


Unsupervised Imputation of Non-Ignorably Missing Data Using Importance-Weighted Autoencoders

David K. Lim, Naim U. Rashid, Junier B. Oliva & Joseph G. Ibrahim

To cite this article: David K. Lim, Naim U. Rashid, Junier B. Oliva & Joseph G. Ibrahim (15 Jul 2024): Unsupervised Imputation of Non-Ignorably Missing Data Using Importance-Weighted Autoencoders, Statistics in Biopharmaceutical Research, DOI: [10.1080/19466315.2024.2368787](https://doi.org/10.1080/19466315.2024.2368787)



To link to this article: <https://doi.org/10.1080/19466315.2024.2368787>

 View supplementary material 

 Published online: 15 Jul 2024.

 Submit your article to this journal 

 Article views: 17

 View related articles 

 View Crossmark data 



Unsupervised Imputation of Non-Ignorably Missing Data Using Importance-Weighted Autoencoders

David K. Lim^a, Naim U. Rashid^a, Junier B. Oliva^b, and Joseph G. Ibrahim^a

^aDepartment of Biostatistics, University of North Carolina at Chapel Hill, Chapel Hill, NC; ^bDepartment of Computer Science, University of North Carolina at Chapel Hill, Chapel Hill, NC

ABSTRACT

Deep Learning (DL) methods have dramatically increased in popularity in recent years. While its initial success was demonstrated in the classification and manipulation of image data, there has been significant growth in the application of DL methods to problems in the biomedical sciences. However, the greater prevalence and complexity of missing data in biomedical datasets present significant challenges for DL methods. Here, we provide a formal treatment of missing data in the context of Variational Autoencoders (VAEs), a popular unsupervised DL architecture commonly used for dimension reduction, imputation, and learning latent representations of complex data. We propose a new VAE architecture, NIMIWAE, that is one of the first to flexibly account for both ignorable and non-ignorable patterns of missingness in input features at training time. Following training, samples can be drawn from the approximate posterior distribution of the missing data can be used for multiple imputation, facilitating downstream analyses on high dimensional incomplete datasets. We demonstrate through statistical simulation that our method outperforms existing approaches for unsupervised learning tasks and imputation accuracy. We conclude with a case study of an EHR dataset pertaining to 12,000 ICU patients containing a large number of diagnostic measurements and clinical outcomes, where many features are only partially observed.

ARTICLE HISTORY

Received October 2023
Accepted June 2024

KEYWORDS

Deep learning; Electronic health records; Missing data; MNAR; Variational autoencoder

1. Introduction

Deep Learning (DL) methods have dramatically increased in popularity in recent years. While its initial success was demonstrated in the classification and manipulation of image data, there has been significant growth in the application of DL methods to an array of problems in the biomedical sciences (Shickel et al. 2018; Lopez et al. 2018). The presence of complex interactions among high-dimensional features in modern biomedical data has motivated the use of Variational Autoencoders (VAEs), a DL architecture commonly used for unsupervised learning tasks such as dimension reduction, representational learning, and generation of synthetic data mimicking real input data, which may be unavailable due to patient confidentiality (Shickel et al. 2018). In prior evaluations, VAEs have shown tremendous performance in data generation and representation learning (Kingma and Welling 2019).

Some of the useful characteristics of the VAE can be derived from its underlying architecture. This structure is similar to that of the Autoencoder (AE), which consists of two major parts. The first part is the encoder neural network, which in the AE deterministically transforms the input data into a lower-dimensional set of factors that is assumed to capture the salient qualities of the data. The second is the decoder, which is a separate neural network that attempts to reconstruct the original data from the lower dimensional space itself (Tschannen, Bachem, and Lucic 2018). The quality of the dimension reduction and

reconstruction is measured by the reconstruction error, which is measured by the difference in the input data and its reconstructed version. This lower dimensional space may be used to define structure in the underlying set of input data (Wang, Yao, and Zhao 2016), similar to PCA. However, the encoder and decoder are often represented by feed forward neural networks, allowing for complex nonlinear interactions between features to be captured in both (LeCun, Bengio, and Hinton 2015). The VAE additionally imposes a probabilistic assumption on the input data and the lower-dimensional (or “latent”) space (Doersch 2016), where explicit encoder and decoder distributions are defined, and neural networks are used to output the parameters of each distribution (Kingma and Welling 2019). The learned decoder allows the VAE to also generate synthetic data that has high fidelity to the original training data, given the learned latent space.

However, accounting for the prevalence and complexity of missing data in modern biomedical datasets presents a significant challenge for training VAEs and other DL architectures. Existing DL methods that claim to handle missing data either require pre-training on studies with fully observed data, use heuristics to replace missing observations, or cannot handle more complicated forms of missingness, such as Missing at Random (MAR), or Missing Not a Random (MNAR) (Li, Akbar, and Oliva 2020; Strauss and Oliva 2021). Intuitively, when observations with missing values differ systematically from those

without missingness, results from standard approaches that do not properly account for this fact may no longer be accurate. For example, the common presence of missing data across patient records in electronic health record data has been shown to present a significant barrier to the generalizability and applicability of deep learning methods (Wells et al. 2013). EHR data also typically contain many features pertaining to a large number of patients (Ross, Wei, and Ohno-Machado 2014), and such features may show complex interactions with each other as well as with clinical outcomes of interest (Vinzamuri and Reddy 2013), compounding the impact of the missingness of these features. For these reasons, the ability to flexibly account for different patterns of missing data in modern biomedical datasets, especially in the more difficult non-ignorable missingness setting, would have great utility.

In this article, we introduce a novel deep learning architecture called Non-Ignorably Missing Importance-weighted Autoencoder (NIMIWAE) that treats missing observations as latent variables in the VAE framework using Importance-Weighted Autoencoders (IWAEs). Here we model the probability of missingness using a feed forward neural network, facilitating more flexible handling of MNAR patterns of missingness across a large set of available features. Through simulations, we show that proper modeling of the missingness mechanism increases the accuracy of missing data imputation, as well as downstream coefficient estimation using these imputed datasets. Lastly, we show that in the Physionet 2012 Challenge EHR dataset, accounting for MNAR missingness results in significant differences in downstream fitted models to predict in-hospital mortality. More generally, NIMIWAE is a flexible framework within the VAE family to handle multiple patterns of missingness commonly found in complex biomedical datasets.

2. Methods

2.1. Variational Autoencoder

Let \mathbf{X} be an $n \times p$ data matrix, where \mathbf{x}_i denotes the observation vector pertaining to the i th observation, $i = 1, \dots, n$, and x_{ij} denotes the value of the j th feature in this vector, $j = 1, \dots, p$. In a VAE, we assume $\mathbf{x}_1, \dots, \mathbf{x}_n$ are iid samples from a multivariate p.d.f or “generative model” $p_\psi(\mathbf{X}|\mathbf{Z})$. Here, \mathbf{Z} is an $n \times d$ matrix, such that $\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_n\}$ and \mathbf{z}_i is a latent vector of length d pertaining to the i th sample (Kingma and Welling 2019). Typically it is assumed that $d \leq p$, such that \mathbf{Z} constitutes a lower-dimensional representation of the original data \mathbf{X} . The parameters ψ and conditional distribution $p_\psi(\mathbf{X}|\mathbf{Z})$ indicate how the observed data \mathbf{X} may be generated from \mathbf{Z} . In this manner, a VAE aims to learn accurate representations of high-dimensional data, and may be used to generate synthetic data with similar qualities as its training data. These aspects are also aided through the use of embedded deep learning neural networks, for example within $p_\psi(\mathbf{X}|\mathbf{Z})$, which also facilitate its applicability to larger dimensions and complex datasets.

2.1.1. Objective Function

Since ψ is unknown, learning is performed by maximizing the marginal log-likelihood of \mathbf{X} with respect to ψ , where we denote this marginal log-likelihood as $\log p_\psi(\mathbf{X}) =$

$\log \int p_\psi(\mathbf{X}, \mathbf{Z}) d\mathbf{Z} = \log \int p_\psi(\mathbf{X}|\mathbf{Z}) p(\mathbf{Z}) d\mathbf{Z}$. However, due to the integral involved, this quantity is often intractable and is difficult to maximize directly. Therefore, VAEs alternatively optimize the so-called “Evidence Lower Bound” (ELBO), which has the following form (Kingma and Welling 2019):

$$\mathcal{L}^{\text{ELBO}}(\theta, \psi) = \mathbb{E}_{\mathbf{Z} \sim q_\theta(\mathbf{Z}|\mathbf{X})} \log \left[\frac{p_\psi(\mathbf{X}|\mathbf{Z}) p(\mathbf{Z})}{q_\theta(\mathbf{Z}|\mathbf{X})} \right] \quad (1)$$

$$\hat{\mathcal{L}}_K^{\text{ELBO}}(\theta, \psi) = \frac{1}{K} \sum_{k=1}^K \log \left[\frac{p_\psi(\mathbf{X}|\tilde{\mathbf{Z}}_k) p(\tilde{\mathbf{Z}}_k)}{q_\theta(\tilde{\mathbf{Z}}_k|\mathbf{X})} \right]. \quad (2)$$

Here, $\mathcal{L}^{\text{ELBO}}(\theta, \psi)$ denotes the ELBO such that $\mathcal{L}^{\text{ELBO}}(\theta, \psi) \leq \log p_\psi(\mathbf{X})$. Also let $\hat{\mathcal{L}}_K^{\text{ELBO}}(\theta, \psi)$ denote the empirical approximation to (1) computed by Monte Carlo integration, such that $\mathcal{L}^{\text{ELBO}}(\theta, \psi) \approx \hat{\mathcal{L}}_K^{\text{ELBO}}(\theta, \psi)$ and $\tilde{\mathbf{Z}}_1, \dots, \tilde{\mathbf{Z}}_K$ are K samples drawn from $q_\theta(\mathbf{Z}|\mathbf{X})$, the variational approximation of the true but intractable posterior $p_\psi(\mathbf{Z}|\mathbf{X})$, also called the “recognition model”. Furthermore, denote $f_\psi(\mathbf{Z})$ and $g_\theta(\mathbf{X})$ as the decoder and encoder feed forward neural networks of the VAE, where ψ and θ are the sets of weights and biases pertaining to each of these neural networks, respectively. Given \mathbf{Z} , $f_\psi(\mathbf{Z})$ outputs the distributional parameters pertaining to $p_\psi(\mathbf{X}|\mathbf{Z})$. Given \mathbf{X} , $g_\theta(\mathbf{X})$ outputs the distributional parameters for $q_\theta(\mathbf{Z}|\mathbf{X})$.

In variational inference, $q_\theta(\mathbf{Z}|\mathbf{X})$ is constrained to be from a class of simple distributions, or “variational family”, to obtain the best candidate from within that class to approximate $p_\psi(\mathbf{Z}|\mathbf{X})$. Variational inference is usually used in tandem with amortization of the parameters where the neural network parameters are shared across observations (Gershman and Goodman 2014), allowing for stochastic gradient descent (SGD) to be used for optimization of Eq. (2) (Kingma and Welling 2019). In practice, both $q_\theta(\mathbf{Z}|\mathbf{X})$ and $p(\mathbf{Z})$ are typically assumed to have simple forms, such as multivariate Gaussians with diagonal covariance structures, and $q_\theta(\mathbf{Z}|\mathbf{X})$ is commonly assumed to be factorizable, such that $q_\theta(\mathbf{Z}|\mathbf{X}) = \prod_{i=1}^n q_\theta(\mathbf{z}_i|\mathbf{x}_i)$ (Kingma and Welling 2019).

