

# DDGANI: Missing Value Imputation using Denoising Diffusion GAN with Modular Plugins (Technical Report)

Shuang Hao, Xiang Huang

School of Computer Science and Technology, Beijing Jiaotong University, China  
Beijing Key Laboratory of Traffic Data Analysis and Mining, China  
{haoshuang,22120382}@bjtu.edu.cn

## ABSTRACT

Dealing with missing values in relational tables is a pervasive challenging preparation phase before proceeding to data analysis and mining. Directly removing missing data often causes the loss of information and thus adversely impacts further data processing. Missing values can also be estimated using extensively studied imputation techniques. However, they seldom consider improving the utility of imputed data for the downstream target task. Recently widely used learning-based approaches rarely try to maintain the inherent data dependencies. In this paper, we introduce DDGANI, a novel missing value imputation framework based on denoising diffusion GAN. Central to our approach is the utilization of denoising diffusion GAN for effective data imputation. Additionally, we incorporate a self-attention mechanism, enabling the model to focus on more relevant information from similar tuples. Our framework also integrates two plugins that can be chosen based on data characteristics and downstream tasks, which ensure the utility of imputed data and the preservation of inherent data dependencies. Extensive experiments on multiple datasets demonstrate the superiority of DDGANI compared to existing state-of-the-art imputation techniques and verify the effectiveness of two modular plugins.

### PVLDB Reference Format:

Shuang Hao, Xiang Huang. DDGANI: Missing Value Imputation using Denoising Diffusion GAN with Modular Plugins (Technical Report). PVLDB, 14(1): XXX-XXX, 2024.  
doi:XX.XX/XXX.XX

### PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://github.com/xianghuang2/DDGANI>.

## 1 INTRODUCTION

Data is undoubtedly an essential element of decision-making, so high-quality data is a necessity. However, in real-world scenarios, data often goes missing, which may arise due to human error during data collection or equipment malfunctions. While some factors leading to missing values are perceived, they are unavoidable, such as table joining for relational data augmentation [10]. Therefore,

it is essential to address the missing data before diving into data analysis and mining.

**Existing Approaches and Limitations.** Manually imputing these missing values is undoubtedly a time-consuming and labor-intensive task. A straightforward approach might be to delete these missing data. However, neglecting such missing values can distort the analysis or lead to inaccuracies, ultimately affecting the downstream decision-making process [2]. Thus, the most common strategy to address missing values is reconstructing a complete dataset through data imputation techniques. [Most of the earlier studies focus on the basic statistical imputation methods, utilizing statistics such as the mean, median, and mode to impute the missing values. However, most of these methods only work well for numerical data. The traditional machine learning \(ML\) imputation solutions train an ML model for missing value prediction. The accuracy of imputation is limited by the predictive power of the ML model. In addition, some of them have strict distributional assumptions such as PCAI \[22\], and some are less effective for categorical data such as KNN \[1\].](#) Heuristic-based methods utilize data imputation rules to deal with missing values, such as the relaxed functional dependencies [6] and differential dependencies [39]. However, they mainly focus on repairing categorical values (e.g., City, State), and the design of high-quality imputation rules requires extensive domain knowledge [28]. Moreover, they sometimes fail to provide a deterministic imputation solution for attributes on the left-hand side of these data dependencies. For example, for the missing value  $t_5[\text{AreaCode}]$  in Table 1, these methods may give two possible imputations as both "501" and "870" satisfy the given imputation rules.

Recently, deep learning based approaches, mainly consisting of methods based on deep generative models and pre-trained language models, have provided an effective way to tackle the above obstacles and have obtained impressive results on the missing value imputation problem. Generative adversarial networks (GAN) and variational autoencoder (VAE) are two common deep generative models adopted in previous work [17, 30, 33, 48]. However, they have difficulty directly generating samples from a complex distribution in one shot. Furthermore, they often neglect to check whether the imputed table violates the original data dependencies. For example, given the similar distribution between  $t_4$  and  $t_5$  in Table 1, most of these methods would impute  $t_4[\text{State}]$  as "AR" which results in a violation of  $\phi_1$ . Actually, "Neward" in  $t_4$  and  $t_5$  are different cities in different states despite having the same name. With the prominence of pre-trained language models (PLMs) on multiple data processing tasks, there are PLMs-based data imputation methods such as TURL [15] and RPT [42]. TURL primarily focuses on learning data representations and aligning these representations with knowledge bases to retrieve missing data. RPT adopts a Transformer-based

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing [info@vldb.org](mailto:info@vldb.org). Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 14, No. 1 ISSN 2150-8097.  
doi:XX.XX/XXX.XX

