbf{x},\mathbf{m}} \left[\left\| \left\{ \mathbf{s}_1(\tilde{\mathbf{x}};\boldsymbol{\theta}) + \frac{1}{\sigma^2} (\tilde{\mathbf{x}} - \mathbf{x}) \right\} \odot \left\{ \frac{1 - \mathbf{m}}{\sqrt{\mathbb{P}[\mathbf{m} = 0|\mathbf{x} = \mathbf{x}]}} \right\} \right\|_2^2 \right]$$

216217 Then we have,

$$\arg\min_{\boldsymbol{\theta}} \mathcal{J}_{\text{DSM}}(\boldsymbol{\theta}) = \arg\min_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\tilde{\mathbf{x}} \mid \mathbf{x}} \left[\left\| \{ \mathbf{s}_1(\tilde{\mathbf{x}}; \boldsymbol{\theta}) - \mathbf{s}_1(\tilde{\mathbf{x}}) \} \right\|_2^2 \right].$$

Furthermore, if there exist unique θ^* such that $\mathbf{s}_1(\tilde{\mathbf{x}}) = \mathbf{s}_1(\tilde{\mathbf{x}}; \theta^*)$, then $\theta^* = \underset{\theta}{\operatorname{arg\,min}} \mathcal{J}_{\mathsf{DSM}}(\theta)$.

We can conclude that the global optimal of DSM with missing mask coincides with the oracle score under MAR. The proof is provided in Appendix C.4. Building on Theorem 2, the second-order score model using DSM with a missing mask under MAR can learn the oracle second-order score, which is the Hessian of $\log p_{data}(\mathbf{x})$ with respect to \mathbf{x} .

Theorem 6 Suppose the first-order score $\mathbf{s}_1(\tilde{\mathbf{x}})$ is given, if the missing mechanism of \mathbf{x} is MAR, and the missing probability of every element lies between 0 and 1, which is $p(m_i m_j = 0) \in [0, 1)$ for all $i, j \in \{1, 2, ..., d\}$, we denote the objective

$$\mathcal{J}_{D_2SM}(\boldsymbol{\theta}) = \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x},\mathbf{m}} \left[\left\| \left\{ \mathbf{s}_2(\tilde{\mathbf{x}};\boldsymbol{\theta}) + \mathbf{s}_1(\tilde{\mathbf{x}})\mathbf{s}_1^\top(\tilde{\mathbf{x}}) + \frac{\mathbf{I} - \mathbf{z}\mathbf{z}^\top}{\sigma^2} \right\} \odot \left\{ \frac{(\mathbf{1} - \mathbf{m})(\mathbf{1} - \mathbf{m})^\top}{\sqrt{\mathbb{P}[\mathbf{m}\mathbf{m}^\top = 0|\mathbf{x} = \mathbf{x}]}} \right\} \right\|_2^2 \right],$$

where $\mathbf{z} = \frac{\tilde{\mathbf{x}} - \mathbf{x}}{\sigma}$, then we have,

$$\arg\min_{\boldsymbol{\theta}} \mathcal{J}_{D_2SM}(\boldsymbol{\theta}) = \arg\min_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\tilde{\mathbf{x}}|\mathbf{x}} \left[\left\| \{ \mathbf{s}_2(\tilde{\mathbf{x}}; \boldsymbol{\theta}) - \mathbf{s}_2(\tilde{\mathbf{x}}) \} \right\|_2^2 \right],$$

If there exist θ^* such that $\mathbf{s}_2(\tilde{\mathbf{x}}) = \mathbf{s}_2(\tilde{\mathbf{x}}; \theta^*)$, then $\theta^* = \arg\min_{\theta} \mathcal{J}_{D_2SM}(\theta)$.

We can conclude that the global optimal of the second-order DSM with missing mask aligns with the oracle Hessian under MAR. The proof is provided in Appendix C.5. We now extend this approach to any desired order in both MCAR and MAR. Theorem 7 indicates that k-th order DSM with missing mask can learn the oracle k-th order score.

Theorem 7 Given score functions $\mathbf{s}_1(\tilde{\mathbf{s}}), \dots, \mathbf{s}_{k-1}(\tilde{\mathbf{s}})$, and the missing probability of every element lies between 0 and 1, which is $p(m_i = 1) \in [0, 1)$ for all $i \in \{1, 2, ..., d\}$. If we correctly model the k-th order derivative $\mathbf{s}_k(\tilde{\mathbf{x}})$, there exists θ^* such that $\mathbf{s}_k(\tilde{\mathbf{x}}, \theta^*) = \mathbf{s}_k(\tilde{\mathbf{x}})$, then

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x},\mathbf{m}} \left[\| \left\{ \otimes^k \mathbf{x} - f_k(\tilde{\mathbf{x}}, \mathbf{s}_1(\tilde{\mathbf{x}}), \dots, \mathbf{s}_{k-1}(\tilde{\mathbf{x}}), \mathbf{s}_k(\tilde{\mathbf{x}}; \boldsymbol{\theta})) \right\} \odot \otimes^k \mathbf{w} \|^2 \right],$$

where $\mathbf{w} = \mathbf{1} - \mathbf{m}$ if the missing mechanism of \mathbf{x} is MCAR, and $\mathbf{w} = \frac{\mathbf{1} - \mathbf{m}}{\mathbb{P}[\mathbf{m}^k = 0 | \mathbf{x} = \mathbf{x}]}$ if the missing mechanism of \mathbf{x} is MAR.

The proof is provided in Appendix C.3. The k-th order score model $s_2(\tilde{x}; \theta)$ can be learned from observed data via optimizing the following objective,

$$\mathcal{L}_{\mathbf{D}_k \mathbf{SM}}(\boldsymbol{\theta}) = \mathbb{E}_{\tilde{\mathbf{x}}_{obs}, \mathbf{x}_{obs}, \mathbf{m}} \left[\| \left\{ \otimes^k \mathbf{x}_{obs} - f_k(\tilde{\mathbf{x}}, \mathbf{s}_1(\tilde{\mathbf{x}}_{obs}), \dots, \mathbf{s}_{k-1}(\tilde{\mathbf{x}}_{obs}), \mathbf{s}_k(\tilde{\mathbf{x}}_{obs}; \boldsymbol{\theta})) \right\} \odot \otimes^k \mathbf{w} \|^2 \right],$$

$$(4)$$

where f_k is a polynomial of $\{\tilde{\mathbf{x}}, \mathbf{s}_1(\tilde{\mathbf{x}}), \dots, \mathbf{s}_k(\tilde{\mathbf{x}})\}$ is defined in Eq. (3), and \mathbf{w} is defined in Theorem 7.

3 TRAINING SCORE MODELS BY HIGH-ORDER DSM WITH MISSING DATA

In this section, we describe the training process for high-order score models in the presence of missing data and evaluate their empirical performance. While our analysis specifically focuses on the first- and second-order scores, the approach can be applied to any order of scores.

Based on Theorem 3 and Theorem 5, the first-order score model $s_1(\tilde{x}; \theta)$ is learned by minimizing the following objective function using the observed data,

$$\mathcal{L}_{\text{DSM}}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}_{\text{obs}},\mathbf{m}} \mathbb{E}_{\tilde{\mathbf{x}}_{\text{obs}}|\mathbf{x}_{\text{obs}},\mathbf{m}} \left[\left\| \left\{ \mathbf{s}_{1}(\tilde{\mathbf{x}}_{\text{obs}};\boldsymbol{\theta}) + \frac{1}{\sigma^{2}}(\tilde{\mathbf{x}}_{\text{obs}} - \mathbf{x}_{\text{obs}}) \right\} \odot \mathbf{w}_{1} \right\|_{2}^{2} \right],$$
(5)

where
$$\mathbf{w}_1 = \mathbf{1} - \mathbf{m}$$
 under MCAR, and $\mathbf{w}_1 = \frac{\mathbf{1} - \mathbf{m}}{\sqrt{\mathbb{P}[\mathbf{m} = 0|\mathbf{x} = \mathbf{x}_{obs}]}}$ under MAR.

272 Similarly, the second-order score model $s_2(\tilde{\mathbf{x}}; \theta)$ is learned using the following objective, as derived 273 in Theorem 4 and Theorem 6, 274

$$\mathcal{L}_{D_2SM}(\boldsymbol{\theta}) = \mathbb{E}_{\tilde{\mathbf{x}}_{obs}, \mathbf{x}_{obs}, \mathbf{m}} \left[\left\| \left\{ \mathbf{s}_2(\tilde{\mathbf{x}}_{obs}; \boldsymbol{\theta}) + \mathbf{s}_1(\tilde{\mathbf{x}}_{obs}) \mathbf{s}_1^\top(\tilde{\mathbf{x}}_{obs}) + \frac{\mathbf{I} - \mathbf{z}_{obs} \mathbf{z}_{obs}^\top}{\sigma^2} \right\} \odot \mathbf{w}_2 \right\|_2^2 \right], \tag{6}$$

where $\mathbf{z}_{obs} = (\tilde{\mathbf{x}} - \mathbf{x})/\sigma$. Under MCAR, $\mathbf{w}_2 = (\mathbf{1} - \mathbf{m})(\mathbf{1} - \mathbf{m})^{\top}$, and under MAR, $\mathbf{w}_2 = (\mathbf{1} - \mathbf{m})(\mathbf{1} - \mathbf{m})^{\top}$. $\sqrt{\mathbb{P}[\mathbf{m}\mathbf{m}^{\top}=0|\mathbf{x}=\mathbf{x}]}$

279

275 276

277 278

284

However, training the second-order score model, $s_2(\tilde{x}; \theta)$, requires knowledge of the first-order 281 score $s_1(\tilde{x})$. Therefore, we adopt a multi-task objective to train both $s_1(\tilde{x}; \theta)$ and $s_2(\tilde{x}; \theta)$ simulta-282 neously, 283

$$\mathcal{L}_{\text{joint}}(\boldsymbol{\theta}) = \mathcal{L}_{\text{DSM}}(\boldsymbol{\theta}) + \omega \mathcal{L}_{\text{D}_2\text{SM}}(\boldsymbol{\theta}), \tag{7}$$

where $\omega \in \mathbb{R}^+$ is a tunable coefficient. $\mathcal{L}_{\text{DSM}}(\theta)$ and $\mathcal{L}_{\text{D}_2\text{SM}}(\theta)$ correspond to Eq. (5) and Eq. (6). 285 $\mathbb{P}[\mathbf{m} = 0 | \mathbf{x} = \mathbf{x}_{obs}]$ and $\mathbb{P}[\mathbf{mm}^{\top} = 0 | \mathbf{x} = \mathbf{x}_{obs}]$ in MAR are estimated using logistic regression 286 models, where the response variable indicates whether the data is missing or observed, and the 287 predictors are the observed variables. In the experiments, missing values are handled by replacing 288 them with 0 for continuous variables and creating a new category for discrete variables. One-hot 289 encoding is then applied to discrete variables. Element-wise multiplication with the mask naturally 290 mitigates the impact of replacing missing values with zeros when computing the objective. The 291 algorithm is provided in Appendix D.

292 293

300

302

304

308

310

3.1 IMPROVING STABILITY WITH VARIANCE REDUCTION

295 It is important to note that, in order to match the score of the true distribution $p_{\text{data}}(\mathbf{x})$, σ needs to be close to zero for both DSM and D₂SM, so that $q_{\sigma}(\tilde{\mathbf{x}})$ closely approximates $p_{\text{data}}(\mathbf{x})$. However, 296 training score models using denoising methods can suffer from high variance when σ approaches 297 zero. This challenge motivates the use of variance reduction techniques. Building on existing vari-298 ance reduction methods for DSM (Song & Kingma, 2021; Meng et al., 2021), we propose tailored 299 variance reduction techniques specifically for training DSM with missing data, as follows

$$\mathcal{L}_{\text{DSM-VR}}(\boldsymbol{\theta}) = \mathcal{L}_{\text{DSM}}(\boldsymbol{\theta}) - \mathbb{E}_{\mathbf{x}_{\text{obs}},\mathbf{m}} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0},\mathbf{I})} \left[\left(\frac{2}{\sigma} \mathbf{s}(\mathbf{x}_{\text{obs}};\boldsymbol{\theta})^{\top} \mathbf{z} \right) \odot \mathbf{g}_{1}(\mathbf{x}_{\text{obs}},\mathbf{m}) + \frac{\|\mathbf{z} \odot \mathbf{g}_{1}(\mathbf{x}_{\text{obs}},\mathbf{m})\|^{2}}{\sigma^{2}} \right],$$

$$(8)$$

where
$$\mathbf{g}_1(\mathbf{x}_{obs}, \mathbf{m}) = \mathbf{1} - \mathbf{m}$$
 under MCAR, and $\mathbf{g}_1(\mathbf{x}_{obs}, \mathbf{m}) = \frac{\mathbf{1} - \mathbf{m}}{\sqrt{\mathbb{P}[\mathbf{m} = 0|\mathbf{x} = \mathbf{x}_{obs}]}}$ under MAR

For the second-order model with missing data, we implement a variance reduction (VR) technique 306 using antithetic sampling (James, 1985; Meng et al., 2021), which involves utilizing two negatively correlated sample vectors centered around \mathbf{x} . The objective function is then formulated as 307

$$\mathcal{L}_{D_{2}SM-VR}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}_{obs},\mathbf{m}} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0},\mathbf{I})} \left[\left\{ \boldsymbol{\psi}(\tilde{\mathbf{x}}_{obs}^{+})^{2} + \boldsymbol{\psi}(\tilde{\mathbf{x}}_{obs}^{-})^{2} + 2\frac{\mathbf{I} - \mathbf{z}\mathbf{z}^{\top}}{\sigma} \odot \boldsymbol{\Psi} \right\} \odot \mathbf{g}_{2}(\mathbf{x}_{obs},\mathbf{m}) \right], \quad (9)$$

311 where the antithetic samples are defined as $\mathbf{x}_{obs}^+ = \mathbf{x}_{obs} + \sigma \mathbf{z}$ and $\mathbf{x}_{obs}^- = \mathbf{x}_{obs} - \sigma \mathbf{z}$. Here, $\boldsymbol{\psi} =$ $\mathbf{s}_2 + \mathbf{s}_1 \mathbf{s}_1^{\top}$, and $\Psi = (\psi(\tilde{\mathbf{x}}_{obs}^+) + \psi(\tilde{\mathbf{x}}_{obs}^-) - 2\psi(\mathbf{x}_{obs}))$. Under the MCAR setting, $\mathbf{g}_2(\mathbf{x}_{obs}, \mathbf{m}) = (\mathbf{1} - \mathbf{m})(\mathbf{1} - \mathbf{m})^{\top}$, while for MAR, $\mathbf{g}_2(\mathbf{x}_{obs}, \mathbf{m}) = \frac{(\mathbf{1} - \mathbf{m})(\mathbf{1} - \mathbf{m})^{\top}}{\sqrt{\mathbb{P}[\mathbf{m}\mathbf{m}^{\top} = 0]|\mathbf{x} = \mathbf{x}]}}$. The formal analysis of the variance reduction are provided in Appendix C.6 and Appendix C.7. 312 313 314 315

316 We perform an empirical analysis to assess the impact of VR on training score models with DSM 317 and D_2SM using incomplete data. The full data is generated from a 2-d Gaussian distribution and we 318 simulate the incomplete data under MCAR, training $s_1(\tilde{x}_{obs}; \theta)$ and $s_2(\tilde{x}_{obs}; \theta)$ using a joint learning 319 objective Eq. (7). Using a sample size of 1000 and a missing ratio of 0.3, In Figure 1, we compare 320 the estimated score and Hessian for the first dimension against the ground truth at noise levels 321 $\sigma = \{0.1, 0.001\}$, both with and without VR. The results indicate that VR is essential for accurate estimation in both DSM and D_2 SM when σ is close to zero, while its importance diminishes at higher 322 values of σ . As σ increases, both methods still achieve reasonable score estimates even without VR. 323 Additionally, when using complete data and varying the missing ratio $\alpha = \{0.1, 0.3, 0.5\}$ with DSM

324 and VR at $\sigma = 0.001$, we observe that the first- and second-order score estimates remain close to the 325 ground truth. While performance degrades as the proportion of missing data increases, the estimates 326 remain generally accurate. 327

328

337

338

339 340 341

342

355

356

357

364 365

366 367

368

369

370

371



Figure 1: Comparison of estimated s_1 and s_2 under different conditions. (a) and (b) show estimates with DSM and D_2SM varying the noise level σ with a fixed missing ratio 0.3. (c) and (d) show estimates with DSM and D_2 SM varying the missing ratio α .