2.1.2. Estimation Procedure and Use Cases

Let $(\hat{\theta}^{(t)}, \hat{\psi}^{(t)})$ be the estimates of (θ, ψ) at update (or iteration) t . For $t = 0$, these values are often initialized to small values centered around 0, although other initialization schemes may be used (Saxe, McClelland, and Ganguli 2014; Murphy 2016). Each subsequent update $t \geq 1$ consists of two general steps to maximize $\mathcal{L}(\theta, \psi)$. First, K samples are drawn from $q_{\hat{\theta}^{(t)}}(\mathbf{Z}|\mathbf{X})$ to compute the quantity in 2), conditional on $(\hat{\theta}^{(t)}, \hat{\psi}^{(t)})$. Then, the so-called “reparameterization trick” is used to facilitate the calculation of gradients of this approximation to obtain $(\hat{\theta}^{(t+1)}, \hat{\psi}^{(t+1)})$ using stochastic gradient descent (Kingma and Welling 2019). This procedure may be repeated for a fixed number of iterations, or may be terminated early via pre-specified early stop criteria (Prechelt 1998). Kingma and Welling (2019) provides additional details on the maximization procedure for VAEs. The networks $f_\psi(\mathbf{Z})$ and $g_\theta(\mathbf{X})$ also allow the VAE to capture complex and nonlinear relationships between features when outputting the distributional parameters for the generative and recognition models, respectively, through the inclusion of hidden layers in each network. The number of hidden layers and nodes per layer for each network are commonly determined via hyperparameter tuning.

After model fitting, the VAE has several useful features. First, synthetic data can be generated by sampling from the learned generative model $p_{\hat{\psi}}(\mathbf{X}|\mathbf{Z})$ after drawing a sample from $q_{\hat{\theta}}(\mathbf{Z}|\mathbf{X})$ (Mattei and Frellsen 2019; Nazabal et al. 2018). Second, the posterior modes in the latent space \mathbf{Z} can be determined from $q_{\hat{\theta}}(\mathbf{Z}|\mathbf{X})$. These posterior modes may be used for purposes such as clustering or substructure discovery (Lim, Jiang, and Yi 2020). In some applications, $p_{\hat{\psi}}(\mathbf{X}|\mathbf{Z})$ may also be used to directly perform imputation (Nazabal et al. 2018), however, the statistical properties of this procedure have not been thoroughly discussed in prior work.

2.2. Importance-Weighted Autoencoder

The IWAE (Burda, Grosse, and Salakhutdinov 2015) is a generalization of the standard VAE, where the resulting IWAE bound, corresponding to the ELBO in (1), can be written as

$$\mathcal{L}_K^{IWAE}(\theta, \psi) = \mathbb{E}_{\mathbf{Z}_k \sim q_{\theta}(\mathbf{Z}|\mathbf{X})} \log \left[\frac{1}{K} \sum_{k=1}^K \frac{p_{\psi}(\mathbf{X}|\mathbf{Z}_k)p(\mathbf{Z}_k)}{q_{\theta}(\mathbf{Z}_k|\mathbf{X})} \right] \quad (3)$$

$$\hat{\mathcal{L}}_K^{IWAE}(\theta, \psi) = \log \left[\frac{1}{K} \sum_{k=1}^K \frac{p_{\psi}(\mathbf{X}|\tilde{\mathbf{Z}}_k)p(\tilde{\mathbf{Z}}_k)}{q_{\theta}(\tilde{\mathbf{Z}}_k|\mathbf{X})} \right]. \quad (4)$$

An important distinction in (3) from (1) is that a VAE assumes a single latent variable \mathbf{Z} in 1) that is sampled K times in the ELBO approximation from 2). In contrast, an IWAE assumes K iid latent variables in the expression for its lower bound, where $\mathbf{Z}_1, \dots, \mathbf{Z}_K \stackrel{\text{iid}}{\sim} q_{\theta}(\mathbf{Z}|\mathbf{X})$. The contribution of each \mathbf{Z}_k in (3) is weighted by $\frac{p(\mathbf{Z}_k)}{q_{\theta}(\mathbf{Z}_k|\mathbf{X})}$. Then, to compute its empirical approximation $\hat{\mathcal{L}}_K^{IWAE}$, typically only one sample is drawn for each \mathbf{Z}_k in (4). For $K > 1$, Burda, Grosse, and Salakhutdinov (2015) showed that $\log p(\mathbf{X}) \geq \mathcal{L}_{K+1}^{IWAE} \geq \mathcal{L}_K^{IWAE}$, such that $\mathcal{L}_K^{IWAE} \rightarrow \log p(\mathbf{X})$ as $K \rightarrow \infty$ if $p_{\psi}(\mathbf{X}, \mathbf{Z})/q_{\theta}(\mathbf{Z}|\mathbf{X})$ is bounded. Thus, the IWAE bound more closely approximates the true marginal log-likelihood when $K > 1$ (Cremer, Morris, and Duvenaud 2017), at the cost of greater computational burden. For $K = 1$, $\mathcal{L}_1^{IWAE} = \mathcal{L}^{VAE}$, the IWAE corresponds exactly to the standard VAE. In this way, the IWAE can be considered to be part of the VAE family, and we refer to methods that use either VAEs or IWAEs broadly as “VAE methods”. Intuitively, using more Monte Carlo samples, that is $K > 1$, allows the IWAE to prevent the lower bound from being greatly susceptible to occasional poor samples from the variational posterior: a limitation of the standard VAE. A visualization of the workflow for an IWAE (without missing data) can be found in Section A of the supplementary materials.

VAE methods have shown excellent performance in representation learning on many types of data. However, the presence of missingness in \mathbf{X} presents significant challenges to the above modeling procedures and the application of VAE methods in general.

2.3. Missing Data

In this section, we first introduce notation for missing data and review the different mechanisms of missingness, as described in the statistical literature. Let the data be factored such that $\mathbf{X} = \{\mathbf{X}^o, \mathbf{X}^m\}$, with \mathbf{X}^o denoting the observed values and \mathbf{X}^m

denoting the missing values. For each observation vector \mathbf{x}_i , denote \mathbf{x}_i^o and \mathbf{x}_i^m respectively to be the observed and missing features of \mathbf{x}_i . Also, let \mathbf{R} be a matrix of the same dimension as \mathbf{X} , with entries $r_{ij} = I(x_{ij} \text{ is observed})$ for the i th observation and j th feature, where $I(\cdot)$ denotes the indicator function. In this way, \mathbf{R} is the “mask” matrix pertaining to \mathbf{X} , such that $\mathbf{x}_i^o = \{x_{ij} : r_{ij} = 1\}$ and $\mathbf{x}_i^m = \{x_{ij} : r_{ij} = 0\}$ for all $i = 1, \dots, n$ and $j = 1, \dots, p$. Importantly, one can consider missing values as latent variables like \mathbf{Z} , as there are no distinctions between unknown parameters and missing variables in Bayesian statistics (Gelman et al. 1995). This allows us to model the missing data \mathbf{X}^m jointly with latent variable \mathbf{Z} and the observed data $\{\mathbf{X}^o, \mathbf{R}\}$.

Missingness was classified into three major categories, or mechanisms, in the seminal work by Little and Rubin (2002). These mechanisms are missing completely at random (MCAR), missing at random (MAR), and missing not at random (MNAR), and they satisfy the following relations: (a) MCAR: $p(\mathbf{r}_i|\mathbf{x}_i, \mathbf{z}_i, \phi) = p(\mathbf{r}_i|\phi)$, (b) MAR: $p(\mathbf{r}_i|\mathbf{x}_i, \mathbf{z}_i, \phi) = p(\mathbf{r}_i|\mathbf{x}_i^o, \phi)$, and (c) MNAR: $p(\mathbf{r}_i|\mathbf{x}_i, \mathbf{z}_i, \phi) = p(\mathbf{r}_i|\mathbf{x}_i^o, \mathbf{x}_i^m, \mathbf{z}_i, \phi)$. Here, ϕ denotes the unknown parameters pertaining to the missingness model $p(\mathbf{r}_i|\mathbf{x}_i, \mathbf{z}_i, \phi)$, where $\mathbf{r}_i = \{r_{i1}, \dots, r_{ip}\}$. We discuss forms of this model in Section 2.3.2. In the presence of missingness mask \mathbf{R} , the marginal log-likelihood may be written as

$$\log p_{\psi, \phi}(\mathbf{X}^o, \mathbf{R}) = \log \iint p_{\psi, \phi}(\mathbf{X}^o, \mathbf{X}^m, \mathbf{Z}, \mathbf{R}) d\mathbf{X}^m d\mathbf{Z}. \quad (5)$$

We may factor $p_{\psi, \phi}(\mathbf{X}^o, \mathbf{X}^m, \mathbf{Z}, \mathbf{R})$ using the selection model factorization (Diggle and Kenward 1994), which is written as $p_{\psi, \phi}(\mathbf{X}^o, \mathbf{X}^m, \mathbf{Z}, \mathbf{R}) = p_{\psi}(\mathbf{X}^o, \mathbf{X}^m|\mathbf{Z})p(\mathbf{Z})p(\mathbf{R}|\mathbf{X}, \mathbf{Z}, \phi)$, where $p(\mathbf{R}|\mathbf{X}, \mathbf{Z}, \phi) = \prod_{i=1}^n p(\mathbf{r}_i|\mathbf{x}_i, \mathbf{z}_i, \phi)$.

2.3.1. Ignorable Missingness

In a likelihood-based analysis under either MCAR or MAR, the missingness mechanism is considered to be “ignorable” such that the missingness mechanism need not be explicitly modeled in these cases (Rubin 1976; Little and Rubin 2002). Under ignorable missingness, the left hand side of (5) can be separated into $\log p_{\psi}(\mathbf{X}^o) + \log p_{\phi}(\mathbf{R}|\mathbf{X}^o)$, where $p_{\psi}(\mathbf{X}^o)$ is the marginal distribution of \mathbf{X}^o . Therefore, $p_{\phi}(\mathbf{R}|\mathbf{X}^o)$ need not be specified because inference on the parameters of interest pertaining to $p_{\psi}(\mathbf{X}^o)$ is independent of $p_{\phi}(\mathbf{R}|\mathbf{X}^o)$. Then, one aims to maximize the quantity

$$\begin{aligned} \log p_{\psi}(\mathbf{X}^o) &= \log \iint p_{\psi}(\mathbf{X}^o, \mathbf{X}^m, \mathbf{Z}) d\mathbf{X}^m d\mathbf{Z} \\ &= \log \int p_{\psi}(\mathbf{X}^o, \mathbf{Z}) d\mathbf{Z}. \end{aligned} \quad (6)$$

This quantity can be bounded below exactly as in Section 2.1.1, conditioning on just the observed data \mathbf{X}^o , rather than the full data \mathbf{X} . Existing methods typically take advantage of this simplification, and are shown to perform well under ignorable missingness. Details for these methods are given in Section A of the supplementary materials.

2.3.2. Non-Ignorable Missingness

In contrast, non-ignorable (or MNAR) missingness refers to the case where the missingness can be dependent on any unobserved values, including the missing entries \mathbf{x}_i^m . MNAR missingness can also be dependent on \mathbf{x}_i^o as well as latent values like \mathbf{Z} ,

and thus MNAR represents the most general and difficult case of missingness in practice. Here, we assume that \mathbf{R} is independent of \mathbf{Z} , as conditioning on such latent factors may be computationally redundant based on the assumed data generating process (Ibrahim 2001). In this setting, the missingness typically requires specification of a model for the missingness $p(\mathbf{R}|\mathbf{X}, \phi)$ (Stubbendick and Ibrahim 2003). Current VAE methods are only able to handle MCAR or MAR missingness, and there is no method to properly deal with the more difficult MNAR case. This issue is especially problematic because missingness in many real world applications have been posited to be non-ignorable (Beaulieu-Jones and 2016; O’Shea 2019).