3.2 SCALABILITY AND NUMERICAL STABILITY

343 We show that the proposed method efficiently and accurately estimates second-order scores across 344 different missing ratios, as summarized in Table 1. To achieve this, we generate 10 synthetic datasets 345 with known ground truth, consisting of 100-dimensional correlated multivariate normal distributions 346 that include varying levels of missing data under MCAR mechanism. The covariance matrix for 347 these distributions is constructed using eigenvalues $t \in \{1, 5\}$ to vary degrees of correlation. We 348 evaluate the performance of the estimated s_1 and the diagonal of s_2 by calculating the mean squared error (MSE) between the estimated scores and the ground truth scores derived from the complete 349 data, across various values of σ . Our results indicate that the jointly optimized $s_1(\tilde{x}_{obs}; \theta)$ and 350 $\mathbf{s}_2(\tilde{\mathbf{x}}_{obs}; \boldsymbol{\theta})$ achieve empirical performance close to the ground truth. As previously mentioned, when 351 σ approaches zero, the estimates without VR become unreliable for both the score and Hessian, 352 likely due to convergence issues. Although performance declines with an increase in missing data, 353 the estimates remain reasonable. 354

Table 1: Mean squared error (MSE) between the estimated first-order and second-order scores and the ground truth is evaluated across 5,000 test samples. We vary the noise scales σ and missing ratios α , with each configuration tested using 10 random seeds.

Metho	ds α =	= 0.0	α =	= 0.1	α =	= 0.3	α =	= 0.5
methous	$\sigma = 0.1$	$\sigma=0.01$	$\sigma = 0.1$	$\sigma=0.01$	$\sigma = 0.1$	$\sigma=0.01$	$\sigma = 0.1$	$\sigma=0.01$
\mathbf{s}_1	0.28 ± 0.01	0.42 ± 0.02 0.07 ± 0.00	0.29 ± 0.01 0.00 ± 0.00	0.42 ± 0.02 0.00 ± 0.00	0.32 ± 0.01 0.12 ± 0.00	0.44 ± 0.01 0.15 ± 0.01	0.37 ± 0.01 0.22 ± 0.00	0.44 ± 0.02 0.27 ± 0.01
$\mathbf{s}_1(\mathbf{v}\mathbf{K})$ \mathbf{s}_2	0.07 ± 0.00 0.16 ± 0.02	15.42 ± 0.47	0.09 ± 0.00 0.16 ± 0.02	0.09 ± 0.00 27.42 ± 2.34	0.13 ± 0.00 0.16 ± 0.02	0.15 ± 0.01 29.74 ± 3.35	0.23 ± 0.00 0.17 ± 0.03	0.27 ± 0.01 26.08 ± 2.03
$\mathbf{s}_2(VR)$	0.04 ± 0.00	0.05 ± 0.00	0.04 ± 0.00	0.04 ± 0.00	0.04 ± 0.00	0.05 ± 0.00	0.05 ± 0.00	0.06 ± 0.01

SAMPLING WITH MISSING DATA VIA SECOND-ORDER SCORE MODELS

In this section, we illustrate how our second-order score model $s_2(\tilde{x}; \theta)$, trained with missing data, enhances both the speed of sample generation and the quality of the synthetic data. We demonstrate the effectiveness of the proposed model through simulations and real-world datasets containing missing values, comparing its performance against various baseline methods.

Langevin dynamics. Langevin dynamics samples data from $p_{\text{data}}(\mathbf{x})$ by utilizing the first-order 372 score function $s_1(x)$ (Bussi & Parrinello, 2007; Song & Ermon, 2019). Starting with a prior distri-373 bution $\pi(\mathbf{x})$, a fixed step size $\epsilon > 0$, and an initial value $\tilde{\mathbf{x}}_0 \sim \pi(\mathbf{x})$, Langevin dynamics iteratively 374 updates the samples as follows:

 $\tilde{\mathbf{x}}_{t+1} = \tilde{\mathbf{x}}_t + \frac{1}{2}\epsilon \mathbf{s}_1(\tilde{\mathbf{x}}_t) + \sqrt{\epsilon}\mathbf{z}_t,$ (10)

where $\mathbf{z}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ represents Gaussian noise.

Ozaki Sampling. Following Meng et al. (2021), Ozaki discretization improves data synthesis by integrating second-order information from $s_2(x)$ to precondition the sampling process. The updates for Ozaki sampling are performed as follows:

381 382

384

 $\tilde{\mathbf{x}}_t = \tilde{\mathbf{x}}_{t-1} + \mathbf{M}_{t-1}\mathbf{s}_1(\tilde{\mathbf{x}}_{t-1}) + \Sigma_{t-1}^{1/2}\mathbf{z}_t,$ (11)

where $\mathbf{z}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \mathbf{M}_{t-1} = e^{\epsilon \mathbf{s}_2(\tilde{\mathbf{x}}_{t-1})} - \mathbf{I}$, and $\Sigma_{t-1} = (e^{2\epsilon \tilde{\mathbf{x}}_{t-1}} - \mathbf{I})\mathbf{s}_2(\tilde{\mathbf{x}}_{t-1})^{-1}$.

Illustration. We use the Swiss-Roll dataset to demonstrate the effectiveness of Ozaki Sampling, focusing on its speed and quality of data generation through second-order information. Both methods employ a step size of $\epsilon = 0.005$ and a missing ratio of 0.5 under the MCAR missing mechanism. As shown in Figure 2, Ozaki Sampling generates comparable data to Langevin dynamics with fewer iterations, and resulting in data that is more concentrated around the original distribution, while Langevin dynamics yields noisier and more dispersed results.

Following Kim et al. (2022); Ouyang et al. (2023), we conduct experiments on a simulated Bayesian Network dataset and a real Census dataset (Kohavi, 1996) to illustrate the efficiency and effectiveness of the data generated by our proposed model trained on missing data.

Baselines. We evaluate the proposed method using both Langevin and Ozaki sampling against several baseline techniques for synthetic data generation on datasets with missing values. Specifically, we implement a vanilla DSM model that (1) removes rows with missing values, and (2) uses mean imputation for missing values in each column. Additionally, we include STaSy (Kim et al., 2022), a state-of-the-art score-based model, which significantly outperforms other approaches for tabular data. As STaSy requires complete datasets for training, we apply mean imputation to handle any missing values in the training data.

401 **Metrics.** Following Kim et al. (2022); Ouyang et al. (2023), we employ two criteria, *fidelity* and *utility*, to assess the quality of the generated synthetic tabular data. For evaluating *fidelity*, 402 we utilize the model-agnostic library SDMetrics. The result ranges from 0 to 100%. A higher 403 score indicates better overall quality of the synthetic data. To measure *utility*, we adopt the same 404 pipeline as Kim et al. (2022), training various models—including Decision Tree, AdaBoost, Logistic 405 Regression, MLP Classifier, Random Forest, and XGBoost on the synthetic data and test them with 406 real data. Our primary metric is classification accuracy, and we also report AUROC and Weighted-407 F1 scores in the Appendix E.4. All experimental results are based on three repetitions. 408

Results. Figure 5 and Table 2 demonstrate the effectiveness of the proposed method on both the simulated Bayesian Network dataset and the Census data, showing superior performance in terms of fidelity and utility compared to other baselines. Specifically, these results confirm that the impute-then-generate approach introduces bias, whereas directly learning from missing data significantly improves the performance of the generative model. Furthermore, the advantages of the proposed model become more pronounced as the missing ratios increase. Additional details and results of the experiments can be found in Appendix E.

Table 2: Fidelity (SDMetric) and Utility (Accuracy) evaluation of MissScore using Langevin and
Ozaki samplings, along with other baselines, on the Census dataset with missing ratio 0.3 under
MCAR.

	Legenvin	Okazi	DSM-delete	DSM-mean	STaSy-mean
Fiedility	86%	88%	73%	77%	82%
Utility	80%	$\mathbf{81\%}$	70%	75%	77%

423 424 425

426

5 CAUSAL DISCOVERY WITH MISSING DATA VIA SECOND-ORDER SCORES

427 **Background.** Causal discovery aims to identify causal relationship from purely observational data. 428 However, the task is ill-posed without additional assumptions. Assuming an additive noise model 429 (ANM) allows for the identification of causal structures. In this context, consider the ANM defined 430 as $x_i = f_i(x_{PA_i}) + z_i$, where f_i is a nonlinear function and z_i is a Gaussian noise. Rolland et al. 431 (2022) proposed an order-based algorithm that uses the second order score of an ANM with a probability distribution $p_{data}(\mathbf{x})$ to identify leaf nodes, and iteratively determine the topological order of



Figure 2: Sampling a Swiss Roll dataset with a step size of 0.005 using Langevin dynamics and Ozaki sampling. Ozaki sampling demonstrates faster convergence and better data quality compared to Langevin dynamics under MCAR with a missing ratio of 0.5. Figure (a) displays the dataset; Figures (b)-(d) show the results from Langevin dynamics; and Figures (e)-(h) present the results from Ozaki sampling. The numbers in parentheses indicate the number of iterations sampling taken.



Figure 3: Fidelity evaluation of MissScore using Langevin and Ozaki samplings, along with other baselines, on the Bayesian Network dataset varies with missing ratio $\alpha = \{0.1, 0.3, 0.5, 0.7, 0.9\}$ under different missing mechanisms.

447

448

449

454

455

456

457

458

459

460

the variables. However, the computation of the Hessian requires complete data, which poses challenges in real-world scenarios such as clinical trials, and biology, where missing data is common.

468 A straightforward approach to address missing data problem is to first impute the incomplete entries 469 using off-the-shelf imputation methods and then apply existing causal discovery methods. However, 470 this two-step approach can be suboptimal, as the imputation process may introduce bias for modeling the underlying data distribution. Our method mitigates this issue by directly training a second-471 order model with incomplete data, thereby reducing potential bias. Since the Hessian only provides 472 information about variable order, we adopt a strategy similar to Rolland et al. (2022); Sanchez et al. 473 (2022), first computing the topological order and then using CAM pruning to derive the final directed 474 acyclic graph (DAG) (Bühlmann et al., 2014). 475

Baselines. We utilize the MissForest imputation method to address missing data, followed by the implementation of DiffAN (Sanchez et al., 2022) (termed MissDiffAN) and DAGMA (Bello et al., 2022) (termed MissForest) for structure learning. DiffAN serves as a diffusion-based adaptation of the approach proposed by Rolland et al. (2022), ensuring a fair comparison. Furthermore, we compare MissScore with MissDAG (Gao et al., 2022) and MissOTM (Vo et al., 2024), both of which are prominent methods that have shown superior performance in causal discovery with missing data relative to various other baselines.

483 Metrics. All quantitative results are averaged over 10 random initiazations. For comparing the
 484 estimated DAG with the ground-truth one, we report commonly used metrics: Order Divergence
 485 and Structural Hamming Distance (SHD). Order Divergence measures the number of errors in the
 ordering, while SHD indicates the minimum number of edge additions, deletions, and reversals

needed to convert the recovered DAG into the true one. Lower values for both metrics are preferred.
 Order Divergence is calculated only for MissScore and MissDiffAN, as these are the only order based methods.

Simulations. We simulate synthetic datasets generating a ground-truth DAG from the graph model Erdős–Rényi (ER). Each function f_i is constructed from a multi-layer perceptron (MLP) and a multiple index model (MIM) with random coefficients. We consider a general scenario of non-euqal variances, sampling 1000 observations according to all missing mechansims: MCAR, MAR, MNAR at 10% and 30% missing rates and complete data. In the main text, we report the SHD and runtime in MCAR cases with missing ratio 0.1 varies with number of dimensions, using Gaussian noise. Specifically, the number of edges is set equal to the dimensionality.

496 **Results.** In Figure 4, our approach demonstrates comparable performance to the state-of-the-497 art methods MissDAG and MissOTM, although it shows slightly lower performance in the high-498 dimensional scenario with d = 50. This may be attributed to the challenges of training a denois-499 ing score matching model with a limited number of samples. However, MissScore significantly 500 reduces computation time compared to MissOTM and MissDAG, offering an efficient alternative 501 while maintaining similar results. With MissForest, our performance is similar, possibly due to the additional constraints enforced during training by DAGMA. However, when comparing our method 502 503 to those that rely on imputation followed by causal discovery approach, we find that our approach outperforms MissDiffAN. This suggests that imputation can introduce bias in downstream tasks. 504 Additional results for various missing ratios, mechanisms, and order divergences are provided, along 505 with further experimental details, in Appendix F.4. 506



Figure 4: Data is generated under MCAR with a missing ratio of 0.1, varying dimensions $d = \{10, 20, 50\}$ using ER graph model. The sample size is 1000, and f_i corresponds to an MLP. Left: SHD; Right: Runtime. The shaded area indicates 95% confidence.

6 CONCLUSION

519

520

525 526

In this work, we introduce a method to directly estimate higher-order data density scores in the 527 presence of missing data, extending denoising score matching to accommodate scores of any order 528 across different missing data mechanisms. Our approach directly handles missing data without re-529 lying on imputation or deletion. Empirical results demonstrate that models trained with this method 530 estimate second-order scores more efficiently and accurately. Moreover, second-order models en-531 hance the sampling quality of Langevin dynamics via Ozaki discretization in missing data scenarios. 532 Our proposed causal discovery method for incomplete data scales effectively with dimensionality 533 and achieves performance comparable to state-of-the-art approaches. However, the effectiveness 534 of this approach diminishes in low-noise environments without variance reduction, particularly for second-order scores, which poses challenges in low-noise or high-missingness scenarios. This limi-536 tation also extends to downstream tasks, such as causal discovery and sampling, where performance 537 tends to decline slightly with increasing dimensionality. Future directions include exploring scorebased models with constraints for causal discovery in missing data contexts, applying higher-order 538 score estimation to a broader range of applications like image and time-series data, and further investigating the use of denoising score matching for handling MNAR data directly.

540 REFERENCES

548

551

552

553

554

555

582

583

- Kevin Bello, Bryon Aragam, and Pradeep Ravikumar. Dagma: Learning dags via m-matrices and a
 log-determinant acyclicity characterization. *Advances in Neural Information Processing Systems*,
 35:8226–8239, 2022.
- Dimitris Bertsimas, Colin Pawlowski, and Ying Daisy Zhuo. From predictive methods to missing data imputation: an optimization approach. *Journal of Machine Learning Research*, 18(196): 1–39, 2018.
- Peter Bühlmann, Jonas Peters, and Jan Ernest. Cam: Causal additive models, high-dimensional order search and penalized regression. *The Annals of Statistics*, 42(6):2526–2556, 2014.
 - Giovanni Bussi and Michele Parrinello. Accurate sampling using langevin dynamics. *Physical Review E—Statistical, Nonlinear, and Soft Matter Physics*, 75(5):056707, 2007.
 - Hyungjin Chung and Jong Chul Ye. Score-based diffusion models for accelerated mri. *Medical image analysis*, 80:102479, 2022.
- Arnak S Dalalyan and Avetik Karagulyan. User-friendly guarantees for the langevin monte carlo with inaccurate gradient. *Stochastic Processes and their Applications*, 129(12):5278–5311, 2019.
- 559 Synthetic Data Metrics. DataCebo, Inc., 2023.
- Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the royal statistical society: series B (methodological)*, 39(1): 1–22, 1977.
- Bradley Efron. Tweedie's formula and selection bias. Journal of the American Statistical Association, 106(496):1602–1614, 2011.
- Cristóbal Esteban, Stephanie L Hyland, and Gunnar Rätsch. Real-valued (medical) time series generation with recurrent conditional gans. *arXiv preprint arXiv:1706.02633*, 2017.
- Alexander Gain and Ilya Shpitser. Structure learning under missing data. In *International conference* on probabilistic graphical models, pp. 121–132. Pmlr, 2018.
- Erdun Gao, Ignavier Ng, Mingming Gong, Li Shen, Wei Huang, Tongliang Liu, Kun Zhang, and Howard Bondell. Missdag: Causal discovery in the presence of missing data with continuous additive noise models. *Advances in Neural Information Processing Systems*, 35:5024–5038, 2022.
- Lovedeep Gondara and Ke Wang. Mida: Multiple imputation using denoising autoencoders. In *Advances in Knowledge Discovery and Data Mining: 22nd Pacific-Asia Conference, PAKDD 2018, Melbourne, VIC, Australia, June 3-6, 2018, Proceedings, Part III 22*, pp. 260–272. Springer, 2018.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
 Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
 - Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Aapo Hyvärinen and Peter Dayan. Estimation of non-normalized statistical models by score match *Journal of Machine Learning Research*, 6(4), 2005.
- Sebastian Jäger, Arndt Allhorn, and Felix Bießmann. A benchmark for data imputation methods.
 Frontiers in big Data, 4:693674, 2021.
- BAP James. Variance reduction techniques. Journal of the Operational Research Society, 36(6): 525–530, 1985.
- Jayoung Kim, Chaejeong Lee, and Noseong Park. Stasy: Score-based tabular data synthesis. *arXiv* preprint arXiv:2210.04018, 2022.

594 595	Diederik P Kingma. Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114, 2013.
595	Ron Kohavi Census Income LICI Machine Learning Repository 1996 DOI:
597	https://doi.org/10.24432/C5GP7S.
598	
599	Haoying Li, Yifan Yang, Meng Chang, Shiqi Chen, Huajun Feng, Zhihai Xu, Qi Li, and Yueting
600	Chen. Srdiff: Single image super-resolution with diffusion probabilistic models. <i>Neurocomputing</i> ,
601	4/9:47-59, 2022.
602	Steven Cheng-Xian Li, Bo Jiang, and Benjamin Marlin. Misgan: Learning from incomplete data
603	with generative adversarial networks. arXiv preprint arXiv:1902.09599, 2019.
604	
605	Qiang Liu, Jason Lee, and Michael Jordan. A kernelized stein discrepancy for goodness-of-fit tests.
606	In International conference on machine learning, pp. 270–284. PNILK, 2016.
607	Wenqin Liu, Biwei Huang, Erdun Gao, Qiuhong Ke, Howard Bondell, and Mingming Gong. Causal
608	discovery with mixed linear and nonlinear additive noise models: A scalable approach. In Causal
609	Learning and Reasoning, pp. 1237–1263. PMLR, 2024.
610	Chang Lu, Kaiwan Zhang, Fan Rao, Jianfai Chan, Changyuan Li, and Jun Zhu, Mayimum likalihood
611	training for score-based diffusion odes by high order denoising score matching. In <i>International</i>
612	Conference on Machine Learning, pp. 14429–14460. PMLR, 2022.
613	
614	Chenlin Meng, Yang Song, Wenzhe Li, and Stefano Ermon. Estimating high order gradients of the
615	data distribution by denoising. Advances in Neural Information Processing Systems, 34:25359–
616	25509, 2021.
617	Pablo Morales-Alvarez, Wenbo Gong, Angus Lamb, Simon Woodhead, Simon Peyton Jones, Nick
610	Pawlowski, Miltiadis Allamanis, and Cheng Zhang. Simultaneous missing value imputation and
620	structure learning with groups. Advances in Neural Information Processing Systems, 35:20011-
621	20024, 2022.
622	Boris Muzellec, Julie Josse, Claire Boyer, and Marco Cuturi, Missing data imputation using optimal
623	transport. In International Conference on Machine Learning, pp. 7130–7140. PMLR, 2020.
624	
625	Yidong Ouyang, Liyan Xie, Chongxuan Li, and Guang Cheng. Missdiff: Training diffusion models
626	on tabular data with missing values. arXiv preprint arXiv:2307.00467, 2023.
627	Noseong Park, Mahmoud Mohammadi, Kshitij Gorde, Sushil Jajodia, Hongkyu Park, and Youngmin
628	Kim. Data synthesis based on generative adversarial networks. arXiv preprint arXiv:1806.03384,
629	2018.
630	George Paterakis Stefanos Fafalios Paulos Charonyktakis Vassilis Christophidas and Joannis
631	Tsamardinos Do we really need imputation in automl predictive modeling? ACM Transactions
632	on Knowledge Discovery from Data, 18(6):1–64, 2024.
633	
634	Jason Poulos and Rafael Valle. Missing data imputation for supervised learning. <i>Applied Artificial</i>
635	Intelligence, 32(2):186–196, 2018.
636	Paul Rolland, Volkan Cevher, Matthäus Kleindessner, Chris Russell, Dominik Janzing, Bernhard
630	Schölkopf, and Francesco Locatello. Score matching enables causal discovery of nonlinear addi-
630	tive noise models. In International Conference on Machine Learning, pp. 18741–18753. PMLR,
640	2022.
641	Donald B Rubin. Inference and missing data. <i>Biometrika</i> 63(3):581–592, 1976
642	2
643	Sotirios Sabanis and Ying Zhang. Higher order langevin monte carlo algorithm. 2019.
644	Pedro Sanchez Xiao Liu Alicon O O'Neil and Sotirios A Teaffaris Diffusion models for equal
645	discovery via topological ordering arXiv preprint arXiv:2210.06201, 2022
646	
647	Shaun R Seaman and Ian R White. Review of inverse probability weighting for dealing with missing data. <i>Statistical methods in medical research</i> , 22(3):278–295, 2013.

648	Tolou Shadbahr, Michael Roberts, Jan Stanczuk, Julian Gilbey, Philip Teare, Sören Dittmer,
649	Matthew Thorpe, Ramon Viñas Torné, Evis Sala, Pietro Lió, et al. The impact of imputation qual-
650	ity on machine learning classifiers for datasets with missing values. Communications Medicine,
651	3(1):139, 2023.
652	

- Ilya Shpitser. Consistent estimation of functions of data missing non-monotonically and not at random. Advances in Neural Information Processing Systems, 29, 2016.
- Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution.
 Advances in neural information processing systems, 32, 2019.
 - Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. Advances in neural information processing systems, 33:12438–12448, 2020.
- Yang Song and Diederik P Kingma. How to train your energy-based models. *arXiv preprint arXiv:2101.03288*, 2021.
- 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*, pp. 574–584. PMLR, 2020.
 - Yang Song, Liyue Shen, Lei Xing, and Stefano Ermon. Solving inverse problems in medical imaging with score-based generative models. *arXiv preprint arXiv:2111.08005*, 2021.
 - Nicolas Städler and Peter Bühlmann. Missing values: sparse inverse covariance estimation and an extension to sparse regression. *Statistics and Computing*, 22:219–235, 2012.
- Ruibo Tu, Cheng Zhang, Paul Ackermann, Karthika Mohan, Hedvig Kjellström, and Kun Zhang.
 Causal discovery in the presence of missing data. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 1762–1770. Pmlr, 2019.
 - Stef Van Buuren and Karin Groothuis-Oudshoorn. mice: Multivariate imputation by chained equations in r. *Journal of statistical software*, 45:1–67, 2011.
- Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011.
 - Vy Vo, He Zhao, Trung Le, Edwin V Bonilla, and Dinh Phung. Optimal transport for structure learning under missing data. *arXiv preprint arXiv:2402.15255*, 2024.
- Zhenhua Wang, Olanrewaju Akande, Jason Poulos, and Fan Li. Are deep learning models superior
 for missing data imputation in large surveys? evidence from an empirical comparison. *arXiv preprint arXiv:2103.09316*, 2021.
 - Jeffrey M Wooldridge. Inverse probability weighted estimation for general missing data problems. *Journal of econometrics*, 141(2):1281–1301, 2007.
 - Jinsung Yoon, James Jordon, and Mihaela Schaar. Gain: Missing data imputation using generative adversarial nets. In *International conference on machine learning*, pp. 5689–5698. PMLR, 2018.
 - Xun Zheng, Bryon Aragam, Pradeep K Ravikumar, and Eric P Xing. Dags with no tears: Continuous optimization for structure learning. *Advances in neural information processing systems*, 31, 2018.

Appendix

Table of Contents

A	Related Work	15
B	Additional Discussions	15
С	Proofs	16
	C.1 Proof of Theorem 3	16
	C.2 Proof of Theorem 4	17
	C.3 Proof of Theorem 7	19
	C.4 Proof of Theorem 5	19
	C.5 Proof of Theorem 6	20
	C.6 Proof of Equation 8	20
	C.7 Proof of Equation 9	21
E	Additional Information on Sampling	24
	E.1 Dataset Description and Processing	24
	E.2 Evaluation Methods	25
	E.3 Model Architecture	25
	E.4 Experimental Results	25
	E.5 Analysis of Sampling Results Using Census Data	26
F	Additional information on Causal Discovery	28
	F.1 Related Work	28
	F.2 Evaluation Metrics	30
	F.3 Model Architecture	30
		20
	F.4 Experimental Results	50

756 A RELATED WORK

758

Missing Data Addressing missing values in training data has been widely studied, leading to the 759 development of numerous imputation techniques. Traditional approaches often involve either re-760 moving rows or columns with missing entries or imputing missing values by substituting them with 761 the mean of observed values for a given feature. More advanced methods apply machine learning 762 and deep generative models for imputation, such as those explored in prior work (Muzellec et al., 763 2020; Van Buuren & Groothuis-Oudshoorn, 2011; Bertsimas et al., 2018). However, imputation can 764 reduce data diversity and introduce biases in downstream tasks. For instance, Ouyang et al. (2023) highlight that in an "impute-then-generate" pipeline, imputation may impair generation quality by 765 introducing biases. In contrast, "impute-then-predict" pipelines often yield better predictive accu-766 racy for classification and regression tasks, as evidenced by works like (Poulos & Valle, 2018; Jäger 767 et al., 2021; Shadbahr et al., 2023; Paterakis et al., 2024). 768

769 Additionally, imputation methods often fail to account for uncertainty in missing data. Multiple imputation (MI) is a valuable alternative that incorporates this uncertainty. While Van Buuren & 770 Groothuis-Oudshoorn (2011) introduce the widely-used MICE algorithm for MI, MI techniques can 771 also be applied with deep learning models such as GAIN and MIDA (Gondara & Wang, 2018; Wang 772 et al., 2021) to produce multiple plausible datasets by re-running models with different random ini-773 tializations, thereby capturing imputation uncertainty. A comparative study by Wang et al. (2021) 774 indicates that MI with classification trees (like MICE) often outperforms deep learning-based impu-775 tation, especially in survey data. 776

Beyond imputation, recent research has explored models that directly learn from incomplete data or synthesize complete datasets using generative architectures such as GANs and VAEs (Li et al., 2019; Yoon et al., 2018; Park et al., 2018). These approaches typically require auxiliary networks and rely on assumptions about the missing data mechanism. In addition, unique challenges posed by tabular data in these contexts remain underexplored. Recent advancements include applying score-based models with self-paced learning and fine-tuning strategies (Kim et al., 2022), as well as diffusion models, which learn from incomplete data through first-order derivatives (Ouyang et al., 2023).

Score Matching. Score matching estimates the gradient of the log-density (known as the score 784 function) of a data distribution, making them highly effective in modeling complex distributions 785 (Hyvärinen & Dayan, 2005). A prominent technique in this family is Denoising Score Matching, 786 which learns the score of a perturbed version of the target distribution by minimizing a regres-787 sion loss (Vincent, 2011). This method is also widely used in training Denoising Diffusion Models 788 and has found success in tasks such as image and audio generation (Ho et al., 2020). While first-789 order score matching is prevalent, higher-order derivatives offer richer information about the data 790 distributions by capturing its local curvature. Although these can be derived using automatic differ-791 entiation of a learned score model, such method is computationally expensive and prone to error in 792 high-dimensional settings. To address these, recent work by Meng et al. (2021) extends denoising score matching to estimate higher-order derivatives by leveraging Tweedie's formula, which relates 793 higher-order moments of the distribution to its scores. Lu et al. (2022) further shows that the negative 794 likelihood of the ODE can be bounded by controlling the high-order score matching errors. 795

796 797

798

B ADDITIONAL DISCUSSIONS

Differences between MissScore and MissDiff MissScore differs from MissDiff by focusing on
 learning high-order score with missing data, specifically leveraging the Hessian to enhance down stream performance. In contrast, MissDiff emphasizes unbiased data synthesis with a first-order
 diffusion-based approach. Additionally, MissScore introduces an oracle estimator under the MAR
 mechanism and introduce a novel causal discovery method that operates effectively in the presence
 of missing data.

805

Scalability and numerical stability for MAR We have included additional synthetic results for
 the MAR mechanism in Table 3, following the same setup as described in Section 3.2. The results
 exhibit a similar pattern, though slightly worse than MCAR, which may be attributed to the inverse
 probability, potentially increasing variance and leading to less stable estimations. Nevertheless,
 the estimates remain reasonable and demonstrate empirical performance closely aligned with the

ground truth, along with variance reduction. Additionally, to illustrate the impact of potential model misspecification by the logistic regression model in MAR, we conduct an experiment comparing the ground truth p with the estimated \hat{p} from the logistic regression model on the same synthetic dataset, with varying missing ratios. Comparing the performance in Tables 3 and 4, potential model misspecification in the logistic regression model does impact the estimation of missing probability. However, this effect remains within a reasonable range in the variance-reduced version.

Table 3: Mean squared error (MSE) between the estimated first-order and second-order scores and the ground truth is evaluated across 5,000 test samples. We vary the noise scales σ and missing ratios α , with each configuration tested using 10 random seeds. MAR with estimated missing probability \hat{p} by logistic model.

Methods	$\alpha = 0.0$ (Complete data)		$\alpha = 0.1$		lpha=0.3		lpha=0.5	
methous	$\sigma = 0.1$	$\sigma=0.01$	$\sigma = 0.1$	$\sigma = 0.01$	$\sigma = 0.1$	$\sigma = 0.01$	$\sigma = 0.1$	$\sigma = 0.01$
\mathbf{s}_1	0.28 ± 0.01	0.42 ± 0.02	0.32 ± 0.00	0.51 ± 0.02	0.36 ± 0.02	0.52 ± 0.00	0.45 ± 0.01	0.56 ± 0.03
$s_1(VR)$	0.07 ± 0.00	0.07 ± 0.00	0.11 ± 0.02	0.13 ± 0.01	0.15 ± 0.01	0.16 ± 0.02	0.22 ± 0.02	0.24 ± 0.04
s_2	0.16 ± 0.02	15.42 ± 0.47	0.28 ± 0.02	31.22 ± 1.08	0.30 ± 0.04	34.24 ± 5.14	0.36 ± 0.04	35.41 ± 1.03
$s_2(VR)$	0.04 ± 0.00	0.05 ± 0.00	0.04 ± 0.00	0.05 ± 0.00	0.06 ± 0.00	0.06 ± 0.00	0.08 ± 0.01	0.07 ± 0.00

Table 4: Mean squared error (MSE) between the estimated first-order and second-order scores and the ground truth is evaluated across 5,000 test samples. We vary the noise scales σ and missing ratios α , with each configuration tested using 10 random seeds. MAR with ground truth missing probability p.