There have been a number of ways to specify $p(\mathbf{R}|\mathbf{X}, \phi)$ in statistical literature. For example, Diggle and Kenward (1994) proposes a binomial model for the missing data mechanism:

$$p(\mathbf{R}|\mathbf{X}, \phi) = \prod_{i=1}^n \prod_{j_m=1}^p \left[p(r_{ij_m} = 1 | \mathbf{x}_i, \phi_{j_m}) \right]^{r_{ij_m}} \times \left[1 - p(r_{ij_m} = 1 | \mathbf{x}_i, \phi_{j_m}) \right]^{1-r_{ij_m}},$$

where $j_m = 1, \dots, p_{\text{miss}}$ indexes the p_{miss} features in \mathbf{X} that contain missingness, ϕ_{j_m} is the sets of coefficients pertaining to j_m th missingness model, and $p(r_{ij_m} = 1 | \mathbf{x}_i, \phi_{j_m})$ can be modeled straightforwardly by a logistic regression model, such that

$$\text{logit}[p(r_{ij_m} = 1 | \mathbf{x}_i, \phi_{j_m})] = \phi_{0j_m} + \mathbf{x}_i^o \phi_{1j_m} + \mathbf{x}_i^m \phi_{2j_m}, \quad (7)$$

where $\phi_{1j_m} = \{\phi_{1j_m,1}, \dots, \phi_{1j_m,p_{\text{obs}}}\}^T$ is a $p_{\text{obs}} \times 1$ vector of coefficients pertaining to the fully-observed features, $\phi_{2j_m} = \{\phi_{2j_m,1}, \dots, \phi_{2j_m,p_{\text{miss}}}\}^T$ is a $p_{\text{miss}} \times 1$ vector of coefficients pertaining to the missing features, and ϕ_{0j_m} is the intercept of the j_m^{th} missingness model.

2.4. NIMIWAE: IWAE with Nonignorable Missingness

We now propose a novel method to perform statistical learning and imputation using an IWAE in the presence of missing data (NIMIWAE), assuming missingness is nonignorable. We later show how this model can be simplified when missingness is assumed to be ignorable (IMIWAE). First, we specify a general form of the lower bound, in which we use the general IWAE framework to form a tighter bound on the marginal log-likelihood than the VAE ELBO. Let us define $q_\theta(\mathbf{Z}, \mathbf{X}^m)$ as the variational joint posterior pertaining to $(\mathbf{Z}, \mathbf{X}^m)$. We can factor this variational joint posterior as $q_\theta(\mathbf{Z}, \mathbf{X}^m) = q_{\theta_1}(\mathbf{Z}|\mathbf{X}^o)q_{\theta_2}(\mathbf{X}^m|\mathbf{Z}, \mathbf{X}^o, \mathbf{R})$. Here, for $k = 1, \dots, K$, we assume $\mathbf{Z}_k \stackrel{\text{iid}}{\sim} q_{\theta_1}(\mathbf{Z}|\mathbf{X}^o)$ similar to the traditional IWAE. We additionally assume iid latent variables $\mathbf{X}_k^m \stackrel{\text{iid}}{\sim} q_{\theta_2}(\mathbf{X}^m|\mathbf{Z}, \mathbf{X}^o, \mathbf{R})$ pertaining to the missing features, where each \mathbf{X}_k^m has dimensionality p_{miss} . Similar to traditional VAEs, we use the class of factorized variational posteriors, except now $q_\theta(\mathbf{Z}, \mathbf{X}^m) = \prod_{i=1}^n q_\theta(\mathbf{z}_i, \mathbf{x}_i^m)$ and $q_\theta(\mathbf{z}_i, \mathbf{x}_i^m) = q_{\theta_1}(\mathbf{z}_i|\mathbf{x}_i^o)q_{\theta_2}(\mathbf{x}_i^m|\mathbf{z}_i, \mathbf{x}_i^o, \mathbf{r}_i)$. Then, denoting \mathbf{z}_{ik} and \mathbf{x}_{ik}^m as the i th observation vectors of \mathbf{Z}_k and \mathbf{X}_k^m , respectively, we have $\mathbf{z}_{i1}, \dots, \mathbf{z}_{iK} \stackrel{\text{iid}}{\sim} q_{\theta_1}(\mathbf{z}_i|\mathbf{x}_i^o)$ and $\mathbf{x}_{i1}^m, \dots, \mathbf{x}_{iK}^m \stackrel{\text{iid}}{\sim} q_{\theta_2}(\mathbf{x}_i^m|\mathbf{z}_i, \mathbf{x}_i^o, \mathbf{r}_i)$. The form of the lower bound, which we call the NonIgnorably Missing Importance-Weighted

Auto Encoder bound, or “NIMIWAE bound”, is derived as follows:

$$\begin{aligned} \log p_{\psi, \phi}(\mathbf{X}^o, \mathbf{R}) &= \sum_{i=1}^n \log p_{\psi, \phi}(\mathbf{x}_i^o, \mathbf{r}_i) \\ &= \sum_{i=1}^n \log \left[\iint p_{\psi, \phi}(\mathbf{x}_i^o, \mathbf{x}_i^m, \mathbf{r}_i, \mathbf{z}_i) d\mathbf{z}_i d\mathbf{x}_i^m \right] \\ &= \sum_{i=1}^n \log \mathbb{E}_{(\mathbf{z}_{ik}, \mathbf{x}_{ik}^m) \sim q_\theta(\mathbf{z}_i, \mathbf{x}_i^m)} \left[\frac{1}{K} \sum_{k=1}^K \frac{p_{\psi, \phi}(\mathbf{x}_i^o, \mathbf{x}_{ik}^m, \mathbf{r}_i, \mathbf{z}_{ik})}{q_{\theta_1}(\mathbf{z}_{ik}|\mathbf{x}_i^o)q_{\theta_2}(\mathbf{x}_{ik}^m|\mathbf{z}_{ik}, \mathbf{x}_i^o, \mathbf{r}_i)} \right] \\ &\geq \sum_{i=1}^n \mathbb{E}_{(\mathbf{z}_{ik}, \mathbf{x}_{ik}^m) \sim q_\theta(\mathbf{z}_i, \mathbf{x}_i^m)} \log \left[\frac{1}{K} \sum_{k=1}^K \frac{p_{\psi, \phi}(\mathbf{x}_i^o, \mathbf{x}_{ik}^m, \mathbf{r}_i, \mathbf{z}_{ik})}{q_{\theta_1}(\mathbf{z}_{ik}|\mathbf{x}_i^o)q_{\theta_2}(\mathbf{x}_{ik}^m|\mathbf{z}_{ik}, \mathbf{x}_i^o, \mathbf{r}_i)} \right] \\ &= \mathcal{L}_K^{\text{NIMIWAE}}. \end{aligned} \quad (8)$$

As explained in Section 2.3, we use the selection model factorization of the joint distribution of $\{\mathbf{x}_i, \mathbf{r}_i, \mathbf{z}_i\}$, such that

$$p_{\psi, \phi}(\mathbf{x}_i^o, \mathbf{x}_i^m, \mathbf{r}_i, \mathbf{z}_i) = p_\psi(\mathbf{x}_i^o, \mathbf{x}_i^m|\mathbf{z}_i)p_\phi(\mathbf{r}_i|\mathbf{x}_i^o, \mathbf{x}_i^m).$$

Here, ψ denotes the weights and biases of the encoder and decoder neural networks, and ϕ denotes the weights and biases of the missingness network that learns the parameters of the missingness model.

Applying the above factorizations to (8), and estimating the expectations in (8) by sampling from $q_\theta(\mathbf{z}_i, \mathbf{x}_i^m)$ we obtain the estimate of the NIMIWAE bound:

$$\hat{\mathcal{L}}_K^{\text{NIMIWAE}} = \sum_{i=1}^n \log \left[\frac{1}{K} \sum_{k=1}^K \frac{p_\psi(\mathbf{x}_i|\tilde{\mathbf{z}}_{ik})p(\tilde{\mathbf{z}}_{ik})p_\phi(\mathbf{r}_i|\mathbf{x}_i^o, \tilde{\mathbf{x}}_{ik}^m)}{q_{\theta_1}(\tilde{\mathbf{z}}_{ik}|\mathbf{x}_i^o)q_{\theta_2}(\tilde{\mathbf{x}}_{ik}^m|\tilde{\mathbf{z}}_{ik}, \mathbf{x}_i^o, \mathbf{r}_i)} \right], \quad (9)$$

where $\{\tilde{\mathbf{z}}_{ik}, \tilde{\mathbf{x}}_{ik}^m\}$ are samples of $\{\mathbf{z}_i, \mathbf{x}_i^m\}$ that are drawn via ancestral sampling (Bishop 2006) from $q_{\theta_1}(\tilde{\mathbf{z}}_{ik}|\mathbf{x}_i^o)$ and $q_{\theta_2}(\tilde{\mathbf{x}}_{ik}^m|\tilde{\mathbf{z}}_{ik}, \mathbf{x}_i^o, \mathbf{r}_i)$, respectively, and the NIMIWAE bound is optimized using the Adam optimizer (Kingma and Ba 2014). In NIMIWAE, we have four neural networks $f_\psi(\mathbf{z}_i)$, $g_{\theta_1}(\mathbf{x}_i^o)$, $g_{\theta_2}(\mathbf{x}_i^o, \mathbf{r}_i, \mathbf{z}_i)$, and $h_\phi(\mathbf{x}_i)$, which respectively output the parameters pertaining to $p_\psi(\mathbf{x}_i|\mathbf{z}_{ik})$, $q_{\theta_1}(\mathbf{z}_{ik}|\mathbf{x}_i^o)$, $q_{\theta_2}(\mathbf{x}_{ik}^m|\mathbf{z}_{ik}, \mathbf{x}_i^o, \mathbf{r}_i)$, and $p_\phi(\mathbf{r}_i|\mathbf{x}_i^o, \mathbf{x}_{ik}^m)$.

The quantity $h_\phi(\mathbf{x}_i)$, which we call the “missingness network”, outputs the distributional parameters to the missingness model $p_\phi(\mathbf{r}_i|\mathbf{x}_i^o, \tilde{\mathbf{x}}_{ik}^m)$ specifically given in (9). Furthermore, by omitting this network and its contribution to the NIMIWAE bound altogether, one can attain an ignorably missing version of the NIMIWAE method (IMIWAE), which would be more suitable for use when assuming MCAR and MAR missingness. We explore the empirical performance of each of these models under misspecification of the missingness mechanism in Section 3. An illustration of $h_\phi(\mathbf{x}_i)$ is given in Figure 1.