Method	$\alpha = 0.0 (C$	$\alpha = 0.0$ (Complete data)		lpha=0.1		lpha=0.3		= 0.5
wiedhod	$\sigma = 0.1$	$\sigma=0.01$	$\sigma = 0.1$	$\sigma=0.01$	$\sigma = 0.1$	$\sigma=0.01$	$\sigma = 0.1$	$\sigma=0.01$
\mathbf{s}_1	0.28 ± 0.01	0.42 ± 0.02	0.24 ± 0.01	0.42 ± 0.02	0.27 ± 0.01	0.41 ± 0.01	0.33 ± 0.01	0.43 ± 0.04
$s_1(VR)$	0.07 ± 0.00	0.07 ± 0.00	0.06 ± 0.00	0.06 ± 0.00	0.09 ± 0.00	0.10 ± 0.00	0.19 ± 0.00	0.22 ± 0.01
\mathbf{s}_2	0.16 ± 0.02	15.42 ± 0.47	0.24 ± 0.01	16.66 ± 4.34	0.27 ± 0.01	41.47 ± 7.28	0.33 ± 0.01	39.89 ± 4.03
$s_2(VR)$	0.04 ± 0.00	0.05 ± 0.00	0.02 ± 0.00	0.03 ± 0.00	0.03 ± 0.00	0.03 ± 0.00	0.05 ± 0.00	0.05 ± 0.00

С PROOFS

C.1 PROOF OF THEOREM 3

We know that

$$\begin{split} \mathbb{E}_{\mathbf{x},\mathbf{m}} \mathbb{E}_{\tilde{\mathbf{x}}|\mathbf{x},\mathbf{m}} \left[\left\| \left\{ \mathbf{s}(\tilde{\mathbf{x}};\boldsymbol{\theta}) + \frac{1}{\sigma^2} (\tilde{\mathbf{x}} - \mathbf{x}) \right\} \odot (\mathbf{1} - \mathbf{m}) \right\|_2^2 \right] \\ = \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x},\mathbf{m}} \left[\left\| \left\{ \mathbf{s}(\tilde{\mathbf{x}};\boldsymbol{\theta}) + \frac{1}{\sigma^2} (\tilde{\mathbf{x}} - \mathbf{x}) \right\} \odot (\mathbf{1} - \mathbf{m}) \right\|_2^2 \right] \\ = \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x}} \mathbb{E}_{m|\tilde{\mathbf{x}},\mathbf{x}} \left[\left\| \left\{ \mathbf{s}(\tilde{\mathbf{x}};\boldsymbol{\theta}) + \frac{1}{\sigma^2} (\tilde{\mathbf{x}} - \mathbf{x}) \right\} \odot (\mathbf{1} - \mathbf{m}) \right\|_2^2 \right] \\ = \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x}} \left[\left\| \left\{ \mathbf{s}(\tilde{\mathbf{x}};\boldsymbol{\theta}) + \frac{1}{\sigma^2} (\tilde{\mathbf{x}} - \mathbf{x}) \right\} \odot \sqrt{\mathbb{E}_{\mathbf{m}|\tilde{\mathbf{x}},\mathbf{x}} (\mathbf{1} - \mathbf{m})} \right\|_2^2 \right] \\ = \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x}} \left[\left\| \left\{ \mathbf{s}(\tilde{\mathbf{x}};\boldsymbol{\theta}) + \frac{1}{\sigma^2} (\tilde{\mathbf{x}} - \mathbf{x}) \right\} \odot \sqrt{\mathbf{w}_1} \right\|_2^2 \right]. \end{split}$$

Denote $q_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x}) = \mathcal{N}(\tilde{\mathbf{x}}; \mathbf{x}, \sigma^2)$, we only need to show there exists some constant C independent of θ such that

$$\mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x}}\left[\left\|\left\{\mathbf{s}(\tilde{\mathbf{x}};\boldsymbol{\theta}) - \nabla_{\tilde{\mathbf{x}}}q_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})\right\} \odot \sqrt{\mathbf{w}_{1}}\right\|_{2}^{2}\right] = \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x}}\left[\left\|\left\{\mathbf{s}(\tilde{\mathbf{x}};\boldsymbol{\theta}) - \mathbf{s}(\tilde{\mathbf{x}})\right\} \odot \sqrt{\mathbf{w}_{1}}\right\|_{2}^{2}\right] + C$$
(12)

For simplicity, we consider the case d = 1. For the right hand side of equation 12, $\lim_{T \to 0} \left[\left[c^2(\tilde{x}; \theta) \right] - 2\mathbb{R}_{-1} \left[s(\tilde{x}; \theta) s(\tilde{x}) \right] \right] + C.$

$$\mathbf{R.H.S} = w_1 \cdot \left\{ \mathbb{E}_{\tilde{x},x} \left[s^2(\tilde{x};\theta) \right] - 2\mathbb{E}_{\tilde{x},x} \left[s(\tilde{x};\theta)s(\tilde{x}) \right] \right\} + \frac{1}{2} \left\{ \mathbb{E}_{\tilde{x},x} \left[s(\tilde{x};\theta)s($$

For the left hand side of equation 12,

$$\mathbf{L.H.S} = w_1 \cdot \left\{ \mathbb{E}_{\tilde{x},x} \left[s^2(\tilde{x};\theta) \right] - 2\mathbb{E}_{\tilde{x},x} \left[s(\tilde{x};\theta) \nabla_{\tilde{x}} \log q_\sigma(\tilde{x}|x) \right] \right\} + C$$

Hence, we only need to show

$$\mathbb{E}_{\tilde{x},x}\left[s(\tilde{x};\theta)s(\tilde{x})\right] = \mathbb{E}_{\tilde{x},x}\left[s(\tilde{x};\theta)\nabla_{\tilde{x}}\log q_{\sigma}(\tilde{x}|x)\right].$$
(13)

We have,

$$\begin{split} \mathbb{E}_{\tilde{x},x} \left[s(\tilde{x};\theta) s(\tilde{x}) \right] &= \int s(\tilde{x};\theta) \nabla_{\tilde{x}} p_{\tilde{x}}(\tilde{x}) d\tilde{x} \\ &= \int s(\tilde{x};\theta) \nabla_{\tilde{x}} \left(\int p(x) q_{\sigma}(\tilde{x}|x) dx \right) d\tilde{x} \\ &= \int \int s(\tilde{x};\theta) p(x) \nabla_{\tilde{x}} \left(q_{\sigma}(\tilde{x}|x) \right) dx d\tilde{x} \\ &= \int \int s(\tilde{x};\theta) p(x) q_{\sigma}(\tilde{x}|x) \nabla_{\tilde{x}} \left(\log q_{\sigma}(\tilde{x}|x) \right) dx d\tilde{x} \\ &= \int \int s(\tilde{x};\theta) p_{\tilde{x},x}(\tilde{x},x) \nabla_{\tilde{x}} \left(\log q_{\sigma}(\tilde{x}|x) \right) dx d\tilde{x} = \mathbb{E}_{\tilde{x},x} \left[s(\tilde{x};\theta) \nabla_{\tilde{x}} \log q_{\sigma}(\tilde{x}|x) \right] \end{split}$$

Hence, we get our desired result.

C.2 PROOF OF THEOREM 4

We have

$$\begin{split} & \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x},\mathbf{m}} \left[\left\| \left\{ \mathbf{s}_{2}(\tilde{\mathbf{x}};\boldsymbol{\theta}) + \mathbf{s}_{1}(\tilde{\mathbf{x}})\mathbf{s}_{1}^{\top}(\tilde{\mathbf{x}}) + \frac{\mathbf{I} - \mathbf{z}\mathbf{z}^{\top}}{\sigma^{2}} \right\} \odot \left\{ (\mathbf{1} - \mathbf{m})(\mathbf{1} - \mathbf{m})^{\top} \right\} \right\|_{2}^{2} \right] \\ & = \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x}} \left[\left\| \left\{ \mathbf{s}_{2}(\tilde{\mathbf{x}};\boldsymbol{\theta}) + \mathbf{s}_{1}(\tilde{\mathbf{x}})\mathbf{s}_{1}^{\top}(\tilde{\mathbf{x}}) + \frac{\mathbf{I} - \mathbf{z}\mathbf{z}^{\top}}{\sigma^{2}} \right\} \odot \sqrt{\mathbb{E}_{\mathbf{m}|\tilde{\mathbf{x}},\mathbf{x}}(\mathbf{1} - \mathbf{m})(\mathbf{1} - \mathbf{m})^{\top}} \right\|_{2}^{2} \right] \\ & = \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x}} \left[\left\| \left\{ \mathbf{s}_{2}(\tilde{\mathbf{x}};\boldsymbol{\theta}) + \mathbf{s}_{1}(\tilde{\mathbf{x}})\mathbf{s}_{1}^{\top}(\tilde{\mathbf{x}}) + \frac{\mathbf{I} - \mathbf{z}\mathbf{z}^{\top}}{\sigma^{2}} \right\} \odot \sqrt{\mathbf{w}_{2}} \right\|_{2}^{2} \right] \end{split}$$

It is sufficient to show

$$\mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x}}\left[\left\|\left\{\mathbf{s}_{2}(\tilde{\mathbf{x}};\boldsymbol{\theta})+\mathbf{s}_{1}(\tilde{\mathbf{x}})\mathbf{s}_{1}^{\top}(\tilde{\mathbf{x}})+\frac{\mathbf{I}-\mathbf{z}\mathbf{z}^{\top}}{\sigma^{2}}\right\}\odot\sqrt{\mathbf{w}_{2}}\right\|_{2}^{2}\right]$$

$$=\mathbb{E}_{\mathbf{x}}\mathbb{E}_{\tilde{\mathbf{x}}\mid\mathbf{x}}\left[\left\|\left\{\mathbf{s}_{2}(\tilde{\mathbf{x}};\boldsymbol{\theta})-\mathbf{s}_{2}(\tilde{\mathbf{x}})\right\}\odot\sqrt{\mathbf{w}_{2}}\right\|_{2}^{2}\right]+C.$$
(14)

with some constant C independent of θ .

For the right hand side of equation 14,

$$\mathbb{E}_{\mathbf{x}} \mathbb{E}_{\tilde{\mathbf{x}}|\mathbf{x}} \left[\left\| \{ \mathbf{s}_{2}(\tilde{\mathbf{x}}; \boldsymbol{\theta}) - \mathbf{s}_{2}(\tilde{\mathbf{x}}) \} \odot \sqrt{\mathbf{w}_{2}} \right\|_{2}^{2} \right] \\
= \sum_{i=1}^{d} \sum_{j=1}^{d} w_{2,ij} \cdot \left\{ \mathbb{E}_{\tilde{x},x} \left[s_{2,ij}^{2}(\tilde{x}; \boldsymbol{\theta}) \right] - 2\mathbb{E}_{\tilde{x},x} \left[s_{2,ij}(\tilde{x}; \boldsymbol{\theta}) s_{2,ij}(\tilde{x}) \right] + \mathbb{E}_{\tilde{x},x} \left[s_{2,ij}^{2}(\tilde{x}) \right] \right\} \\
= \sum_{i=1}^{d} \sum_{j=1}^{d} w_{2,ij} \cdot \left\{ \mathbb{E}_{\tilde{x},x} \left[s_{2,ij}^{2}(\tilde{x}; \boldsymbol{\theta}) \right] - 2\mathbb{E}_{\tilde{x},x} \left[s_{2,ij}(\tilde{x}; \boldsymbol{\theta}) s_{2,ij}(\tilde{x}) \right] + \mathbb{E}_{\tilde{x},x} \left[s_{2,ij}^{2}(\tilde{x}) \right] \right\} + C_{1}$$

where C_1 is some constant independent of $\boldsymbol{\theta}$.

For the left hand side of equation 14, note that $\frac{\mathbf{I}-\mathbf{z}\mathbf{z}^{\top}}{\sigma^{2}} = -\left\{\nabla_{\tilde{x}}^{2}q_{\sigma}(\tilde{x}|x)\right\} - \nabla_{\tilde{x}}q_{\sigma}(\tilde{x}|x) + \left\{\nabla_{\tilde{x}}q_{\sigma}(\tilde{x}|x)\right\}^{\top}$, then we know that

$$\mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x}}\left[\left\|\left\{\mathbf{s}_{2}(\tilde{\mathbf{x}};\boldsymbol{\theta})+\mathbf{s}_{1}(\tilde{\mathbf{x}})\mathbf{s}_{1}^{\top}(\tilde{\mathbf{x}})+\frac{\mathbf{I}-\mathbf{z}\mathbf{z}^{\top}}{\sigma^{2}}\right\}\odot\sqrt{\mathbf{w}_{2}}\right\|_{2}^{2}\right]$$

$$=\sum_{i=1}^{d}\sum_{j=1}^{d}w_{2,ij}\left(\mathbb{E}_{\tilde{x},x}\left[s_{2,ij}^{2}(\tilde{x};\theta)\right]+2\mathbb{E}\left[s_{2,ij}^{2}(\tilde{x};\theta)s_{1,i}(\tilde{x})s_{1,j}(\tilde{x})\right]\right.\\\left.-2\mathbb{E}\left[s_{2,ij}(\tilde{x};\theta)\nabla_{\tilde{x}_{i}}\nabla_{\tilde{x}_{j}}\log q_{\sigma}(\tilde{x}|x)\right]\right.\\\left.-2\mathbb{E}\left[s_{2,ij}(\tilde{x};\theta)\nabla_{\tilde{x}_{i}}q_{\sigma}(\tilde{x}|x)\nabla_{\tilde{x}_{j}}q_{\sigma}(\tilde{x}|x)\right]\right)+C_{2}.$$

Comparing left and right hand side of equation 14, it is sufficient to show

$$\mathbb{E}_{\tilde{x},x}\left[s_{2,ij}(\tilde{x};\theta)s_{2,ij}(\tilde{x})\right] = \mathbb{E}\left[s_{2,ij}(\tilde{x};\theta)\nabla_{\tilde{x}_i}\nabla_{\tilde{x}_j}\log q_{\sigma}(\tilde{x}|x)\right] \\ + \mathbb{E}\left[s_2(\tilde{x};\theta)\nabla_{\tilde{x}_i}q_{\sigma}(\tilde{x}|x)\nabla_{\tilde{x}_j}q_{\sigma}(\tilde{x}|x)\right] - \mathbb{E}\left[s_2(\tilde{x};\theta)s_1^2(\tilde{x})\right].$$
(15)

We have

$$\mathbb{E}_{\tilde{x},x}\left[s_{2,ij}(\tilde{x};\theta)s_{2,ij}(\tilde{x})\right] = \int s_{2,ij}(\tilde{x};\theta)s_{2,ij}(\tilde{x})p_{\tilde{x}}(\tilde{x})d\tilde{x}$$

$$= \int s_{2,ij}(\tilde{x};\theta)\left\{\nabla_{\tilde{x}_i}\left\{\frac{\nabla_{\tilde{x}_j}p(\tilde{x})}{p(\tilde{x})}\right\}\right)p_{\tilde{x}}(\tilde{x})d\tilde{x}dx$$

$$= \int s_{2,ij}(\tilde{x};\theta)\left\{\frac{\nabla_{\tilde{x}_i}\nabla_{\tilde{x}_j}p(\tilde{x})}{p(\tilde{x})} - \frac{\nabla_{\tilde{x}_i}p(\tilde{x})\cdot\nabla_{\tilde{x}_j}p(\tilde{x})}{p^2(\tilde{x})}\right\}p_{\tilde{x}}(\tilde{x})d\tilde{x}$$

$$= \int s_2(\tilde{x};\theta)\nabla_{\tilde{x}_i}\nabla_{\tilde{x}_j}p(\tilde{x})d\tilde{x} \qquad (16)$$

$$- \int s_2(\tilde{x};\theta)\frac{\nabla_{\tilde{x}_i}p(\tilde{x})\cdot\nabla_{\tilde{x}_j}p(\tilde{x})}{p(\tilde{x})}d\tilde{x}. \qquad (17)$$