Historically, the set of features for the p_{miss} logistic regression models from (7) need to be carefully pre-specified, usually based upon prior information (Little and Rubin 2002). Prior work has shown that overparameterization of the missingness model can lead to identifiability issues and divergence in EM-based maximization procedures (Ibrahim and Molenberghs 2009). Our proposed method allows users to similarly pre-specify a subset of features in the missingness network, which can be used

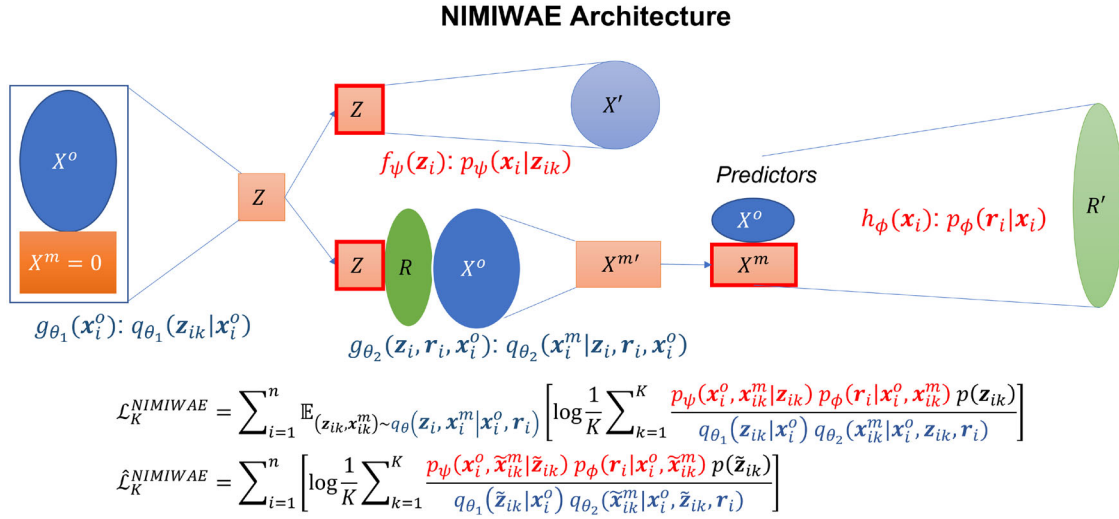


Figure 1. Architecture of proposed NIMIWAE method. Dark colored nodes ($X^o, R, X^m = 0$) represent deterministic values, lightly colored nodes (Z, X^m) represent learned distributional parameters, and outlined (in red) nodes represent sampled values. Orange cells correspond to latent variables Z and X^m . Z_1, \dots, Z_K and X_1^m, \dots, X_K^m are sampled from their respective variational posteriors $q_{\theta_1}(Z|X^o)$ and $q_{\theta_2}(X^m|Z, R, X^o)$. Below is the NIMIWAE bound ($\mathcal{L}_K^{NIMIWAE}$) and the estimate of the NIMIWAE bound ($\hat{\mathcal{L}}_K^{NIMIWAE}$), which is optimized via stochastic gradient descent.

for sensitivity analyses, which is very commonly done for real-life settings. When prior information on the missingness is not available, however, it is unclear how overparameterization will affect model fitting and performance in this setting. Therefore, in Section 3, we evaluate empirically the use of all p features in the missingness network when assuming MNAR missingness using statistical simulation.

Additional details regarding the NIMIWAE algorithm are outlined in Section A of the supplementary materials.

2.4.1. Multiple Imputation

Following training, NIMIWAE can provide point estimates for $\mathbb{E}[\mathbf{x}_i^m | \mathbf{x}_i^o, \mathbf{r}_i]$, defined as the expected value of the missing features given the observed data and the mask for the i th observation under MNAR. The same NIMIWAE model can also perform multiple imputation of the incomplete training dataset to facilitate inference in downstream statistical models. We first note that

$$\begin{aligned} \mathbb{E}[\mathbf{x}_i^m | \mathbf{x}_i^o, \mathbf{r}_i] &= \int \mathbf{x}_i^m p_{\psi, \phi}(\mathbf{x}_i^m | \mathbf{x}_i^o, \mathbf{r}_i) d\mathbf{x}_i^m \\ &= \iint \mathbf{x}_i^m p_{\psi, \phi}(\mathbf{x}_i^m, \mathbf{z}_i | \mathbf{x}_i^o, \mathbf{r}_i) d\mathbf{z}_i d\mathbf{x}_i^m \\ &= \iint \mathbf{x}_i^m \frac{p_{\psi, \phi}(\mathbf{x}_i^m, \mathbf{z}_i | \mathbf{x}_i^o, \mathbf{r}_i)}{p_{\psi, \phi}(\mathbf{x}_i^o, \mathbf{r}_i)} d\mathbf{z}_i d\mathbf{x}_i^m \\ &= \iint \mathbf{x}_i^m \frac{p_{\psi}(\mathbf{x}_i | \mathbf{z}_i) p(\mathbf{z}_i) p_{\phi}(\mathbf{r}_i | \mathbf{x}_i^o, \mathbf{x}_i^m)}{p_{\psi, \phi}(\mathbf{x}_i^o, \mathbf{r}_i)} d\mathbf{z}_i d\mathbf{x}_i^m. \end{aligned}$$

Then, we may estimate this integral by self-normalized importance sampling. We use the proposal density $q_{\theta_1}(\mathbf{z}_i | \mathbf{x}_i^o) q_{\theta_2}(\mathbf{x}_i^m | \mathbf{z}_i, \mathbf{x}_i^o, \mathbf{r}_i)$, and consider $p_{\psi, \phi}(\mathbf{x}_i^o, \mathbf{r}_i)$ as some unknown constant. Then, we define the quantities

$$w_{ik} = \frac{s_{ik}}{s_{i1} + \dots + s_{iK}}, \text{ and } s_{ik} = \frac{p_{\psi}(\mathbf{x}_i | \tilde{\mathbf{z}}_{ik}) p(\tilde{\mathbf{z}}_{ik}) p_{\phi}(\mathbf{r}_i | \mathbf{x}_i^o, \tilde{\mathbf{x}}_{ik}^m)}{q_{\theta_1}(\tilde{\mathbf{z}}_{ik} | \mathbf{x}_i^o) q_{\theta_2}(\tilde{\mathbf{x}}_{ik}^m | \tilde{\mathbf{z}}_{ik}, \mathbf{x}_i^o, \mathbf{r}_i)}$$

for $k = 1, \dots, K$, with 1 sample drawn from the variational posterior of each latent variable \mathbf{z}_{ik} and \mathbf{x}_{ik}^m to compute s_{ik} ,

where w_{ik} is defined as standardized “importance weights” (Mattei and Frellsen 2019). Using these weights we may estimate $\mathbb{E}[\mathbf{x}_i^m | \mathbf{x}_i^o, \mathbf{r}_i] \approx \sum_{k=1}^K w_{ik} \tilde{\mathbf{x}}_{ik}^m$. Then, the process can be repeated for each observation $i = 1, \dots, n$.

In the MCAR or MAR case, one can similarly estimate $\mathbb{E}[\mathbf{x}_i^m | \mathbf{x}_i^o]$ using the fitted IMIWAE model. By following a similar derivation using the proposal density $q_{\theta_1}(\mathbf{z}_i | \mathbf{x}_i^o) q_{\theta_2}(\mathbf{x}_i^m | \mathbf{z}_i, \mathbf{x}_i^o)$, we obtain the same approximation $\mathbb{E}[\mathbf{x}_i^m | \mathbf{x}_i^o] \approx \sum_{k=1}^K w_{ik} \tilde{\mathbf{x}}_{ik}^m$, with w_{ik} defined as before, but with a slightly different form for s_{ik} :

$$s_{ik} = \frac{p_{\psi}(\mathbf{x}_i | \tilde{\mathbf{z}}_{ik}) p(\tilde{\mathbf{z}}_{ik})}{q_{\theta_1}(\tilde{\mathbf{z}}_{ik} | \mathbf{x}_i^o) q_{\theta_2}(\tilde{\mathbf{x}}_{ik}^m | \tilde{\mathbf{z}}_{ik}, \mathbf{x}_i^o)}.$$

Given these weights w_{ik} , we may now also produce Q multiply-imputed datasets using the sampling importance resampling (SIR) algorithm (Smith and Gelfand 1992). We first construct a set of K candidate draws and corresponding weights using the procedure described above (default $K = 10 \times Q$ draws), and then perform a weighted resample of size Q with replacement from this set of draws to obtain approximate draws from $p_{\psi}(\mathbf{x}_i^m | \mathbf{x}_i^o, \mathbf{r}_i)$, for each observation $i = 1, \dots, n$. We may then use techniques (Rubin 2004) to pool the estimates obtained from a candidate regression model fit on each imputed dataset, and obtain pooled point estimates and standard errors that account for the uncertainty due to the imputation. In our analyses, we used NIMIWAE to construct $Q = 50$ multiply-imputed datasets by drawing $K = 500$ times from $q_{\theta_2}(\mathbf{x}_i^m | \mathbf{z}_i, \mathbf{x}_i^o, \mathbf{r}_i)$ for each $i = 1, \dots, n$ after the model was trained.

3. Numerical Results

3.1. Simulated Data

We use statistical simulation to evaluate the imputation performance of our proposed NIMIWAE and IMIWAE methods and under the assumption of MCAR, MAR, and MNAR missingness. We note that for each of these simulations, we use all

p features in NIMIWAE's missingness network. We also compared this performance to state-of-the-art missing data methods in machine learning that claim to handle ignorable missingness patterns: HIVAE (Nazabal et al. 2018), VAEAC (Ivanov, Figurnov, and Vetrov 2019), MIWAE (Mattei and Frellsen 2019), in addition to the popular MICE method (Van Buuren and Groothuis-Oudshoorn 2011) and a naïve mean imputation method. We also included MissForest (Stekhoven and Buhlmann 2011) in analyses where the model could be fit in a CPU with 32 GB of memory. For all simulations, we divided the full data into training and validation sets with ratio 8:2. For methods that require hyperparameter tuning, each method was trained on the training set for a given set of hyperparameters, the performance of this model was then evaluated on the validation set, and finally the training set was imputed using the model pertaining to the optimal set of hyperparameters (based up on validation set performance). For methods that required no hyperparameter tuning, we directly imputed just the training set, for consistency across all methods. In Sections 3.1.1–3.1.2, we evaluate performance on fully synthetic data. Then, in Section 3.1.3, we simulate missingness into existing datasets from the UCI machine learning repository to preserve non-linearity and interactions between features previously observed. The simulation setup and performance criteria are described in the subsequent sections.

3.1.1. Simulation Setup

We first evaluate the performance of each method on completely synthetic data. Here we assume \mathbf{X} is generated such that $\mathbf{X} = \mathbf{Z}\mathbf{W} + \mathbf{B}$, where \mathbf{W} and \mathbf{B} are matrices of dimensions $d \times p$ and $n \times p$, respectively, $\mathbf{Z} \sim N_d(\mathbf{0}, \mathbf{I})$, $W_{lj} \sim N(0, 0.5)$, and $B_{ij} \sim N(0, 1)$ for $i = 1, \dots, n, j = 1, \dots, p$, and $l = 1, \dots, d$.

We then simulate the missingness mask matrix \mathbf{R} such that 30% of features are partially observed, and 50% of the observations for each of these features are missing. We generate r_{ij} from the Bernoulli distribution with probability equal to $p(r_{ijm} = 1 | \mathbf{x}_i, \boldsymbol{\phi})$, such that $\text{logit}[p(r_{ijm} = 1 | \mathbf{x}_i, \boldsymbol{\phi})] = \phi_0 + \phi_1 \mathbf{x}_i^o + \phi_2 \mathbf{x}_i^m$. Here, we assume that $j_m = 1, \dots, p_{\text{miss}}$ index the missing features, $\phi_1 = \{\phi_{11}, \dots, \phi_{1, p_{\text{obs}}}\}$ are the coefficients pertaining to the fully observed features, and $\phi_2 = \{\phi_{21}, \dots, \phi_{2, p_{\text{miss}}}\}$ are those pertaining to the partially observed features, and p_{obs} and p_{miss} are the total number of features that are fully and partially observed, respectively. We set $p_{\text{miss}} = \lfloor 0.3 * p \rfloor$ and $p_{\text{obs}} = p - p_{\text{miss}}$. We drew nonzero values of ϕ_1 and ϕ_2 from the log-normal distribution with mean $\mu_\phi = 5$, with log standard deviation $\sigma_\phi = 0.2$.