For equation 16:

$$\begin{split} \int s_2(\tilde{x};\theta) \nabla_{\tilde{x}_i} \nabla_{\tilde{x}_j} p_{\tilde{x}}(\tilde{x}) d\tilde{x} &= \int s_2(\tilde{x};\theta) \nabla_{\tilde{x}_i} \nabla_{\tilde{x}_j} \left\{ \int q_\sigma(\tilde{x}|x) p_x(x) dx \right\} d\tilde{x} \\ &= \iint s_2(\tilde{x};\theta) \left\{ \nabla_{\tilde{x}_i} \nabla_{\tilde{x}_j} \log q_\sigma(\tilde{x}|x) \right\} \cdot p_x(x) d\tilde{x} dx \\ &= \iint s_2(\tilde{x};\theta) \left\{ \nabla_{\tilde{x}_i} q_\sigma(\tilde{x}|x) \nabla_{\tilde{x}_j} q_\sigma(\tilde{x}|x) \right\} \cdot q_\sigma(\tilde{x}|x) p_x(x) d\tilde{x} dx \\ &+ \iint s_2(\tilde{x};\theta) \frac{\nabla_{\tilde{x}_i} q_\sigma(\tilde{x}|x) \nabla_{\tilde{x}_j} q_\sigma(\tilde{x}|x)}{q_\sigma(\tilde{x}|x)} p_x(x) d\tilde{x} dx \\ &= \iint s_2(\tilde{x};\theta) \left\{ \nabla_{\tilde{x}_i} \nabla_{\tilde{x}_j} \log q_\sigma(\tilde{x}|x) \right\} p_{x,\tilde{x}}(x,\tilde{x}) d\tilde{x} dx \\ &+ \iint s_2(\tilde{x};\theta) \left\{ \frac{\nabla_{\tilde{x}_i} q_\sigma(\tilde{x}|x)}{q_\sigma(\tilde{x}|x)} \right\} \left\{ \frac{\nabla_{\tilde{x}_j} q_\sigma(\tilde{x}|x)}{q_\sigma(\tilde{x}|x)} \right\} p_{x,\tilde{x}}(x,\tilde{x}) d\tilde{x} dx \\ &= \mathbb{E} \left[s_2(\tilde{x};\theta) \nabla_{\tilde{x}_i} \nabla_{\tilde{x}_j} \log q_\sigma(\tilde{x}|x) \right] \\ &+ \mathbb{E} \left[s_2(\tilde{x};\theta) \nabla_{\tilde{x}_i} q_\sigma(\tilde{x}|x) \right] . \end{split}$$

For equation 17:

$$\int s_2(\tilde{x};\theta) \frac{\nabla_{\tilde{x}_i} p(\tilde{x}) \cdot \nabla_{\tilde{x}_j} p(\tilde{x})}{p(\tilde{x})} d\tilde{x} = \int s_2(\tilde{x};\theta) \frac{\nabla_{\tilde{x}_i} p(\tilde{x})}{p(\tilde{x})} \frac{\nabla_{\tilde{x}_j} p(\tilde{x})}{p(\tilde{x})} p(\tilde{x}) d\tilde{x}$$
$$= \mathbb{E} \left[s_2(\tilde{x};\theta) s_1^2(\tilde{x}) \right].$$

Combine all the results, we get the desired result.

We first consider the MCAR case. Recall that $\mathbf{w} = \mathbf{1} - \mathbf{m}$

972 C.3 PROOF OF THEOREM 7

 $= \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x}} \left[\| \left\{ \otimes^k \mathbf{x} - f_k(\tilde{\mathbf{x}}, \mathbf{s}_1(\tilde{\mathbf{x}}), \dots, \mathbf{s}_{k-1}(\tilde{\mathbf{x}}), \mathbf{s}_k(\tilde{\mathbf{x}}; \boldsymbol{\theta})) \right\} \odot \left\{ \mathbb{E}_{\mathbf{m}} \otimes^k (\mathbf{1} - \mathbf{m}) \right\} \|^2 \right]$ Noting that $\mathbb{E}_{\mathbf{m}} \otimes^k (\mathbf{1} - \mathbf{m})$ is a constant. We can show that the solution for the weighted least

 $=\mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x},\mathbf{m}}\left[\|\left\{\otimes^{k}\mathbf{x}-f_{k}(\tilde{\mathbf{x}},\mathbf{s}_{1}(\tilde{\mathbf{x}}),\ldots,\mathbf{s}_{k-1}(\tilde{\mathbf{x}}),\mathbf{s}_{k}(\tilde{\mathbf{x}};\boldsymbol{\theta})\right)\right\}\odot\otimes^{k}(\mathbf{1}-\mathbf{m})\|^{2}\right]$

 $= \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x}} \mathbb{E}_{\mathbf{m}|\tilde{\mathbf{x}},\mathbf{x}} \left[\| \left\{ \otimes^{k} \mathbf{x} - f_{k}(\tilde{\mathbf{x}}, \mathbf{s}_{1}(\tilde{\mathbf{x}}), \dots, \mathbf{s}_{k-1}(\tilde{\mathbf{x}}), \mathbf{s}_{k}(\tilde{\mathbf{x}}; \boldsymbol{\theta})) \right\} \odot \otimes^{k} (\mathbf{1} - \mathbf{m}) \|^{2} \right]$

 $\mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x},\mathbf{m}}\left[\left\|\left\{\otimes^{k}\mathbf{x}-f_{k}(\tilde{\mathbf{x}},\mathbf{s}_{1}(\tilde{\mathbf{x}}),\ldots,\mathbf{s}_{k-1}(\tilde{\mathbf{x}}),\mathbf{s}_{k}(\tilde{\mathbf{x}};\boldsymbol{\theta})\right)\right\}\odot\otimes^{k}\mathbf{w}\right\|^{2}\right]$

square equation for any constant matrix c is

$$\underset{h}{\operatorname{arg\,min}} \mathbb{E}\left[\left\|\left\{\otimes^{k}\mathbf{x}-h(\tilde{\mathbf{x}})\right\}\odot\mathbf{c}\right\|^{2}\right]=\mathbb{E}\left\{\otimes^{k}\mathbf{x}\mid\tilde{\mathbf{x}}\right\}.$$

Such equation holds since

$$\mathbb{E}\left[\left\|\left\{\otimes^{k}\mathbf{x}-h(\tilde{\mathbf{x}})\right\}\odot\mathbf{c}\right\|^{2}\right] = \mathbb{E}\left[\left\|\left\{\otimes^{k}\mathbf{x}-\mathbb{E}\left[\otimes^{k}\mathbf{x}\mid\tilde{\mathbf{x}}\right]\right\}\odot\mathbf{c}\right\|^{2}\right]\right.\\ \left.+\mathbb{E}\left[\left\|\left\{\mathbb{E}\left[\otimes^{k}\mathbf{x}\mid\tilde{\mathbf{x}}\right]-h(\tilde{\mathbf{x}})\right\}\odot\mathbf{c}\right\|^{2}\right]\right.\\ \left.\geq\mathbb{E}\left[\left\|\left\{\otimes^{k}\mathbf{x}-\mathbb{E}\left[\otimes^{k}\mathbf{x}\mid\tilde{\mathbf{x}}\right]\right\}\odot\mathbf{c}\right\|^{2}\right]\right.\right]$$

By Theorem 2, we know that $\mathbb{E}\left\{\otimes^{k}\mathbf{x} \mid \tilde{\mathbf{x}}\right\} = f_{k}(\tilde{\mathbf{x}}, \mathbf{s}_{1}(\tilde{\mathbf{x}}), \dots, \mathbf{s}_{k-1}(\tilde{\mathbf{x}}), \mathbf{s}_{k}(\tilde{\mathbf{x}}))$. Hence, the desired result follows.

For the MAR case, recall that $\mathbf{w} = \frac{1-\mathbf{m}}{\mathbb{P}[\mathbf{m}^k = 0 | \mathbf{x} = \mathbf{x}]}$

$$\begin{split} & \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x},\mathbf{m}} \left[\left\| \left\{ \otimes^{k} \mathbf{x} - f_{k}(\tilde{\mathbf{x}},\mathbf{s}_{1}(\tilde{\mathbf{x}}),\ldots,\mathbf{s}_{k-1}(\tilde{\mathbf{x}}),\mathbf{s}_{k}(\tilde{\mathbf{x}};\boldsymbol{\theta})) \right\} \odot \otimes^{k} \mathbf{w} \right\|^{2} \right] \\ & = \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x},\mathbf{m}} \left[\left\| \left\{ \otimes^{k} \mathbf{x} - f_{k}(\tilde{\mathbf{x}},\mathbf{s}_{1}(\tilde{\mathbf{x}}),\ldots,\mathbf{s}_{k-1}(\tilde{\mathbf{x}}),\mathbf{s}_{k}(\tilde{\mathbf{x}};\boldsymbol{\theta})) \right\} \odot \otimes^{k} \frac{1-\mathbf{m}}{\mathbb{P}[\mathbf{m}^{k}=0|\mathbf{x}=\mathbf{x}]} \right\|^{2} \right] \\ & = \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x}} \mathbb{E}_{\mathbf{m}|\mathbf{x}} \left[\left\| \left\{ \otimes^{k} \mathbf{x} - f_{k}(\tilde{\mathbf{x}},\mathbf{s}_{1}(\tilde{\mathbf{x}}),\ldots,\mathbf{s}_{k-1}(\tilde{\mathbf{x}}),\mathbf{s}_{k}(\tilde{\mathbf{x}};\boldsymbol{\theta})) \right\} \odot \otimes^{k} \frac{1-\mathbf{m}}{\mathbb{P}[\mathbf{m}^{k}=0|\mathbf{x}=\mathbf{x}]} \right\|^{2} \right] \\ & = \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x}} \left[\left\| \left\{ \otimes^{k} \mathbf{x} - f_{k}(\tilde{\mathbf{x}},\mathbf{s}_{1}(\tilde{\mathbf{x}}),\ldots,\mathbf{s}_{k-1}(\tilde{\mathbf{x}}),\mathbf{s}_{k}(\tilde{\mathbf{x}};\boldsymbol{\theta})) \right\} \right\|^{2} \right]. \end{split}$$

$$\arg\min_{\theta} \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x},\mathbf{m}} \left[\| \left\{ \otimes^{k} \mathbf{x} - f_{k}(\tilde{\mathbf{x}},\mathbf{s}_{1}(\tilde{\mathbf{x}}),\ldots,\mathbf{s}_{k-1}(\tilde{\mathbf{x}}),\mathbf{s}_{k}(\tilde{\mathbf{x}};\boldsymbol{\theta})) \right\} \odot \otimes^{k} \mathbf{w} \|^{2} \right] = \theta^{*}.$$

C.4 PROOF OF THEOREM 5

By similar argument in the proof of Theorem 3,

$$\begin{split} \mathbb{E}_{\mathbf{x},\mathbf{m}} \mathbb{E}_{\tilde{\mathbf{x}}|\mathbf{x},\mathbf{m}} \left[\left\| \left\{ \mathbf{s}(\tilde{\mathbf{x}};\boldsymbol{\theta}) + \frac{1}{\sigma^2} (\tilde{\mathbf{x}} - \mathbf{x}) \right\} \odot \frac{\mathbf{1} - \mathbf{m}}{\sqrt{\mathbb{P}[\mathbf{m} = 0 | \mathbf{x} = \mathbf{x}]}} \right\|_2^2 \right] \\ = \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x},\mathbf{m}} \left[\left\| \left\{ \mathbf{s}(\tilde{\mathbf{x}};\boldsymbol{\theta}) + \frac{1}{\sigma^2} (\tilde{\mathbf{x}} - \mathbf{x}) \right\} \odot \frac{\mathbf{1} - \mathbf{m}}{\sqrt{\mathbb{P}[\mathbf{m} = 0 | \mathbf{x} = \mathbf{x}]}} \right\|_2^2 \right] \\ = \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x}} \mathbb{E}_{m|\tilde{\mathbf{x}},\mathbf{x}} \left[\left\| \left\{ \mathbf{s}(\tilde{\mathbf{x}};\boldsymbol{\theta}) + \frac{1}{\sigma^2} (\tilde{\mathbf{x}} - \mathbf{x}) \right\} \odot \frac{\mathbf{1} - \mathbf{m}}{\sqrt{\mathbb{P}[\mathbf{m} = 0 | \mathbf{x} = \mathbf{x}]}} \right\|_2^2 \right] \\ = \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x}} \left[\left\| \left\{ \mathbf{s}(\tilde{\mathbf{x}};\boldsymbol{\theta}) + \frac{1}{\sigma^2} (\tilde{\mathbf{x}} - \mathbf{x}) \right\} \odot \frac{\sqrt{\mathbb{E}_{\mathbf{m}|\tilde{\mathbf{x}},\mathbf{x}}(\mathbf{1} - \mathbf{m})}}{\sqrt{\mathbb{P}[\mathbf{m} = 0 | \mathbf{x} = \mathbf{x}]}} \right\|_2^2 \right] \\ = \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x}} \left[\left\| \left\{ \mathbf{s}(\tilde{\mathbf{x}};\boldsymbol{\theta}) + \frac{1}{\sigma^2} (\tilde{\mathbf{x}} - \mathbf{x}) \right\} \right\|_2^2 \right]. \end{split}$$

where the last equality comes from $\mathbb{E}_{\mathbf{m}|\tilde{\mathbf{x}},\mathbf{x}}(1-\mathbf{m}) = \mathbb{E}_{\mathbf{m}|\mathbf{x}}(1-\mathbf{m}) = \mathbb{P}[\mathbf{m} = 0|\mathbf{x} = \mathbf{x}]$. Then following the steps in section C.1, taking w_1 in C.1 as a vector of 1, we get our desired result.

C.5 PROOF OF THEOREM 6

By similar argument in the proof of Theorem 4,

$$\begin{split} & \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x},\mathbf{m}}\left[\left\|\left\{\mathbf{s}_{2}(\tilde{\mathbf{x}};\boldsymbol{\theta})+\mathbf{s}_{1}(\tilde{\mathbf{x}})\mathbf{s}_{1}^{\top}(\tilde{\mathbf{x}})+\frac{\mathbf{I}-\mathbf{z}\mathbf{z}^{\top}}{\sigma^{2}}\right\}\odot\frac{\left\{(\mathbf{I}-\mathbf{m})(\mathbf{I}-\mathbf{m})^{\top}\right\}}{\sqrt{\mathbb{P}[\mathbf{m}\mathbf{m}^{\top}=0|\mathbf{x}=\mathbf{x}]}}\right\|_{2}^{2}\right]\\ &=&\mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x}}\left[\left\|\left\{\mathbf{s}_{2}(\tilde{\mathbf{x}};\boldsymbol{\theta})+\mathbf{s}_{1}(\tilde{\mathbf{x}})\mathbf{s}_{1}^{\top}(\tilde{\mathbf{x}})+\frac{\mathbf{I}-\mathbf{z}\mathbf{z}^{\top}}{\sigma^{2}}\right\}\odot\frac{\sqrt{\mathbb{E}_{\mathbf{m}|\tilde{\mathbf{x}},\mathbf{x}}(\mathbf{I}-\mathbf{m})(\mathbf{I}-\mathbf{m})^{\top}}}{\sqrt{\mathbb{P}[\mathbf{m}\mathbf{m}^{\top}=0|\mathbf{x}=\mathbf{x}]}}\right\|_{2}^{2}\right]\\ &=&\mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x}}\left[\left\|\left\{\mathbf{s}_{2}(\tilde{\mathbf{x}};\boldsymbol{\theta})+\mathbf{s}_{1}(\tilde{\mathbf{x}})\mathbf{s}_{1}^{\top}(\tilde{\mathbf{x}})+\frac{\mathbf{I}-\mathbf{z}\mathbf{z}^{\top}}{\sigma^{2}}\right\}\right\|_{2}^{2}\right]. \end{split}$$

LII ($J_{||_{2}|}$

where the last equality comes from $\mathbb{E}_{\mathbf{m}|\tilde{\mathbf{x}},\mathbf{x}}(1-\mathbf{m})(1-\mathbf{m})^\top = \mathbb{E}_{\mathbf{m}|\mathbf{x}}(1-\mathbf{m})(1-\mathbf{m})^\top =$ $\mathbb{P}[\mathbf{mm}^{\top} = 0 | \mathbf{x} = \mathbf{x}]$. Then following the steps in section C.2, taking \mathbf{w}_2 in C.2 as 1, we get our desired result.

C.6 PROOF OF EQUATION 8

Denote the oracle criterion function as

$$\widetilde{\mathcal{L}}_{\text{DSM}}(oldsymbol{ heta}) := \mathbb{E}_{\mathbf{x},\mathbf{m}} \mathbb{E}_{ ilde{\mathbf{x}} \mid \mathbf{x},\mathbf{m}} \left[\left\| \left\{ \mathbf{s}_1(ilde{\mathbf{x}};oldsymbol{ heta}) + rac{1}{\sigma^2}(ilde{\mathbf{x}} - \mathbf{x})
ight\} \odot \mathbf{g}_1(\mathbf{x},\mathbf{m})
ight\|_2^2
ight].$$

where $\mathbf{g}_1(\mathbf{x}, \mathbf{m}) = 1 - \mathbf{m}$ under MAR and $\mathbf{g}_1(\mathbf{x}, \mathbf{m}) = \frac{1 - \mathbf{m}}{\sqrt{\mathbb{P}[\mathbf{m} = 0 | \mathbf{x} = \mathbf{x}]}}$.

If we want to match the score of true data distribution $p(\mathbf{x})$, σ should be approximately zero for both DSM and D₂SM so that $q_{\sigma}(\tilde{\mathbf{x}})$ is close to $p(\mathbf{x})$. According to Taylor expansion we have,

$$\begin{split} & \widetilde{\mathcal{L}}_{\text{DSM}}(\boldsymbol{\theta}) = \mathbb{E}_{\tilde{\mathbf{x}},\mathbf{x},\mathbf{m}} \left[\left\| \left\{ \mathbf{s}_{1}(\tilde{\mathbf{x}};\boldsymbol{\theta}) + \frac{1}{\sigma^{2}}(\tilde{\mathbf{x}}-\mathbf{x}) \right\} \odot \mathbf{g}_{1}(\mathbf{x},\mathbf{m}) \right\|_{2}^{2} \right] \\ & 1060 \\ & = \mathbb{E}_{\mathbf{x},\mathbf{m}} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0},\mathbf{I})} \left[\left\| \left\{ \mathbf{s}_{1}(\mathbf{x}+\sigma\mathbf{z};\boldsymbol{\theta}) + \frac{\mathbf{z}}{\sigma} \right\} \odot \mathbf{g}_{1}(\mathbf{x},\mathbf{m}) \right\|_{2}^{2} \right] \\ & 1061 \\ & = \mathbb{E}_{\mathbf{x},\mathbf{m}} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0},\mathbf{I})} \left[\left\| \left\{ \mathbf{s}_{1}(\mathbf{x};\boldsymbol{\theta}) + \sigma \nabla_{\mathbf{x}} \mathbf{s}_{1}(\mathbf{x};\boldsymbol{\theta}) \mathbf{z} + \frac{\mathbf{z}}{\sigma} \right\} \odot \mathbf{g}_{1}(\mathbf{x},\mathbf{m}) \right\|_{2}^{2} \right] + \mathcal{O}(1) \\ & 1065 \\ & = \mathbb{E}_{\mathbf{x},\mathbf{m}} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0},\mathbf{I})} \left[\left\| \left\{ \mathbf{s}_{1}(\mathbf{x};\boldsymbol{\theta}) + \frac{\mathbf{z}}{\sigma} \right\} \odot \mathbf{g}_{1}(\mathbf{x},\mathbf{m}) \right\|_{2}^{2} \right] + \mathcal{O}(1) \\ & 1066 \\ & = \mathbb{E}_{\mathbf{x},\mathbf{m}} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0},\mathbf{I})} \left[\left\| \left\{ \mathbf{s}_{1}(\mathbf{x};\boldsymbol{\theta}) + \frac{\mathbf{z}}{\sigma} \right\} \odot \mathbf{g}_{1}(\mathbf{x},\mathbf{m}) \right\|_{2}^{2} \right] + \mathcal{O}(1) \\ & 1067 \\ & 1068 \\ & = \mathbb{E}_{\mathbf{x},\mathbf{m}} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0},\mathbf{I})} \left[\left\{ \left\| \mathbf{s}_{1}(\mathbf{x};\boldsymbol{\theta}) \odot \mathbf{g}_{1}(\mathbf{x}) \right\|_{2}^{2} + \frac{\left\| \mathbf{z} \odot \mathbf{g}_{1}(\mathbf{x},\mathbf{m}) \right\|_{2}^{2}}{\sigma^{2}} + \left(\frac{2}{\sigma} \mathbf{s}_{1}(\mathbf{x};\boldsymbol{\theta})^{\top} \mathbf{z} \right) \odot \mathbf{g}_{1}(\mathbf{x},\mathbf{m}) \right\} \right] + \mathcal{O}(1) \\ & 1069 \\ \end{split}$$

where $\mathbf{z} = \frac{\mathbf{x} - \mathbf{x}}{\sigma}$, where $\mathcal{O}(1)$ is bounded as σ approaches zero. However, when evaluating the expectation above from samples, the variances of $\frac{\|\mathbf{z} \odot \mathbf{g}(\mathbf{x}, \mathbf{m})\|^2}{\sigma^2}$ and $\frac{\mathbf{s}(\mathbf{x}; \boldsymbol{\theta})^\top \mathbf{z}}{\sigma}$ both increase without bound as σ nears zero, due to the terms involving σ and σ^2 in the denominator. This leads to a significant increase in the variance of the DSM loss, complicating the optimization process. As a consequence, DSM may become unstable and fail to converge when σ is small, highlighting the need for methods to reduce variance.

We have,

1078
1079
$$\mathbb{E}_{\mathbf{z}\sim\mathcal{N}(\mathbf{0},\mathbf{I})}\left[\frac{\|\mathbf{z}\odot\mathbf{g}_{1}(\mathbf{x},\mathbf{m})\|^{2}}{\sigma^{2}} + \frac{2}{\sigma}\mathbf{s}_{1}(\mathbf{x};\boldsymbol{\theta})^{\top}\mathbf{z}\right] = \frac{\|\mathbf{g}_{1}(\mathbf{x},\mathbf{m})\|^{2}}{\sigma^{2}},$$

where d is the dimension of the data distribution $p(\mathbf{x})$. Therefore, we can construct a variable that is, for sufficiently small σ , positively correlated with \mathcal{L}_{DSM} while having an expected value of zero:

$$c_{\boldsymbol{\theta}}(\mathbf{x}; \mathbf{z}) = \left(\frac{2}{\sigma} \mathbf{s}_1(\mathbf{x}; \boldsymbol{\theta})^\top \mathbf{z}\right) \odot \mathbf{g}_1(\mathbf{x}, \mathbf{m}) + \frac{\|\mathbf{z} \odot \mathbf{g}_1(\mathbf{x}, \mathbf{m})\|^2}{\sigma^2} - \frac{\|\mathbf{g}_1(\mathbf{x}, \mathbf{m})\|^2}{\sigma^2}$$

Subtracting it from \mathcal{L}_{DSM} will yield an estimator with reduced variance for DSM training with missing data:

$$\mathcal{L}_{\text{DSM-VR}}(\boldsymbol{\theta}) = \mathcal{L}_{\text{DSM}}(\boldsymbol{\theta}) - \mathbb{E}_{\mathbf{x},\mathbf{m}} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0},\mathbf{I})} \left[\left(\frac{2}{\sigma} \mathbf{s}(\mathbf{x};\boldsymbol{\theta})^{\top} \mathbf{z} \right) \odot \mathbf{g}_{1}(\mathbf{x},\mathbf{m}) + \frac{\|\mathbf{z} \odot \mathbf{g}_{1}(\mathbf{x},\mathbf{m})\|^{2}}{\sigma^{2}} \right].$$

Here we omit the part $\frac{\|\mathbf{g}_1(\mathbf{x},\mathbf{m})\|^2}{\sigma^2}$ since it is independent of $\boldsymbol{\theta}$.

C.7 PROOF OF EQUATION 9

Similar to proof C.6, consider the oracle criterion function

$$\widetilde{L}_{D_2SM}(\boldsymbol{\theta}) = \mathbb{E}_{\tilde{\mathbf{x}}, \mathbf{x}, \mathbf{m}} \left[\left\| \left\{ \mathbf{s}_2(\tilde{\mathbf{x}}; \boldsymbol{\theta}) + \mathbf{s}_1(\tilde{\mathbf{x}}) \mathbf{s}_1^\top(\tilde{\mathbf{x}}) + \frac{\mathbf{I} - \mathbf{z}\mathbf{z}^\top}{\sigma^2} \right\} \odot \mathbf{g}_2(\mathbf{x}, \mathbf{m}) \right\|_2^2 \right],$$

where $\mathbf{g}_2(\mathbf{x}_{obs}, \mathbf{m}) = (\mathbf{1} - \mathbf{m})(\mathbf{1} - \mathbf{m})^\top$ under MCAR, and $\mathbf{g}_2(\mathbf{x}_{obs}, \mathbf{m}) = \frac{(\mathbf{1} - \mathbf{m})(\mathbf{1} - \mathbf{m})^\top}{\sqrt{\mathbb{E}[\mathbf{m}\mathbf{m}^\top | \mathbf{x} = \mathbf{x}_{obs}]}}$ under MAR.

Denote $\psi(\tilde{\mathbf{x}}; \boldsymbol{\theta}) = \mathbf{s}_2(\tilde{\mathbf{x}}_{obs}; \boldsymbol{\theta}) + \mathbf{s}_1(\tilde{\mathbf{x}}_{obs})\mathbf{s}_1^{\top}(\tilde{\mathbf{x}}_{obs})$, we have

$$= \mathbb{E}_{\mathbf{x},\mathbf{m}} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0},\mathbf{I})} \left[\left\{ \| \boldsymbol{\psi}(\mathbf{x} + \sigma \mathbf{z}; \boldsymbol{\theta}) \|_{2}^{2} + \left\| \frac{\mathbf{I} - \mathbf{z} \mathbf{z}^{\top}}{\sigma^{2}} \right\|_{2}^{2} + 2\boldsymbol{\psi}(\mathbf{x} + \sigma \mathbf{z}; \boldsymbol{\theta}) \frac{\mathbf{I} - \mathbf{z} \mathbf{z}^{\top}}{\sigma^{2}} \right\} \odot \mathbf{g}_{2}(\mathbf{x}, \mathbf{m}) \right]$$

$$= \mathbb{E}_{\mathbf{x},\mathbf{m}} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0},\mathbf{I})} \left[\left\{ \| \boldsymbol{\psi}(\mathbf{x} + \sigma \mathbf{z}; \boldsymbol{\theta}) \|_{2}^{2} + \left\| \frac{\mathbf{I} - \mathbf{z} \mathbf{z}^{\top}}{\sigma^{2}} \right\|_{2}^{2} + 2\boldsymbol{\psi}(\mathbf{x} + \sigma \mathbf{z}; \boldsymbol{\theta}) \frac{\mathbf{I} - \mathbf{z} \mathbf{z}^{\top}}{\sigma^{2}} \right\} \odot \mathbf{g}_{2}(\mathbf{x}, \mathbf{m}) \right]$$

Denote $\psi_{ij}(\tilde{x};\theta)$ as the *ij*th term of $\psi(\tilde{\mathbf{x}};\theta)$, $\phi_{ij} = \mathbf{I}_{ij} - \mathbf{z}_i \mathbf{z}_j$ and g_{ij} as the *ij*th term of $\mathbf{g}_2(\mathbf{x},\mathbf{m})$ and according to Taylor expansion, we have,

$$\mathbb{E}_{x}\mathbb{E}_{z\sim\mathcal{N}(0,I)}\left[\left\{\psi_{ij}(x+\sigma z;\theta)^{2}+\frac{\phi_{ij}}{\sigma^{2}}+2\psi_{ij}(x+\sigma z;\theta)\frac{\phi_{ij}}{\sigma^{2}}\right\}\odot g_{ij}(x,m)\right]$$
$$=\mathbb{E}_{x}\mathbb{E}_{z\sim\mathcal{N}(0,I)}\left[\left\{\psi_{ij}(x;\theta)^{2}+2\psi_{ij}(x;\theta)\frac{\phi_{ij}}{\sigma^{2}}+2\nabla\psi_{ij}(x;\theta)\frac{\phi_{ij}}{\sigma}+C\right\}\odot g_{ij}(x,m)\right]+\mathcal{O}(1)$$

where $z = \frac{\tilde{x}-x}{\sigma}$, with $C = \left(\frac{\phi_{ij}}{\sigma^2}\right)^2$ and $\nabla \psi_{ij}(x+\sigma z; \theta)$ representing the derivative of $\psi_{ij}(x+\sigma z; \theta)$ with respect to x. $\left(\frac{\phi_{ij}}{\sigma^2}\right)^2$ can be treated as a constant that does not depend on θ , and $\mathcal{O}(1)$ remains bounded when $\sigma \rightarrow 0$. However, when calculating the expectation from samples, the variances of $\frac{\phi_{ij}}{\sigma^2}$ and $\nabla \psi_{ij}(x;\theta) \frac{\phi_{ij}}{\sigma^2}$ increase without bound as $\sigma \to 0$, due to the presence of σ and σ^2 in the denominator. This causes a significant rise in variance for D_2SM , making the optimization process more difficult. Consequently, D_2SM can become unstable and fail to converge as σ approaches zero, necessitating the use of variance reduction techniques.

In this case, we can employ the same variance reduction method outlined in the proof of C.6. How-ever, to bypass the need for estimating $\nabla \psi_{ij}(x;\theta) \frac{\phi_{ij}}{\sigma^2}$, we utilize the antithetic sampling technique same as Meng et al. (2021) to reduce variance.

1134 Denote $\tilde{x}_{+} = x + \sigma z$ and $\tilde{x}_{-} = x + \sigma z$ as the antithetic samples, according to Taylor expansion, the *ij*th term of the D_2 SM(θ) then becomes,

1150 where $C = 2\left(\frac{\phi_{ij}}{\sigma^2}\right)^2$, a constant with respect to the optimization.

1152 Therefore, we have the variance reduction for $\mathcal{L}_{D_2SM-VR}(\theta)$ which is equivalent to optimizing Eq. (4) up to a control variate. Moreover, when σ approaches zero, optimizing Eq. (9) is more stable.

$$\mathcal{L}_{D_2SM-VR}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x},\mathbf{m}} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0},\mathbf{I})} \left[\left\{ \boldsymbol{\psi}(\tilde{\mathbf{x}}^+)^2 + \boldsymbol{\psi}(\tilde{\mathbf{x}}^-)^2 + 2\frac{\mathbf{I} - \mathbf{z}\mathbf{z}^\top}{\sigma} \odot \boldsymbol{\Psi}(\cdot) \right\} \odot \mathbf{g}_2(\mathbf{x},\mathbf{m}) \right],$$
(18)

where the antithetic samples are defined as $\mathbf{x}^+ = \mathbf{x} + \sigma \mathbf{z}$ and $\mathbf{x}^- = \mathbf{x} - \sigma \mathbf{z}$. Here, $\boldsymbol{\psi} = \mathbf{s}_2 + \mathbf{s}_1 \mathbf{s}_1^\top$, and $\boldsymbol{\Psi} = (\boldsymbol{\psi}(\tilde{\mathbf{x}}^+) + \boldsymbol{\psi}(\tilde{\mathbf{x}}^-) - 2\boldsymbol{\psi}(\mathbf{x}))$.