To evaluate the impact of the misspecification of the missingness mechanism on model performance, r_{ijm} was simulated under each mechanism as follows: (a) MCAR: $\{\phi_1, \phi_2\} = \mathbf{0}$, (b) MAR: Same as MCAR except $\phi_{1j_o} \neq 0$ for one randomly selected completely-observed feature j_o where $j_o = p_{\text{miss}} + 1, \dots, p$, and (c) MNAR: Same as MCAR except $\phi_{2j_m} \neq 0$. In this way, for each MAR or MNAR feature, the missingness is dependent on just one feature. In each case, we used ϕ_0 to adjust the overall expected rate of missingness of each feature to approximately 50%.

Lastly, we simulated a binary response variable assuming $\Pr(\mathbf{y} = 1 | \mathbf{X}) = \text{Sigmoid}(\beta_0 + \mathbf{X}\boldsymbol{\beta})$, where \mathbf{y} is a binary response

variable, $\boldsymbol{\beta} = \{\beta_1, \dots, \beta_p\}$ are the set of regression coefficients, and β_0 is the intercept. In many applications, it is of interest to use the features in \mathbf{X} to predict some outcome variable \mathbf{y} when \mathbf{X} is only partially observed. Therefore, the ability accurately estimate $\boldsymbol{\beta}$ is also of importance in the presence of missingness. Multiple imputation methods like MICE can be used to perform coefficient estimation by pooling the coefficient estimates from logistic regression models fitted on each individual multiply-imputed dataset. We similarly perform coefficient estimation using multiply-imputed datasets from the SIR algorithm within NIMIWAE, allowing for direct comparisons in coefficient estimation with MICE in estimating $\boldsymbol{\beta}$.

We vary n , p , and d such that $n = \{10,000; 100,000\}$, $p = \{25; 100\}$ features, and $d = \{2; 8\}$. We simulated 10 datasets per simulation condition, spanning various missingness mechanisms and values for $\{n, p, d\}$. We fix the values of $\boldsymbol{\beta}$ at 0.25 for each feature, and adjusted β_0 to ensure equal proportions for the two binary classes in \mathbf{Y} . We measured imputation performance by calculating the average L1 distance between true and imputed masked values in \mathbf{X} . Letting $\hat{\mathbf{X}}^m$ denote the imputed masked values of the true \mathbf{X}^m values of the missing entries, we denote the average L1 distance is simply $\frac{|\hat{\mathbf{X}}^m - \mathbf{X}^m|}{N_{\text{miss}}}$, where N_{miss} is the total number of missing entries in the dataset. To assess the ability of multiple imputation methods in performing coefficient estimation, we reported the percent bias (PB) of these pooled estimates compared to the truth, averaged across the p features, that is $PB = \frac{1}{p} \sum_{j=1}^p \frac{|\beta_j - \hat{\beta}_j|}{|\beta_j|}$.

We also evaluated the methods' performance under nonlinearity of the underlying missingness mechanism. In particular, we simulated data with missingness model $\text{logit}[p(r_{ijm} = 1 | \mathbf{x}_i, \boldsymbol{\phi})] = \phi_0 + \phi_1 (\mathbf{x}_i^o)^2 + \phi_2 (\mathbf{x}_i^m)^2$, where the covariate in each missingness model is now the square of the respective feature. We evaluated performance under 10 simulated datasets with $n = 100,000$ and $p = 25$.

3.1.2. Simulation Results

Figure 2 shows the results pertaining to $n = 100,000$. Error bars represented the variability in performance across 10 simulations. Despite the overparameterized missingness model, NIMIWAE consistently yields improved imputation performance compared to other methods under MNAR missingness, while yielding an average L1 that is comparable to other methods in the MCAR and MAR cases, under these conditions assuming large sample sizes. This can be exceptionally useful for many real data applications, where the true covariates of the missingness model may not be known a priori. In estimating $\boldsymbol{\beta}$, we see that NIMIWAE yields a lower average percent bias under MNAR missingness, while MICE, IMIWAE, and NIMIWAE perform similarly under MCAR or MAR missingness.

Figure 3 shows the results pertaining to $n = 10,000$ using the alternative initialization method described in Section 2.4, due to the smaller sample size in this setting. We see that NIMIWAE performs best in imputing MNAR values, and in estimating the coefficients under MNAR missingness, and still performs comparably well to other methods in the MCAR and MAR cases. The results of the analyses on $n = 10,000$ simulations with default initialization can be found in Section B of the supplementary materials. We found that the default method

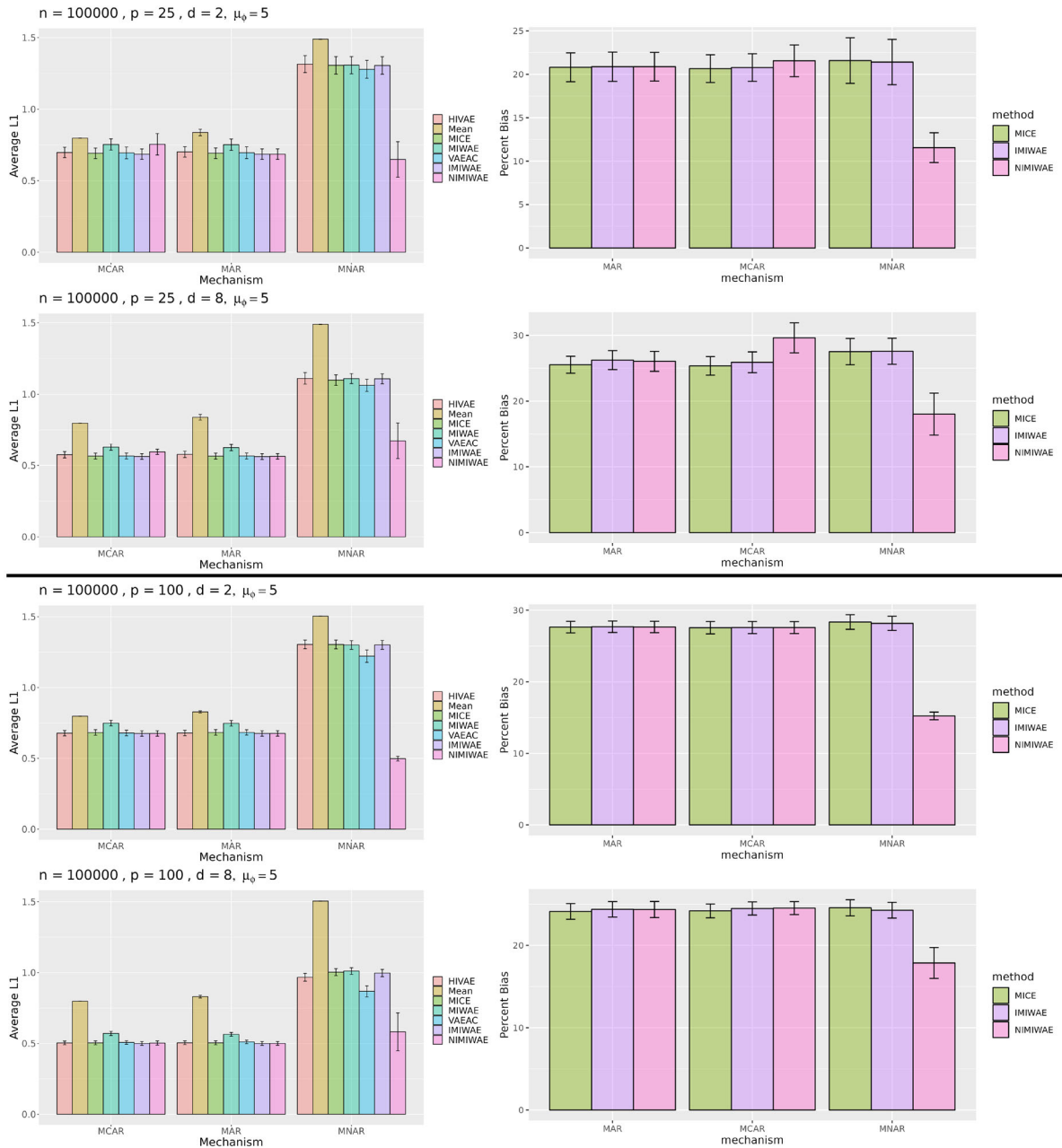


Figure 2. Average L1 distance between true and imputed values for missing entries (left) and percent bias of pooled coefficient estimates (right) for $p = 25$ (top 4) and $p = 100$ (bottom 4) features, stratified by d . NIMIWAE outperforms all methods in imputing MNAR missing values, while performing comparably to other methods in imputing MCAR and MAR values. Here, $n = 100,000$, $\mu_\phi = 5$, and error bars show the variability of each metric across 10 reps. Weights and biases were initialized by using the default semi-orthogonal matrix method.

initialized weights for missing features too small for the network to recover in the MNAR case, under the lower sample size and high dimensionality setting. Based upon these results, we recommend the alternative weight initialization for smaller sample sizes ($n \leq 10,000$) with large dimensionality ($p \geq 100$), while the default initialization may suffice when the sample size is large, or if the dimensionality of the data is smaller. Alternatively, one may greatly improve imputation performance by narrowing down the features that are input into the missingness network using some prior knowledge or assumptions on the mechanism of missingness. In practical applications, it is unknown which mechanism may truly underly the data, and, as usual, sensitivity analyses, where one varies the assumed mechanism, is still

important. For example, imputing data via IMIWAE may be helpful and be more efficient under MCAR or MAR than via NIMIWAE.

It is also worth noting that the computational time for NIMIWAE and IMIWAE were comparable to that of the other deep learning methods (HIVAE, VAEAC, MIWAE). The majority of the computing time for each method was consumed during hyperparameter tuning. Significantly, as these compared methods are all generative models, multiple imputation and downstream statistical analysis was very fast, compared to parameter tuning and model training.

Overall, these simulations confirm that existing methods can impute MCAR and MAR missingness with a reasonable degree

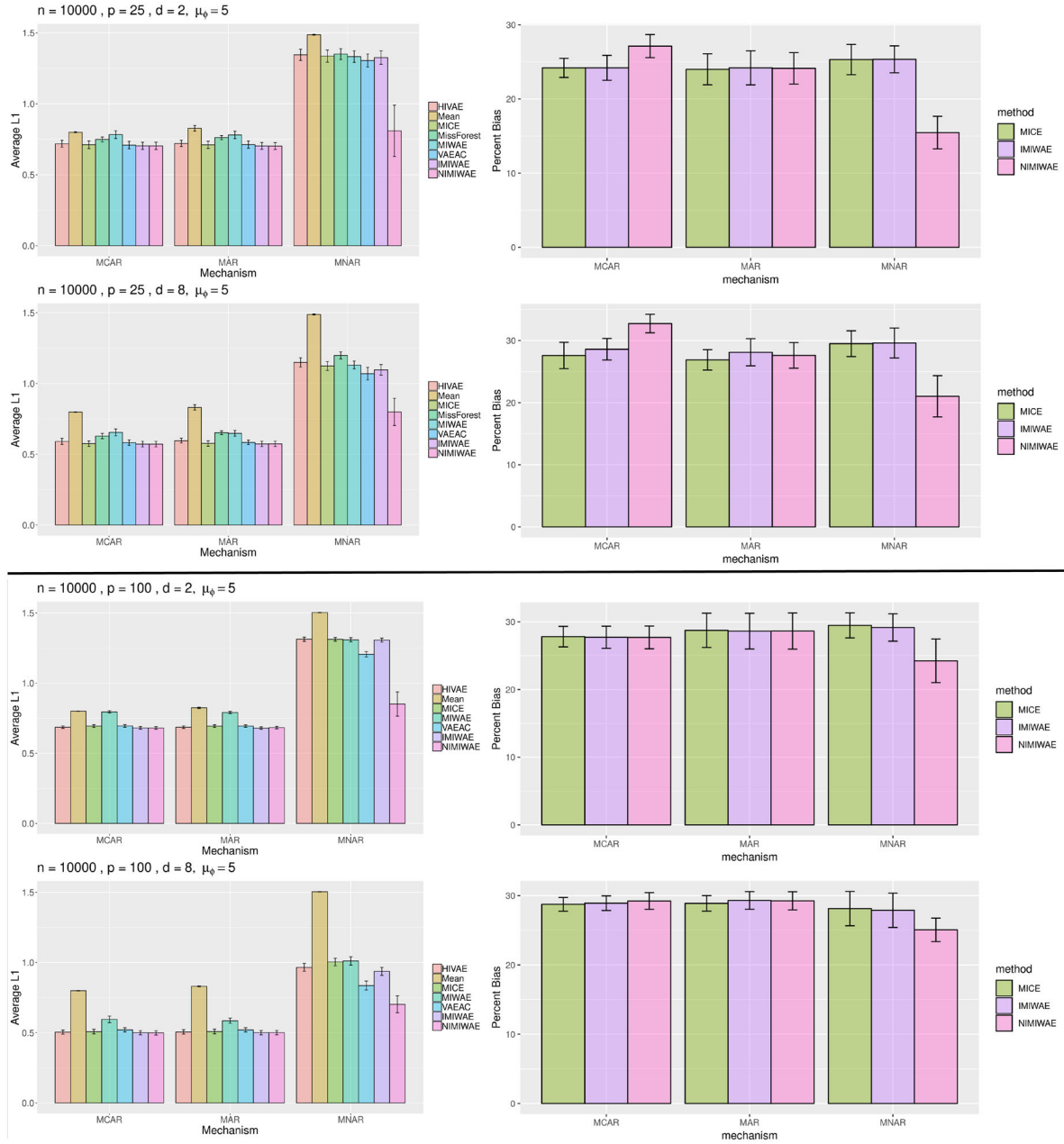


Figure 3. Average L1 distance between true and imputed values for missing entries (left) and percent bias of pooled coefficient estimates (right) for $p = 25$ (top 4) and $p = 100$ (bottom 4) features, stratified by d . NIMIWAE outperforms all methods in imputing MNAR missing values, while performing comparably to other methods in imputing MCAR and MAR values. Here, $n = 10,000$, $\mu_\phi = 5$, and error bars show the variability of each metric across 10 reps. Weights and biases were initialized by using the alternative method, as described in Section 2.4

of accuracy, but they break down under MNAR missingness. Our NIMIWAE method is able to adequately impute MNAR missing values, while still imputing MCAR and MAR missing values with a comparable degree of accuracy as existing methods geared specifically toward the ignorable mechanisms of missingness. Additionally, we note that MissForest was able to be run only on the $n = 10,000$ and $p = 25$ simulation case due to memory constraints, and yielded very poor imputation performance in the MNAR case, like the other ignorably missing methods. We also show empirically that using all p features in NIMIWAE’s missingness network can yield reasonable performance, especially when the number of samples is large. We postulate that this may be due to the fact that neural networks have been shown

to generalize well despite severe overparameterization (Poggio, Banburski, and Liao 2020). Still, the specification of a smaller model that is closer to the truth may improve the accuracy of imputations, especially under smaller sample sizes (Du et al. 2021).

Finally, we evaluated the compared methods when the underlying missingness mechanism was nonlinear. Figure 4 shows the results pertaining to $n = 100,000$ and $p = 25$. NIMIWAE’s missingness network can capture such nonlinear dependencies between the missingness and the features, allowing NIMIWAE to consistently outperform all compared methods, even under this more complex underlying missingness mechanism.

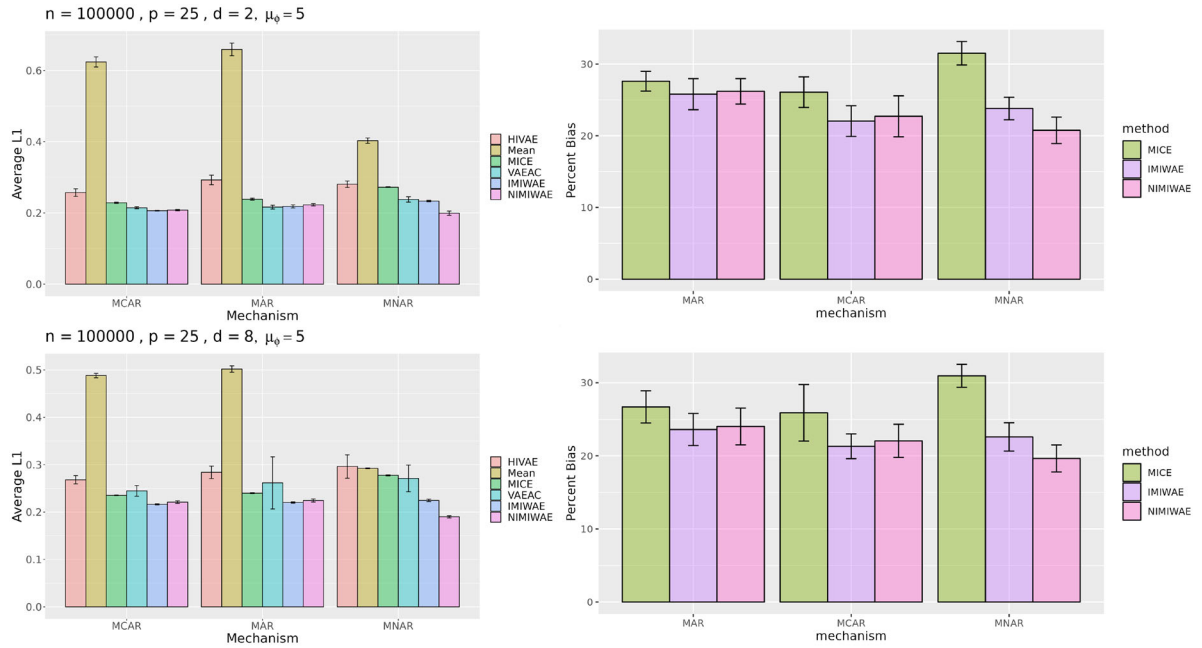


Figure 4. Average L1 distance between true and imputed values for missing entries (left) and percent bias of pooled coefficient estimates (right) for $p = 25$ features, stratified by d , with an underlying nonlinear missingness mechanism. NIMIWE still outperforms all methods in imputing MNAR missing values, while performing comparably to other methods in imputing MCAR and MAR values. Here, $n = 100,000$, $\mu_\phi = 5$, and error bars show the variability of each metric across 10 reps.

3.1.3. UCI Machine Learning Datasets

Next, we analyzed how accurately these methods can impute missing values that may be present in real data. Since we are unable to obtain the true missing values in real datasets, we selected 6 different completely observed datasets from the UCI Machine Learning Repository, and simulated missingness according to each mechanism as in the fully synthetic datasets. Here, we masked half of the features, such that $p_{\text{miss}} = \text{floor}(0.5 * p)$. Additional details on these datasets and how they may be obtained can be found in Section B of the supplementary materials, and the values of the hyperparameters that were tuned can be found in Section C of the supplementary materials.

Table 1 shows the average L1 distance between true and imputed missing values for each of the six UCI datasets, with each simulated mechanism of missingness. We see that NIMIWE again performs best across all methods in accurately imputing MNAR missing values. Overall, MissForest performed best in imputing MCAR missing values in the smaller UCI datasets, but this method was too memory-consuming to be trained on the significantly larger *hepmass* and *power* datasets. As in the simulated data, mean imputation is consistently one of the least accurate methods of imputation. Also, whereas MICE performed very well under MCAR and MAR in the simulated data, deep learning methods like VAEAC and IMIWAE consistently yield more accurate imputations here. This may be due to the fact that MICE uses a fully conditional linear model for imputation, it may therefore perform suboptimally when the true relationships between features are nonlinear. Also, we see that NIMIWE imputes values slightly less accurately than IMIWAE when the missingness is MCAR or MAR, since in these cases NIMIWE explicitly estimates the missingness network when it is not necessary.

3.2. Physionet 2012 Challenge Dataset

Finally, we analyzed the Physionet 2012 Challenge data using each of the compared methods. This data consisted of 114 features of 12,000 patients after pre-processing, and the details of the acquisition and preparation of the data can be found in Section C of the supplementary materials. We performed a qualitative analysis of each imputed dataset, highlighting differences between the results of downstream regression models fitted on mortality, assuming non-ignorable (NIMIWE) and ignorable missingness (Ibrahim et al. 2005; Ibrahim and Molenberghs 2009). This is because, in contrast to our simulations, the true values of the missing entries are not available to directly assess imputation performance. Additionally, the missingness mechanism itself is generally not “testable” by the observed data in practice (Ibrahim, Lipsitz, and Chen 1999). Here, we used the alternative initialization scheme of NIMIWE, due to the limited sample size and large dimension of the data. Additionally, we added “supervised” versions of NIMIWE, IMIWAE, and MICE in our comparisons, where the response of interest (mortality) is included with the features of the data during training and imputation, such that the response variable may inform the multiple imputations. After model fitting, we pool the estimates of the fitted models on the multiply imputed datasets, and compute standard errors across imputations.

Based on the imputed datasets from these methods, we fit a logistic regression model with post-baseline mortality as the binary response, with the $p = 114$ baseline features of this dataset as the covariates. We report details of the covariates with the top 10 largest effects on mortality when using the multiply imputed datasets from the supervised non-ignorably missing NIMIWE model in Figure 5. We found some significant differences in the results of the logistic regression

Table 1. Average L1 distance between true and imputed missing values in various datasets, under different mechanisms of simulated missingness.

Dataset		HIVAE	Mean	MICE	MF	MIWAE	VAEAC	NIMI	IMI
banknote	MCAR	1.62	2.00	1.63	0.95	1.37	1.32	1.49	1.19
$n = 1372$	MAR	2.00	2.27	2.43	1.88	2.90	1.76	1.89	1.96
	MNAR	3.39	3.92	3.78	3.18	3.10	3.46	1.46	3.30
concrete	MCAR	47.54	51.28	29.97	25.57	40.53	33.03	42.12	33.69
$n = 1030$	MAR	67.37	60.00	44.77	53.19	66.35	61.09	55.82	57.56
	MNAR	59.48	95.85	68.24	79.87	76.79	70.12	47.46	74.44
hepmass	MCAR	0.75	0.82	0.74	NA	0.80	0.68	0.78	0.69
$n = 525,123$	MAR	0.76	0.84	0.75	NA	0.84	0.72	0.72	0.71
	MNAR	1.41	1.54	1.40	NA	1.36	1.30	0.99	1.36
power	MCAR	0.54	0.66	0.50	NA	0.59	0.48	0.56	0.49
$n = 1M$	MAR	0.61	0.75	0.56	NA	0.67	0.54	0.78	0.57
	MNAR	0.75	1.14	0.84	NA	0.86	0.79	0.73	0.79
red	MCAR	1.15	1.64	1.04	0.86	1.10	1.10	1.08	0.98
$n = 1599$	MAR	1.23	1.67	1.16	0.98	1.30	1.29	1.09	1.06
	MNAR	2.12	3.24	2.46	2.06	1.87	2.73	0.90	1.74
white	MCAR	2.14	2.61	1.97	1.60	2.18	1.93	2.16	1.88
$n = 4898$	MAR	2.28	2.63	1.99	1.69	2.15	2.11	1.89	1.89
	MNAR	4.42	5.36	4.21	4.27	4.46	4.10	3.47	4.70

NOTE: Best imputation performance (lowest average L1) in each row is highlighted in red. Proportion of missing entries was fixed at 50% per feature, with 50% of the features containing missingness. We see that NIMIWAE (NIMI) consistently performs best in imputing MNAR missingness, while performance of the “Ignorable” IMIWAE (IMI) model is comparable to other methods under MCAR and MAR. Although MissForest (MF) claims superiority in MCAR and some MAR cases in the smaller datasets, it was not scalable to larger datasets like hepmass and power.