1188 1189	D DATA GENERATION UNDER DIFFERENT MISSING MECHANISMS
1190	Missing data machanisms can vary significantly, but they are typically astegorized into three main
1191	types as defined by (Rubin 1976): missing completely at random (MCAR) missing at random
1192	(MAR), and missing not at random (MNAR). In our experiments, we simulate missing data based
1193	on these mechanisms as follows:
1194	MCAR (Missing Completely at Random): Missing values are generated uniformly with each
1195	data point having an equal probability of being missing determined by a predefined missing rate α
1196	Specifically, missing values are generated using a Bernoulli distribution. Ber(α), where each entry
1197	is missing independently with probability α .
1198	MAD (Missing at Dandom): In this scenario, missing values are generated using a logistic model
1199	A random subset of the variables is selected to remain fully observed, while the remaining variables
1200	have missing values depending on the fully observed ones. The missingness is determined by a
1201	logistic model with random coefficients, scaled to achieve the target proportion of missing data for
1202	the variables influenced by the fully observed subset.
1203	MNAD (Missing Not at Dandam): The MNAD mechanism is modeled using a logistic marking
1204	model. It implements two mechanisms and in either case, weights are random and the intercept is
1205	selected to attain the desired proportion of missing values.
1206	
1207	• Missing probabilities for each variable are determined by a logistic model that takes all the
1208	variables (including those with missing data) as inputs;
1209	• Variables are split into two sets: a set of input variables for the logistic model and a set of
1210	variables whose missingness is determined by the logistic model. The input variables are
1211	masked using an MCAR process, meaning the missingness in the second set depends on
1212	the missingness in the input set.
1213	In all annualization of a MAD mission markenism, and he is the manualization of a structure the libra
1214	lihood of each data point being observed. For MNAP, we utilize the same training objective as
1215	MCAR, while recognizing that this method may introduce some bias.
1217	The algorithm for training the first- and second-order models under different missing mechanisms
1218	is outlined in Algorithm 1:
1219	
1220	Algorithm 1 Mi as Capro
1221	
1222	1: Input: Observed data \mathbf{x}_{obs} , score models $\mathbf{s}_1(\cdot; \boldsymbol{\theta})$, $\mathbf{s}_2(\cdot; \boldsymbol{\theta})$, noise level σ , coefficient ω
1223	2: Infer the missingness mask $\mathbf{m} = \mathbf{I}_{[\mathbf{x}_{obs}=na]}$
1224	5: repeat 4: Sample poise $\mathbf{z} \sim \mathcal{N}(0 \mathbf{I})$
1220	5: Compute perturbed data: $\tilde{\mathbf{x}}_{abc} = \mathbf{x}_{abc} + \sigma \mathbf{z}$
1227	6: if missing mechanism is MAR then
1228	7: Estimate $\mathbb{P}[\mathbf{m} = 0 \mathbf{x} = \mathbf{x}_{obs}]$ and $\mathbb{P}[\mathbf{m}\mathbf{m}^{\top} = 0 \mathbf{x} = \mathbf{x}_{obs}]$ using fitted logistic models
1229	8: end if
1230	9: Update parameters using gradient descent on $\nabla_{\theta} (\text{Eq.}(5) + \omega \text{Eq.}(6))$
1231	10: until convergence
1232	
1233	
1234	
1235	
1236	
1237	
1238	
1239	
1240	
1241	

1242 Е ADDITIONAL INFORMATION ON SAMPLING 1243

1244

In the sampling experiments with the Swiss-Roll dataset under MCAR, we use a small perturbation 1245 $\sigma = 0.01$ and jointly optimize Eq. (8) and Eq.(9), where $s_1(\tilde{\mathbf{x}}) \approx s_1(\mathbf{x})$ and $s_2(\tilde{\mathbf{x}}) \approx s_2(\mathbf{x})$. The 1246 sample size is set to 5000. Both $s_1(\tilde{x}; \theta)$ and $s_2(\tilde{x}; \theta)$ are modeled using a 3-layer MLP with a 1247 latent size of 128 and a Softplus activation function. We use a learning rate of 0.001, a batch size 1248 of 64 and train for 100 epochs, which takes approximately 4 minutes on an Intel(R) Xeon(R) Gold 1249 6448H CPU. The experiments in Section 3 also utilize the same model configuration for training. In 1250 the Ozaki sampling experiments, we only use the diagonal of $s_2(\tilde{x}; \theta)$ to avoid the computational 1251 costs associated with the inversion, exponentiation, and decomposition of $s_2(\tilde{x}; \theta)$. The algorithm 1252 is presented in Algorithm 2.

1253 1254

1: I	nput: Score models $\mathbf{s}_1(\cdot; \boldsymbol{\theta}^*)$, $\mathbf{s}_2(\cdot; \boldsymbol{\theta}^*)$; step size ϵ ; number of iterations T
2: 1 3. f	Initialize: $\mathbf{x}_0 \sim \pi(\mathbf{x})$
3. F 4:	Sample noise $\mathbf{z}_t \sim \mathcal{N}(0, \mathbf{I})$
5:	if Ozaki sampling then
6:	Compute $\mathbf{M}_{t-1} = (e^{\epsilon \mathbf{s}_2(\tilde{\mathbf{x}}_{t-1})} - \mathbf{I}) \mathbf{s}_2(\tilde{\mathbf{x}}_{t-1})^{-1}$
7:	Compute $\Sigma_{t-1} = (e^{2\epsilon \tilde{\mathbf{x}}_{t-1}} - \mathbf{I}) \mathbf{s}_2(\tilde{\mathbf{x}}_{t-1})^{-1}$
8:	Update $ ilde{\mathbf{x}}_t = ilde{\mathbf{x}}_{t-1} + \mathbf{M}_{t-1}\mathbf{s}_1(ilde{\mathbf{x}}_{t-1}) + \Sigma_{t-1}^{1/2}\mathbf{z}_t$
9:	else
10:	Update $\tilde{\mathbf{x}}_t = \tilde{\mathbf{x}}_{t-1} + \frac{1}{2}\epsilon \mathbf{s}_1(\tilde{\mathbf{x}}_{t-1}) + \sqrt{\epsilon}\mathbf{z}_t$
11: 12: 0	ena li and for
12. C	Return: $\tilde{\mathbf{x}}_{m}$
The f Netw	following sections provide experimental details on data generation with the simulated york and real Census data.
The f Netw E.1 Baye Section	Following sections provide experimental details on data generation with the simulated ork and real Census data. DATASET DESCRIPTION AND PROCESSING sian Network Details regarding the data generated from a Bayesian Network can be on B.1 in Ouyang et al. (2023).
The f Netw E.1 Baye Section Censs based	Collowing sections provide experimental details on data generation with the simulated cork and real Census data. DATASET DESCRIPTION AND PROCESSING esian Network Details regarding the data generated from a Bayesian Network can be on B.1 in Ouyang et al. (2023). Sus Census dataset is a binary classification dataset that predict whether income exceeded on census data (Kohavi, 1996). Also known as Adult dataset.
The f Netw E.1 Baye Section Cens based The s tinuo categ	Tollowing sections provide experimental details on data generation with the simulated fork and real Census data. DATASET DESCRIPTION AND PROCESSING esian Network Details regarding the data generated from a Bayesian Network can be on B.1 in Ouyang et al. (2023). sus Census dataset is a binary classification dataset that predict whether income exceed on census data (Kohavi, 1996). Also known as Adult dataset. statistical information of datasets used in our experiments is in Table 5. #train, #te us, and #categorical mean the number of training data, testing data, continuous colusorical columns, respectively. Table 5: Synthetic and Real-World Datasets Used in Experiments.
The f Netw E.1 Baye Section Cens basec The s tinuo categ	Toolowing sections provide experimental details on data generation with the simulated fork and real Census data. DATASET DESCRIPTION AND PROCESSING esian Network Details regarding the data generated from a Bayesian Network can be on B.1 in Ouyang et al. (2023). sus Census dataset is a binary classification dataset that predict whether income exceed on census data (Kohavi, 1996). Also known as Adult dataset. statistical information of datasets used in our experiments is in Table 5. #train, #to us, and #categorical mean the number of training data, testing data, continuous colusorical columns, respectively. Table 5: Synthetic and Real-World Datasets Used in Experiments. Dataset #Train #Categorical #Continuous
The f Netw E.1 Baye Secti basec The s tinuo categ	Following sections provide experimental details on data generation with the simulated vork and real Census data. DATASET DESCRIPTION AND PROCESSING esian Network Details regarding the data generated from a Bayesian Network can be on B.1 in Ouyang et al. (2023). esus Census dataset is a binary classification dataset that predict whether income exceed on census data (Kohavi, 1996). Also known as Adult dataset. statistical information of datasets used in our experiments is in Table 5. #train, #to us, and #categorical mean the number of training data, testing data, continuous colutorical columns, respectively. Table 5: Synthetic and Real-World Datasets Used in Experiments. Dataset #Train #Test #Categorical #Continuous Bayesian Network 2000 20000 3 2

ing during generation. For discrete variables, we use one-hot encoding and apply a rounding function 1295 after the softmax function during generation.

1296 E.2 EVALUATION METHODS

We adopt the "train on synthetic, test on real (TSTR)" framework (Esteban et al., 2017), a widely
used method for assessing the quality of sampling data from generative model (Kim et al., 2022;
Ouyang et al., 2023; Li et al., 2019). The experimental results for sampling in this paper are calculated as follows:

1302 1303

1304

1305

1309

1310

1311

1312

1313

1317

1318

1319

1320

1321

1326

1328

1330

1332

1333

- 1. We first download a dataset and use its existing train-test split.
- 2. Then we generate synthetic records equal in number to the original training set using various synthetic data generation methods.
 - 3. Using the synthetic training records from Step 2, we train base classifiers to make predictions. We conduct a hyperparameter search for each classifier, considering Decision Tree, AdaBoost, Logistic Regression, MLP Classifiers, Random Forest, and XGBoost for the classification tasks. The hyperparameters and their candidate settings follow those described in Kim et al. (2022); Ouyang et al. (2023), and are summarized in Table 26 of Kim et al. (2022).
 - 4. Finally, we evaluate the classifiers using the testing dataset, applying a range of evaluation metrics for comprehensive assessment.

Steps 2 to 4 are repeated three times for each dataset, and the average scores for each method across all evaluation metrics are calculated. The detailed metrics used in our experiment include:

- 1. Accuracy: This is calculated using the accuracy_score function from the sklearn.metrics module.
- 2. Weighted-F1:

Weighted-F1 =
$$\sum_{i=0}^{N} w_i s_i$$

where N is the total number of classes. The weight for the *i*-th class, $w_i = \frac{1-p_i}{N-1}$, with p_i representing the proportion of the *i*-th class's size relative to the total dataset. Here, s_i is the F1 score for the *i*-th class, calculated using the One-vs-Rest strategy. This weighting approach is designed to prioritize the evaluation of synthesized tables by giving more importance to smaller classes, which are often prone to being overlooked by the model, thus addressing mode collapse.

- 3. AUROC: This is calculated using the roc_auc_score function from the sklearn.metrics module.
 - 4. **SDMetrics**: This metric evaluates synthetic data by comparing it against the real data, as described in (Dat, 2023).
- Among all metrics, a higher score indicates better overall quality of the synthetic data.
- 1335 1336
- E.3 MODEL ARCHITECTURE

1337 We use a perturbation of $\sigma = 0.1$ and jointly optimize Eq. (8) and Eq. (9). In the Bayesian 1338 Network experiment, we follow the same configuration as described earlier, with each run taking 1339 approximately 20 minutes. For the census dataset, we employ a simple MLP consisting of 5 Linear 1340 layers, LeakyReLU activation, Layer Normalization, and Dropout with a probability of 0.2 in the 1341 first layer. The learning rate is set to 0.001. The first two layers use a latent size of 128, while the 1342 last three layers use a latent size of 1024. We train with a batch size of 256 for 250 epochs, with each experiment taking approximately 4 hours. All experiments are performed on an Intel(R) Xeon(R) 1344 Gold 6448H CPU. For the downstream classifier, we use the same base hyperparameters as listed in 1345 Table 26 of Kim et al. (2022).

1346

- 1347 E.4 EXPERIMENTAL RESULTS
- In the following experimental results, we use a missing data ratio of $\alpha = 0.3$ and apply XGBoost for the downstream tasks, without delving into specific implementation details. Table 6 presents the

utility evaluation of MissScore using both Langevin and Ozaki samplings, compared to other baseline methods, on the Census dataset with a missing ratio of 0.3. Additionally, Table 7 summarizes
the Accuracy, AUROC, and Weighted-F1 metrics as the missing ratio varies. Figure 5 illustrates
the fidelity evaluation of MissScore, again using Langevin and Ozaki samplings alongside other
baselines, on the Bayesian dataset.

1356Table 6: Utility evaluation of MissScore using Langevin and Ozaki samplings, along with other1357baselines, on the Census dataset with missing ratio 0.3.

Criterion	Mechanism	Langevin	Ozaki	DSM-delete	DSM-mean	STaSy-mean
	MCAR	0.80	0.81	0.70	0.75	0.77
Accuracy	MAR	0.82	0.82	0.69	0.77	0.74
-	MNAR	0.81	0.80	0.59	0.80	0.75
	MCAR	0.84	0.84	0.57	0.67	0.62
AUROC	MAR	0.85	0.86	0.46	0.75	0.61
nenee	MNAR	0.86	0.86	0.52	0.76	0.63
	MCAR	0.52	0.52	0.24	0.32	0.41
Weighted-F1	MAR	0.61	0.60	0.32	0.38	0.38
Weighted 11	MNAR	0.69	0.68	0.41	0.52	0.42

Table 7: Evaluation of MissScore using Ozaki samplings on the Census dataset with varying missing ratios $\{0.1, 0.3, 0.5, 0.7, 0.9\}$.

Methods		$\alpha = 0.1$	$\alpha = 0.3$	$\alpha = 0.5$	$\alpha = 0.7$	$\alpha = 0.9$
	MCAR	0.81	0.81	0.80	0.79	0.77
Accuracy	MAR	0.71	0.82	0.82	0.82	0.79
·	MNAR	0.80	0.80	0.83	0.74	0.72
	MCAR	0.85	0.84	0.84	0.86	0.61
AUROC	MAR	0.85	0.86	0.87	0.85	0.83
	MNAR	0.85	0.87	0.86	0.85	0.80
	MCAR	0.54	0.52	0.41	0.66	0.22
Weighted-F1	MAR	0.64	0.60	0.67	0.65	0.63
-	MNAR	0.46	0.68	0.61	0.64	0.63



Figure 5: Fidelity evaluation of MissScore using Langevin and Ozaki samplings, along with other baselines, on the Census dataset varies with missing ratio $\alpha = \{0.1, 0.3, 0.5, 0.7\}$ under different missing mechanisms.

1401 E.5 ANALYSIS OF SAMPLING RESULTS USING CENSUS DATA

1403 In this section, we analyze the sampling performance using the first- and second-order models at a missing ratio of 0.3 with real Census data. Our analysis consists of two parts: numerical and

1434

1435

1404 categorical data in the census dataset. For the numerical data, which includes variables such as 1405 age, final weight (fnlwgt), years of education completed (education-num), and hours worked per 1406 week, we examine the distribution of both the original data and the sampled data shown in Figures 6 1407 to 9. This comparison is done using MissScore with first-order information via Langevin dynamics 1408 (equivalent to MissDiff) and MissScore with second-order information through Ozaki sampling. We observe that MissScore with second-order information captures finer details in the distribution. For 1409 instance, in the age range of 20-40, the fnlwgt range of 0-0.25, education-num range of 10-11, and 1410 hours per week range of 40-45, the second-order MissScore better approximates the shape of the 1411 original distribution. 1412

For the categorical data, we analyze the education variable and again find that the second-order MissScore aligns more closely with the original distribution, as shown in Figure 10. These results indicate that MissScore with second-order information captures finer-grained details in the sampled dataset, compared to MissDiff (MissScore with first-order information).



Figure 6: Analysis of Age in the Census Dataset: The first plot represents the actual data, the second plot shows the sampled data using MissScore with Langevin dynamics, and the third plot displays the sampled data using MissScore with Ozaki sampling.