	NIMIWAE _{sup}	NIMIWAE _{unsup}	IMIWAE _{sup}	IMIWAE _{unsup}	MICE _{sup}	MICE _{unsup}	MIWAE	HIVAE	VAEAC	MEAN	MF
CSRU	-0.822 (0.169)	-0.862 (0.169)	-0.958 (0.165)	-0.953 (0.166)	-0.987 (0.213)	-1.011 (0.172)	-0.881 (0.165)	-0.926 (0.16)	-0.904 (0.17)	-0.89 (0.16)	-1.03 (0.166)
CCU	-0.359 (0.131)	-0.401 (0.132)	-0.404 (0.124)	-0.429 (0.127)	-0.497 (0.155)	-0.454 (0.129)	-0.424 (0.124)	-0.446 (0.122)	-0.481 (0.132)	-0.391 (0.121)	-0.419 (0.124)
Temp _{highest}	0.225 (0.074)	0.24 (0.074)	0.167 (0.075)	0.149 (0.075)	0.105 (0.089)	0.151 (0.077)	0.142 (0.075)	0.16 (0.077)	0.165 (0.076)	0.165 (0.075)	0.168 (0.075)
GCS _{last}	-0.184 (0.021)	-0.186 (0.021)	-0.148 (0.021)	-0.154 (0.021)	-0.17 (0.028)	-0.174 (0.021)	-0.169 (0.02)	-0.174 (0.021)	-0.17 (0.021)	-0.171 (0.021)	-0.167 (0.021)
Creatinine _{last}	-0.17 (0.079)	-0.171 (0.079)	-0.175 (0.078)	-0.189 (0.079)	-0.159 (0.095)	-0.171 (0.082)	-0.187 (0.078)	-0.172 (0.08)	-0.198 (0.08)	-0.177 (0.08)	-0.183 (0.079)
Lactate _{last}	0.162 (0.031)	0.137 (0.031)	0.14 (0.031)	0.149 (0.033)	0.198 (0.114)	0.119 (0.039)	0.149 (0.032)	0.156 (0.036)	0.163 (0.036)	0.164 (0.035)	0.152 (0.034)
Bilirubin _{last}	0.149 (0.035)	0.123 (0.032)	0.056 (0.021)	0.073 (0.022)	0.161 (0.077)	0.085 (0.034)	0.085 (0.024)	0.156 (0.037)	0.098 (0.031)	0.15 (0.037)	0.144 (0.033)
K _{last}	0.141 (0.07)	0.136 (0.071)	0.089 (0.077)	0.082 (0.077)	0.074 (0.094)	0.061 (0.079)	0.097 (0.077)	0.068 (0.082)	0.065 (0.08)	0.085 (0.077)	0.046 (0.077)
pH _{last}	0.134 (0.091)	-0.004 (0.067)	1.012 (0.85)	0.828 (0.871)	-0.542 (1.882)	0.746 (1.087)	0.062 (0.905)	0.61 (0.961)	0.296 (0.786)	0.736 (0.941)	0.759 (1.026)
Gender	-0.131 (0.077)	-0.116 (0.077)	-0.042 (0.079)	-0.023 (0.079)	-0.093 (0.153)	-0.044 (0.09)	0.012 (0.087)	-0.054 (0.084)	-0.046 (0.086)	-0.037 (0.081)	-0.035 (0.079)

Figure 5. Table of coefficient estimates (and standard errors) of covariates with the top 10 magnitudes of estimates via $NIMIWAE_{sup}$, from fitting a logistic regression model with imputed datasets from each method. Results from multiple imputation methods NIMIWAE, IMIWAE, and MICE (first six columns) are based on 50 multiply imputed datasets, and reflect pooled coefficient estimates and standard errors. For fair comparison, we also included results from single imputation methods (last five columns). Here, IMIWAE is the ignorable version of NIMIWAE.

based on the different imputed datasets. For example, we found that $NIMIWAE_{sup}$ uncovered a stronger effect of the variables $Temp_{highest}$, K_{last} and $Gender$ compared to other methods. These factors have been studied for their association with mortality in ICU patients. Specifically, gender has been studied for having a potentially significant effect on mortality in critically-ill patients (Mahmood, Eldeirawi, and Wahidi 2012; Larsson et al. 2019). Additionally, body temperature (Schell-Chaple et al. 2015) and irregular potassium levels (Tongyoo, Viarasilpa, and Permpikul 2018) have both been known to be associated with an increased risk of in-hospital mortality.

Figure 6 shows the imputed values of two variables that had significantly different missingness rates in surviving versus deceased patients, $FiO2_{last}$ and $RespRate_{last}$. Imputed values from each of the methods are shown in the boxplots on the right, and the rates of missingness in the deceased versus surviving patients for each variable are shown on the left. We found that a larger proportion of entries of $FiO2_{last}$ and a smaller proportion of entries of $RespRate_{last}$ were observed in the surviving patients than in the deceased patients. Of the methods we used to impute values of $FiO2_{last}$, we found that $NIMIWAE_{sup}$ and $NIMIWAE_{unsup}$ imputed values were generally smaller than

those from other methods. We also found that the supervised and unsupervised IMIWAE models yielded similar values to HIVAE, MIWAE, and MICE. These are all ignorably missing methods, and may not impute accurate values when the missingness is MNAR. MissForest and VAEAC imputed values of $FiO2$ that were significantly higher than the observed mean, and values of $RespRate$ that were significantly lower than the observed mean. Some studies have shown that abnormally high or low respiratory rates may be associated with higher mortality in critically-ill patients (Strauß et al. 2014; Ljunggren et al. 2016), suggesting that the missing values of $RespRate_{last}$, which were more prevalent in deceased patients, may have truly been further away from the normal range of 12–20. Additionally, Esteban (2002) found that higher levels of administered $FiO2$ were associated linearly with higher mortality. Thus, the missing values of $FiO2_{last}$, which were more prevalent in surviving patients, may have been lower than the observed values. This suggests that the mechanism of missingness of $FiO2_{last}$ and $RespRate_{last}$ may be non-ignorable, since NIMIWAE imputed more realistic values for these variables according to the mortality rates of patients with missingness in these variables.

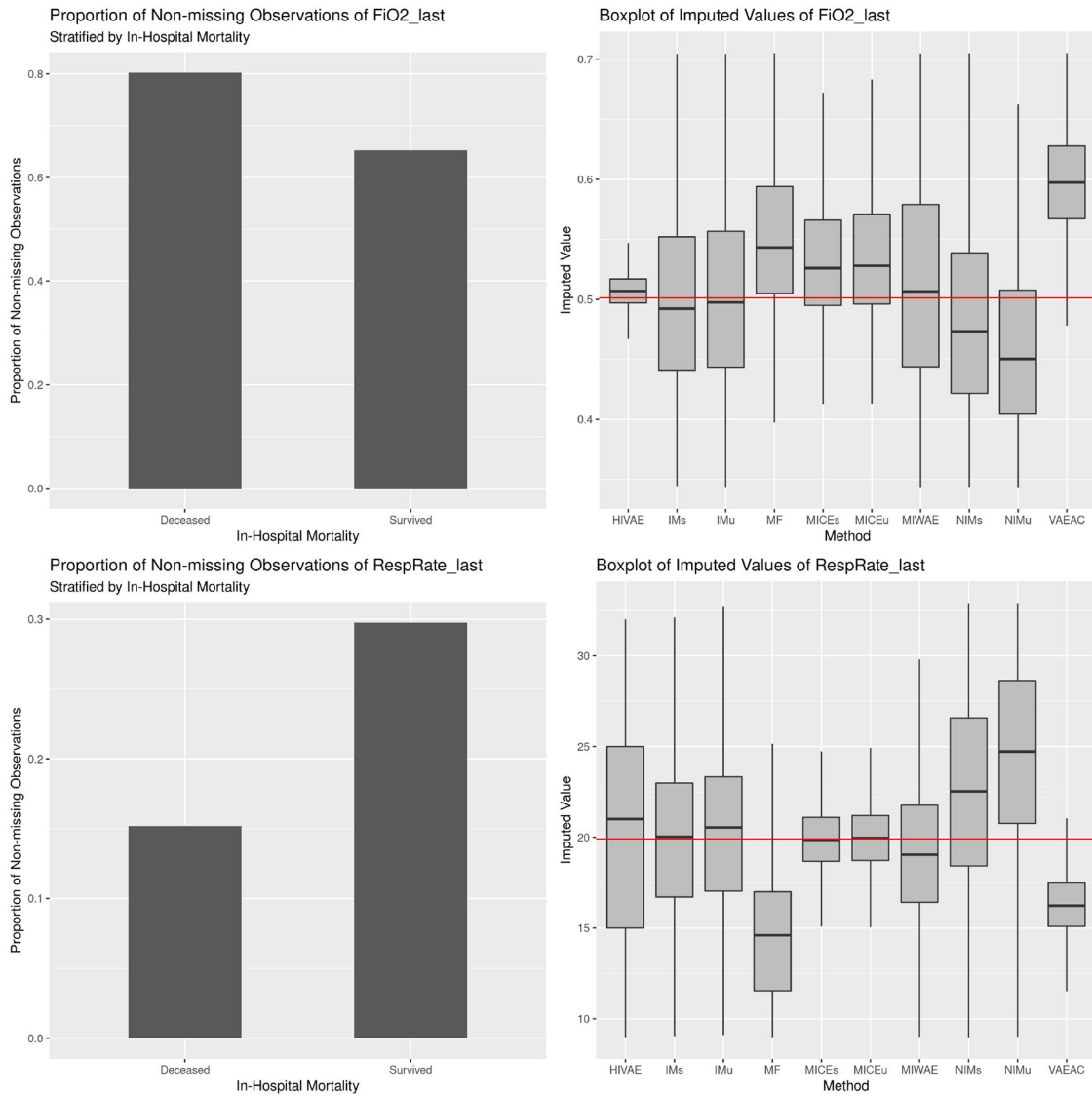


Figure 6. (left) Proportion of non-missing observations of *FiO2_last* (top) and *RespRate_last* (bottom) in surviving and deceased ICU patients, and (right) imputed values of non-missing entries by HIVAE, supervised and unsupervised versions of the ignorable NIMIWAE (IMs, IMu), MissForest (MF), supervised and unsupervised versions of MICE (MICEs, MICEu), MIWAE, supervised and unsupervised versions of NIMIWAE (NIMs, NIMu), and VAEAC. The mean of the observed values is given by the red horizontal line.

4. Discussion

In this article we introduce NIMIWAE, one of the first methods to handle up to MNAR patterns of missingness in the VAE/IWAE class of methods, to address complex patterns of missingness observed in the Physionet EHR data. Using statistical simulations, we show that NIMIWAE performs well in imputing missing features under MNAR, and has reasonable performance under the MCAR and MAR settings. Performance in imputing MCAR and MAR missingness can be further improved in NIMIWAE by using the ignorable version of this model (IMIWAE), where we omit the missingness network. We also found that the results of the analysis on the Physionet data are highly dependent on the choice of missingness model, which specifies the assumption of the missingness mechanism. However, NIMIWAE is able to impute missing values well in simulations regardless of the underlying missingness mechanism, flexibly modeling the mechanism using a deeply-learned neural network. The NIMIWAE-imputed dataset resulted in

more realistic imputed values with respect to what we may expect in the Physionet Challenge patients, since NIMIWAE takes into account possible non-ignorable missingness in the data. Additionally, the NIMIWAE architecture learns a lower-dimensional representation of the data, which can be used for tasks such as patient subgroup identification or visualization of data.

Learning algorithms that can be applied to EHR data like the Physionet Challenge dataset can be valuable tools that clinicians can use to aid decisions in hospital settings and understand patterns within these health records. For example, properly handling missingness in EHRs when imputing the missing entries can improve the performance of prediction algorithms that can assess risk of death or other outcomes of interest, like disease. Informative missingness is a common problem in analyzing EHR data, and accounting for such missingness can be helpful in obtaining accurate, unbiased estimates of the true missing values. We note that although we have used our NIMIWAE method

primarily to analyze the Physionet 2012 Challenge dataset, it can more generally be applied to settings where one wishes to train a VAE when missingness is present among input features.

Supplementary Materials

Supplementary Materials: Contains details of the NIMIWAE algorithm in Section A, additional simulation results and computational details in Sections B, and links and details of analyzed datasets in Section C. (pdf)

R-package for NIMIWAE: R-package NIMIWAE containing code to perform the diagnostic methods described in the article. The package can also be found at <https://www.github.com/DavidKLim/NIMIWAE>. (GNU zipped tar file)

Code for Reproducibility: Repository of code to reproduce all results, tables, and figures in the article. This repository can also be found at https://www.github.com/DavidKLim/NIMIWAE_Paper. (GNU zipped tar file)

Disclosure Statement

No potential conflict of interest was reported by the author(s).

Funding

The authors gratefully acknowledge NIH grants U01-CA274298, P50-CA257911, P50-CA058223, T32-CA106209, 1R01AA02687901A1, and 1OT2OD032581-02-321, and NSF grants IIS2133595 and DMS2324394 for funding this research.

References

- Beaulieu-Jones, B. K., and and, J. H. M. (2016), "Missing Data Imputation in the Electronic Health Record Using Deeply Learned Autoencoders," in *Biocomputing 2017*, World Scientific. [4]
- Bishop, C. M. (2006), *Pattern Recognition and Machine Learning*, New York: Springer-Verlag. [4]
- Burda, Y., Grosse, R., and Salakhutdinov, R. (2015), "Importance Weighted Autoencoders," arXiv e-prints p. arXiv:1509.00519. [3]
- Cremer, C., Morris, Q., and Duvenaud, D. (2017), "Reinterpreting Importance-Weighted Autoencoders," arXiv e-prints p. arXiv:1704.02916. [3]
- Diggle, P., and Kenward, M. G. (1994), "Informative Drop-Out in Longitudinal Data Analysis," *Applied Statistics*, 43, 49–73. [3,4]
- Doersch, C. (2016), "Tutorial on Variational Autoencoders," arXiv e-prints p. arXiv:1606.05908. [1]
- Du, H., Enders, C., Keller, B. T., Bradbury, T. N., and Karney, B. R. (2021), "A Bayesian Latent Variable Selection Model for Nonignorable Missingness," *Multivariate Behavioral Research*, 57, 478–512. [8]
- Esteban, A. (2002), "Characteristics and Outcomes in Adult Patients Receiving Mechanical Ventilation. A 28-day International Study," *JAMA*, 287, 345–355. [10]
- Gelman, A., Carlin, J. B., Stern, H. S., and Rubin, D. B. (1995), *Bayesian Data Analysis*, Boca Raton, FL: Chapman and Hall/CRC. [3]
- Gershman, S. J., and Goodman, N. D. (2014), "Amortized Inference in Probabilistic Reasoning," in *CogSci*. [2]
- Ibrahim, J. G. (2001), "Missing Responses in Generalised Linear Mixed Models When the Missing Data Mechanism is Nonignorable," *Biometrika*, 88, 551–564. [4]
- Ibrahim, J. G., Chen, M.-H., Lipsitz, S. R., and Herring, A. H. (2005), "Missing-Data Methods for Generalized Linear Models," *Journal of the American Statistical Association*, 100, 332–346. [9]
- Ibrahim, J. G., Lipsitz, S. R., and Chen, M.-H. (1999), "Missing Covariates in Generalized Linear Models When the Missing Data Mechanism is Nonignorable," *Journal of the Royal Statistical Society, Series B*, 61, 173–190. [9]
- Ibrahim, J. G., and Molenberghs, G. (2009), "Missing Data Methods in Longitudinal Studies: A Review," *TEST*, 18, 1–43. [4,9]
- Ivanov, O., Figurnov, M., and Vetrov, D. (2019), "Variational Autoencoder with Arbitrary Conditioning," in *International Conference on Learning Representations*. <https://openreview.net/forum?id=SyxtJh0qYm> [6]
- Kingma, D. P., and Ba, J. (2014), "Adam: A Method for Stochastic Optimization," arXiv e-prints p. arXiv:1412.6980. [4]
- Kingma, D. P., and Welling, M. (2019), "An Introduction to Variational Autoencoders," arXiv e-prints p. arXiv:1906.02691. [1,2]
- Larsson, E., Lindström, A.-C., Eriksson, M., and Oldner, A. (2019), "Impact of Gender on Post-Traumatic Intensive Care and Outcomes," *Scandinavian Journal of Trauma, Resuscitation and Emergency Medicine*, 27, 115. [10]
- LeCun, Y., Bengio, Y., and Hinton, G. (2015), "Deep Learning," *Nature*, 521, 436–444. [1]
- Li, Y., Akbar, S., and Oliva, J. (2020), "Acfow: Flow Models for Arbitrary Conditional Likelihoods," in *International Conference on Machine Learning*, PMLR, pp. 5831–5841. [1]
- Lim, K.-L., Jiang, X., and Yi, C. (2020), "Deep Clustering with Variational Autoencoder," *IEEE Signal Processing Letters*, 27, 231–235. [3]
- Little, R. J. A., and Rubin, D. B. (2002), *Statistical Analysis with Missing Data*, Hoboken, NJ: Wiley. [3,4]
- Ljunggren, M., Castrén, M., Nordberg, M., and Kurland, L. (2016), "The Association between Vital Signs and Mortality in a Retrospective Cohort Study of An Unselected Emergency Department Population," *Scandinavian Journal of Trauma, Resuscitation and Emergency Medicine*, 24, 21. [10]
- Lopez, R., Regier, J., Cole, M. B., Jordan, M. I., and Yosef, N. (2018), "Deep Generative Modeling for Single-Cell Transcriptomics," *Nature Methods*, 15, 1053–1058. [1]
- Mahmood, K., Eldeirawi, K., and Wahidi, M. M. (2012), "Association of Gender with Outcomes in Critically Ill Patients," *Critical Care*, 16, R92. [10]
- Mattei, P.-A., and Frellsen, J. (2019), "MIWAE: Deep Generative Modelling and Imputation of Incomplete Data Sets," in *Proceedings of the 36th International Conference on Machine Learning*, Volume 97 of *Proceedings of Machine Learning Research*, PMLR, Long Beach, California, USA, eds. K. Chaudhuri and R. Salakhutdinov, pp. 4413–4423. <http://proceedings.mlr.press/v97/mattei19a.html> [3,5,6]
- Murphy, J. (2016), "An Overview of Convolutional Neural Network Architectures for Deep Learning," *Microway Inc*, 1–22. [2]
- Nazabal, A., Olmos, P. M., Ghahramani, Z., and Valera, I. (2018), "Handling Incomplete Heterogeneous Data using VAEs," arXiv e-prints p. arXiv:1807.03653. [3,6]
- O'Shea, R. (2019), "Interpreting Missing Data Patterns in the ICU," arXiv e-prints p. arXiv:1912.08612. [4]
- Poggio, T., Banburski, A., and Liao, Q. (2020), "Theoretical Issues in Deep Networks," *Proceedings of the National Academy of Sciences*, 117, 30039–30045. [8]
- Prechelt, L. (1998), "Early Stopping-but When?" in *Neural Networks: Tricks of the Trade*, pp. 55–69, Berlin: Springer. [2]
- Ross, M. K., Wei, W., and Ohno-Machado, L. (2014), "'big data" and the Electronic Health Record," *Yearbook of Medical Informatics*, 23, 97–104. [2]
- Rubin, D. B. (1976), "Inference and Missing Data," *Biometrika*, 63, 581–592. [3]
- (2004), *Multiple Imputation for Nonresponse in Surveys* (Vol. 81), New York: Wiley. [5]
- Saxe, A. M., McClelland, J. L., and Ganguli, S. (2014), "Exact Solutions to the Nonlinear Dynamics of Learning in Deep Linear Neural Network," in *International Conference on Learning Representations*. [2]
- Schell-Chaple, H. M., Puntillo, K. A., Matthay, M. A., and and, Liu, K. D. (2015), "Body Temperature and Mortality in Patients with Acute Respiratory Distress Syndrome," *American Journal of Critical Care*, 24, 15–23. [10]
- Shickel, B., Tighe, P. J., Bihorac, A., and Rashidi, P. (2018), "Deep EHR: A Survey of Recent Advances in Deep Learning Techniques for Electronic Health Record (EHR) Analysis," *IEEE Journal of Biomedical and Health Informatics*, 22, 1589–1604. [1]

- Smith, A. F. M., and Gelfand, A. E. (1992), “Bayesian Statistics Without Tears: A Sampling–Resampling Perspective,” *The American Statistician*, 46, 84–88. [5]
- Stekhoven, D. J., and Buhlmann, P. (2011), “MissForest–Non-Parametric Missing Value Imputation for Mixed-Type Data,” *Bioinformatics*, 28, 112–118. [6]
- Strauß, R., Ewig, S., Richter, K., König, T., Heller, G., and Bauer, T. T. (2014), “The Prognostic Significance of Respiratory Rate in Patients with Pneumonia,” *Deutsches Ärzteblatt international*, 111, 503–508. [10]
- Strauss, R., and Oliva, J. B. (2021), “Arbitrary Conditional Distributions with Energy,” in *Advances in Neural Information Processing Systems* (Vol. 34). [1]
- Stubbendick, A. L., and Ibrahim, J. G. (2003), “Maximum Likelihood Methods for Nonignorable Missing Responses and Covariates in Random Effects Models,” *Biometrics*, 59, 1140–1150. [4]
- Tongyoo, S., Viarasilpa, T., and Permpikul, C. (2018), “Serum Potassium Levels and Outcomes in Critically Ill Patients in the Medical Intensive Care Unit,” *Journal of International Medical Research*, 46, 1254–1262. [10]
- Tschannen, M., Bachem, O., and Lucic, M. (2018), “Recent Advances in Autoencoder-Based Representation Learning,” arXiv e-prints p. arXiv:1812.05069. [1]
- Van Buuren, S., and Groothuis-Oudshoorn, K. (2011), “mice: Multivariate Imputation by Chained Equations in r,” *Journal of Statistical Software*, 45, 1–67. [6]
- Vinzamuri, B., and Reddy, C. K. (2013), “Cox Regression with Correlation based Regularization for Electronic Health Records,” in *2013 IEEE 13th International Conference on Data Mining*, pp. 757–766. [2]
- Wang, Y., Yao, H., and Zhao, S. (2016), “Auto-Encoder based Dimensionality Reduction,” *Neurocomputing*, 184, 232–242. [1]
- Wells, B. J., Nowacki, A. S., Chagin, K., and Kattan, M. W. (2013), “Strategies for Handling Missing Data in Electronic Health Record Derived Data,” *eGEMs (Generating Evidence & Methods to improve patient outcomes)*, 1, 1035. [2]