Figure 7: Analysis of Final Weights in the Census Dataset: The first plot represents the actual data, the second plot shows the sampled data using MissScore with Langevin dynamics, and the third plot displays the sampled data using MissScore with Ozaki sampling.





Figure 9: Analysis of Hours-per-week in the Census Dataset: The first plot represents the actual data, the second plot shows the sampled data using MissScore with Langevin dynamics, and the third plot displays the sampled data using MissScore with Ozaki sampling.



Figure 8: Analysis of The number of years of education completed in the Census Dataset: The first plot represents the actual data, the second plot shows the sampled data using MissScore with Langevin dynamics, and the third plot displays the sampled data using MissScore with Ozaki sampling.

¹⁴⁹⁸ F Additional information on Causal Discovery

1501 F.1 RELATED WORK

1473

1474

1475

1496 1497

1500

1502 **Causal Discovery with Complete data.** Causal discovery aims to uncover the underlying causal 1503 relationships among variables of interest from purely observational data, specifically identifying a 1504 causal Directed Acyclic Graph (DAG) for a given dataset. This problem lies at the heart of causal inference, as knowledge of the causal graph enables prediction of the effects of interventions. How-1506 ever, causal discovery from observational data is inherently ill-posed, necessitating additional as-1507 sumptions, such as imposing functional assumptions on the data-generating process. We adopt the notion of structural causal model (SCM) to characterize the causal relations among variables. Each SCM $\mathcal{M} = \langle \mathcal{Z}, \mathcal{X}, \mathcal{F} \rangle$ consists of the exogenous variable set $\mathcal{Z} = \{Z_1, Z_2, \dots, Z_d\}$, the endoge-1509 nous variable set $\mathcal{X} = \{X_1, X_2, \dots, X_d\}$, and the function set $\mathcal{F} = \{f_1, f_2, \dots, f_d\}$. Here, each 1510 function f_i computes the variable X_i from its parents (or causes) X_{PA_i} and an exogenous variable 1511 Z_i , i.e., $X_i = f_i(X_{PA_i}, Z_i)$. We focus on a specific class of SCMs, called the additive noise models



Figure 10: Analysis of Education in the Census Dataset: The first plot represents the actual data, the second plot shows the sampled data using MissScore with Langevin dynamics, and the third plot displays the sampled data using MissScore with Ozaki sampling.

1531 (ANMs), given by $X_i = f_i(X_{PA_i}) + Z_i$, i = 1, 2, ..., d, where Z_i , interpreted as the additive 1532 noise variable, is assumed to be independent of variables in X_{PA_i} and mutually independent with 1533 variables in $Z \setminus Z_i$.

Rolland et al. (2022) proposed an order-based algorithm for this model, further assuming that f_i is a twice-differentiable nonlinear function and Z_i is Gaussian noise. This method enables the identification of leaf nodes based on the diagonal of the Hessian of the log-likelihood. Before proceeding to the method for identifying leaves, we first derive an analytical expression for the score following Lemma 2 in Rolland et al. (2022). The score is written as follows:

1539 1540 1541

 $\nabla_{x_j} \log p(\mathbf{x}) = \nabla_{x_j} \log \prod_{i=1}^d p(x_i \mid x_{\mathsf{PA}_i})$ $= \nabla_{x_j} \sum_{i=1}^d \log p(x_i \mid x_{\mathsf{PA}_i})$

1542 1543

 $= \frac{\partial \log p(x_j - f_j(x_{\mathsf{PA}_j}))}{\partial x_j} - \sum_{i \in \mathsf{CH}_i} \frac{\partial f_i}{\partial x_j} \frac{\partial \log p(x_i - f_i(x_{\mathsf{PA}_i}))}{\partial x}.$

1549 1550

where CH_j represents the children of the variable j. As a result, $\frac{\partial}{\partial x_j} \nabla_{x_j} \log p(\mathbf{x}) = a$, where a is a constant, Consequently, the variance of the diagonal elements of the Hessian is zero (i.e. Var_{**x**}[H_{j,j}(log p(**x**))] = 0) if and only if node j is a leaf node.

 $= \nabla_{x_j} \sum_{i=1}^d \log p(x_i - f_i) \quad \text{(where } z_i = x_i - f_i(x_{\text{PA}_i})\text{)}$

However, existing computational methods for calculating the Hessian struggle to scale efficiently as
the number of variables and samples increases, limiting the scalability of Rolland et al. (2022). To
address this, Sanchez et al. (2022) introduced a diffusion-based model that efficiently computes the
Hessian, enabling the method to scale to larger datasets, both in terms of sample size and number of
variables, while maintaining comparable performance to Rolland et al. (2022).

Causal Discovery with Incomplete Data. Several extensions of the PC algorithm have been developed to learn causal graphs from incomplete data (Tu et al., 2019; Gain & Shpitser, 2018), utilizing
only the fully observed samples while mitigating biases in conditional independence tests. Another
prominent family of methods relies on Expectation-Maximization (Dempster et al., 1977), where
missing values are iteratively inferred while simultaneously learning the causal structure. Building
on the continuous optimization techniques introduced by NOTEARS (Zheng et al., 2018), MissDAG
(Gao et al., 2022) extends this approach to continuous identifiable Additive Noise Models (ANMs),

using approximate posterior inference via Monte Carlo and rejection sampling when the exact posterior is unavailable. MissOTM (Vo et al., 2024) introduces a score-based method that leverages optimal transport to learn causal structures from incomplete data. A distinct approach is taken by VISL (Morales-Alvarez et al., 2022), which employs amortized variational inference in a Bayesian framework. Unlike MissOTM and MissDAG, VISL assumes a latent low-dimensional factor that captures the essential structure of the data based on observed variables. The latent factors are then used to reconstruct the complete data and discover the underlying causal graph.

1573 1574

1575

F.2 EVALUATION METRICS

1576 For each method, we compute the

SHD. Structural Hamming distance between the output and the true causal graph, which counts the number of missing, falsely detected, or reversed edges.

Order Divergence. Rolland et al. (2022) propose this quantity for measuring how well the topological order is estimated. For an ordering π , and a target adjacency matrix A, we define the topological order divergence $D_{\text{top}}(\pi, \mathbf{A})$ as

1585

 $D_{\text{top}}(\pi, \mathbf{A}) = \sum_{i=1}^{d} \sum_{j:\pi_i > \pi_j} \mathbf{A}_{ij}$ (19)

86

1587 F.3 MODEL ARCHITECTURE

We apply a perturbation of $\sigma = 0.1$ and jointly optimize Eq. (8) and Eq. (9). The model is a simple MLP with 5 Linear layers, LeakyReLU activation, Layer Normalization, and a Dropout rate of 0.2 in the first layer. The learning rate is set to 0.001. The first two layers have a latent size of max(128, 3 × d), while the last three use a latent size of max(1024, 5 × d). Training is conducted with a batch size of 128 for 150 epochs. The time efficiency is shown in the figure 4 with ER graph model across various dimensions. All experiments are executed on an Intel(R) Xeon(R) Gold 6448H CPU. The algorithm is presented in Algorithm 3.

1596

1596

Algorithm 3 MissScore-Causal Discovery 1598 1: Input: Observed data \mathbf{x}_{obs} ; score models $\mathbf{s}_1(\cdot; \boldsymbol{\theta}), \mathbf{s}_2(\cdot; \boldsymbol{\theta})$ 2: Initialize: $\pi = []; \text{ nodes} = \{1, ..., d\}$ 3: $n, d \leftarrow$ shape of \mathbf{x}_{obs} 4: **for** k = 1 to d **do** 5: Jointly train the score models $s_1(\theta)$ and $s_2(\theta)$ using x_{obs} with Algorithm 1 Generate n samples $\tilde{\mathbf{x}}_{new}$ using Algorithm 2 with bootstrapping 6: 1604 Estimate the second-order score $s_2(\tilde{x}_{new})$ using $s_2(\theta)$ 7: $V_j = \operatorname{Var}_X[\operatorname{diag}(\mathbf{s}_2(\tilde{\mathbf{x}}_{\operatorname{new}}))]$ 8: 9: $\ell \leftarrow \arg\min_{j \in \text{nodes}} V_j$ \triangleright The leaf node 1607 10: $\pi \leftarrow [\ell, \pi]$ ▷ Update topological order 11: nodes \leftarrow nodes $- \{\ell\}$ \triangleright Remove node ℓ 1609 12: Remove the ℓ -th column from \mathbf{x}_{obs} 1610 13: end for 14: Obtain the final DAG using CAM pruning associated with the topological order π . 1611 1612 1613 1614

1615 F.4 EXPERIMENTAL RESULTS 1616

We begin by conducting experiments using complete data, evaluating our approach alongside various
missing mechanisms and different missing ratios. Additionally, we include order divergence in
the Table. Our findings reveal that, with complete data, MissScore performs on par with DiffAN, as illustrated in Figure 11. Notably, among all settings, MissScore achieves performance similar



Figure 11: The data is generated using an ER graph model with different dimensions $d = \{10, 20, 50\}$ and an equal number of edges. Each dataset consists of 1000 samples. Left: f_i is an MLP; Right: f_i corresponds to MIM.



Figure 12: The data is generated under MCAR with missing ratios of 0.1 and 0.3, using an ER graph model with different dimensions $d = \{10, 20, 50\}$ and an equal number of edges. Each dataset consists of 1000 samples. f_i corresponds to MLP. Left: SHD with missing ratio 0.1; Right: SHD with missing ratio 0.3.

Table 8: Order divergence with missing ratios of $\alpha = \{0.1, 0.3\}$ across different missing data mechanisms. The ER graph model is considered with varying dimensions $d = \{10, 20, 50\}$, and an equal number of edges. Each dataset consists of 1000 samples and f_i corresponds to MLP. Lower order divergence indicating better performance.

Dimensions	Methods -	MCAR		MAR		MNAR	
Dimensions		$\alpha = 0.1$	$\alpha = 0.3$	$\alpha = 0.1$	$\alpha = 0.3$	$\alpha = 0.1$	$\alpha = 0$
d=10	MissScore MissDiffAN	$\begin{array}{c} 1.62 \pm 0.99 \\ 2.00 \pm 1.18 \end{array}$	$\begin{array}{c} 1.60 \pm 0.66 \\ 2.90 \pm 1.87 \end{array}$	$\begin{array}{c} 1.33 \pm 0.47 \\ 2.75 \pm 1.30 \end{array}$	$\begin{array}{c} 1.89 \pm 0.75 \\ 3.60 \pm 1.28 \end{array}$	$\begin{array}{c} 1.33 \pm 0.47 \\ 2.60 \pm 1.20 \end{array}$	$1.44 \pm 2.30 \pm$
d=20	MissScore MissDiffAN	$\begin{array}{c} 2.10 \pm 0.94 \\ 4.30 \pm 2.00 \end{array}$	$\begin{array}{c} 2.67 \pm 2.00 \\ 4.22 \pm 2.35 \end{array}$	2.70 ± 1.85 4.33 ± 2.11	$\begin{array}{c} 1.94 \pm 0.29 \\ 3.00 \pm 1.00 \end{array}$	$\begin{array}{c} 1.70 \pm 1.78 \\ 2.00 \pm 0.89 \end{array}$	$\begin{array}{c} 0.70 \pm \\ 3.40 \pm \end{array}$
d=50	MissScore MissDiffAN	$\begin{array}{c} 3.40 \pm 2.33 \\ 4.50 \pm 2.91 \end{array}$	$\begin{array}{c} 4.00 \pm 3.10 \\ 3.60 \pm 3.32 \end{array}$	$3.10 \pm 1.70 \\ 8.33 \pm 3.65$	$\begin{array}{c} 4.60 \pm 2.91 \\ 7.40 \pm 2.84 \end{array}$	$\begin{array}{c} 2.80 \pm 0.98 \\ 3.50 \pm 1.75 \end{array}$	$4.10 \pm 3.20 \pm$



Figure 13: The data is generated under MCAR with missing ratios of 0.1 and 0.3, using an ER graph model with different dimensions $d = \{10, 20, 50\}$ and an equal number of edges. Each dataset consists of 1000 samples. f_i corresponds to MIM. Left: SHD with missing ratio 0.1; Right: SHD with missing ratio 0.3.



Figure 14: The data is generated under MAR with missing ratios of 0.1 and 0.3, using an ER graph model with different dimensions $d = \{10, 20, 50\}$ and an equal number of edges. Each dataset consists of 1000 samples. f_i corresponds to MLP. Left: SHD with missing ratio 0.1; Right: SHD with missing ratio 0.3.



Figure 15: The data is generated under MAR with missing ratios of 0.1 and 0.3, using an ER graph model with different dimensions $d = \{10, 20, 50\}$ and an equal number of edges. Each dataset consists of 1000 samples. f_i corresponds to MIM. Left: SHD with missing ratio 0.1; Right: SHD with missing ratio 0.3.



Figure 16: The data is generated under MNAR with missing ratios of 0.1 and 0.3, using an ER graph model with different dimensions $d = \{10, 20, 50\}$ and an equal number of edges. Each dataset consists of 1000 samples. f_i corresponds to MLP. Left: SHD with missing ratio 0.1; Right: SHD with missing ratio 0.3.



Figure 17: The data is generated under MNAR with missing ratios of 0.1 and 0.3, using an ER graph model with different dimensions $d = \{10, 20, 50\}$ and an equal number of edges. Each dataset consists of 1000 samples. f_i corresponds to MIM. Left: SHD with missing ratio 0.1; Right: SHD with missing ratio 0.3.

Table 9: Order divergence with missing ratios of $\alpha = \{0.1, 0.3\}$ across different missing data mechanisms. The ER graph model is considered with varying dimensions $d = \{10, 20, 50\}$, and an equal number of edges. Each dataset consists of 1000 samples and f_i corresponds to MIM. Lower order divergence indicating better performance.

Dimensions	Methods	MCAR		MAR		MNAR	
		$\alpha = 0.1$	$\alpha = 0.3$	$\alpha = 0.1$	$\alpha = 0.3$	$\alpha = 0.1$	$\alpha = 0$
d=10	MissScore MissDiffAN	$\begin{array}{c} 1.60 \pm 1.20 \\ 2.22 \pm 2.10 \end{array}$	$\begin{array}{c} 1.20 \pm 0.60 \\ 1.50 \pm 0.67 \end{array}$	$\begin{array}{c} 1.35 \pm 0.47 \\ 2.50 \pm 1.36 \end{array}$	$\begin{array}{c} 1.88 \pm 0.93 \\ 3.70 \pm 1.10 \end{array}$	$\begin{array}{c} 1.62 \pm 0.99 \\ 2.40 \pm 2.20 \end{array}$	$1.38 \pm 2.20 \pm$
d=20	MissScore MissDiffAN	$\begin{array}{c} 1.70 \pm 0.90 \\ 2.70 \pm 1.49 \end{array}$	$\begin{array}{c} 2.20 \pm 1.33 \\ 3.8 \pm 1.47 \end{array}$	$\begin{array}{c} 1.82 \pm 1.29 \\ 3.70 \pm 1.35 \end{array}$	$\begin{array}{c} 1.90 \pm 1.30 \\ 3.10 \pm 2.21 \end{array}$	$\begin{array}{c} 1.56 \pm 0.68 \\ 2.44 \pm 1.07 \end{array}$	$0.56 \pm 2.70 \pm$
d=50	MissScore MissDiffAN	$\begin{array}{c} 4.90 \pm 2.12 \\ 5.56 \pm 2.45 \end{array}$	$\begin{array}{c} 4.00 \pm 1.41 \\ 4.40 \pm 1.56 \end{array}$	$4.56 \pm 2.27 \\ 7.90 \pm 2.39$	$\begin{array}{c} 4.60 \pm 2.24 \\ 7.50 \pm 2.69 \end{array}$	$\begin{array}{c} 4.20 \pm 1.66 \\ 5.00 \pm 1.61 \end{array}$	$4.40 \pm 6.20 \pm$