**Table 1: Example of a Relational Table with Missing Values**  
(A gray cell means it is missing and an **turquoise** colored word indicates the true value.)  
( $\varphi_1 : \text{Areacode} \rightarrow \text{State}$ ,  $\varphi_2 : \text{Areacode}, \text{MaritalStatus} \rightarrow \text{SingleExemp}$ )

Tuple ID	Fname	AreaCode	City	State	MaritalStatus	SingleExemp	Salary	Rate	HasChild (label)
$t_1$	Muneo	201	Hampton	NJ	M	0	10000	1.40	Yes
$t_2$	Hyujin	201	Hampton	NJ	M	0	<b>[90000]</b>	2.50	Yes
$t_3$	Jarno	201	Hampton	NJ	S	<b>[1000]</b>	80000	2.34	No
$t_4$	Coke	201	Newark	<b>[NJ]</b>	M	0	5000	<b>[1.40]</b>	No
$t_5$	Ada	<b>[501]</b>	Newark	AR	M	0	5000	1.38	Yes
$t_6$	Liesbeth	501	Newark	AR	S	20	<b>[5000]</b>	1.38	No
$t_7$	Refik	501	<b>[Newark]</b>	AR	S	20	10000	2.58	Yes
$t_8$	Sreerupa	870	Sage	AR	<b>[S]</b>	20	<b>[5000]</b>	1.38	No

neural translation architecture, pre-trained by corrupting the input tuple and then learning to reconstruct the original tuple to support the missing value imputation task. However, PLMs-based imputation methods do not handle numerical data well (e.g., Rate). Furthermore, they often fail to ensure the applicability of the imputed data for subsequent tasks, which is a significant challenge at present [32].

**Our Proposal.** To address the aforementioned limitations, we propose a novel missing value imputation framework DDGANI as depicted in Figure 1. It contains a data imputation module, which is the mainstay of the framework, and two optional plug-ins: a data utility plug-in and a data dependency plug-in. We employ the denoising diffusion GAN [47] in the data imputation module to reconstruct the complete dataset, which takes the relational table with missing values as input and outputs the imputed table. This choice is motivated by its ability to model complex data distributions through several conditional denoising diffusion steps, allowing for the proficient generation of various types of data and achieving high data quality while also obtaining good mode coverage. Additionally, we incorporate the self-attention mechanism, which enables the generative model to capture more relevant contextual information when dealing with missing data.

We introduce a data utility plug-in to ensure the utility of the imputed data for the subsequent classification task via a meta-learning paradigm. We fully use the downstream classifier and validation set, with the goal that the classifier trained with imputed data could minimize the loss on the validation set. It is achieved by back-propagating the validation loss to the data generation module to guarantee the utility of the imputed table. If the downstream classifier or validation set is unavailable, the data utility plug-in builds a two-layer fully connected neural network acting as the downstream classifier and selects some tuples from the incomplete table where all attributes are observed to serve as a validation set. This can additionally contribute to keeping the same data distribution between the imputed and observed data.

We also design a data dependency plug-in to allow the imputed data to adhere to the original data dependencies among the categorical attributes. In this paper, we describe how it works using function dependencies (FDs) as an example. This plug-in could also support other constraints such as conditional functional dependencies (CFDs). This module initially mines potential FDs from the table with missing values. Each FD  $\varphi : Z \rightarrow A$  is then modeled by

a classifier, which takes the values of attributes  $Z$  in the left-hand side of the FD as input and predicts the attribute value of  $A$  (These models also could be pre-trained using the master data or knowledge bases). We minimize the difference between the predicted values given by the dependency models and the imputed values given by the data imputation module to keep the imputed data from violating the existing data constraints. For example, given  $\varphi_1 : \text{AreaCode} \rightarrow \text{State}$ , the imputation of  $t_4[\text{State}]$  in Table 1 with "AR" would result in a penalty in the loss of the generative model. Consequently, DDGANI would lean towards imputing it with "NJ". However, the automatically generated FDs from an incomplete dataset may not be fully trusted. As the imputation becomes more accurate, we must go back to determine whether these FDs are valid based on the confidence of constraints and violated tuples. Intuitively, if the violated tuples are more trusted, we refine the FDs and their corresponding dependency models.

Each component in our framework, though appearing to have its distinct goal, works together to minimize imputation error. The attention module and two plug-ins further improve the ability of the base model to capture the correlations between attributes and between attributes and the label in different ways. Meanwhile, this modular design can make the whole framework more flexible for different data characteristics and downstream tasks.

**Contributions.** Our contributions are summarized as follows:

- We propose a novel missing value imputation framework with a self-attention-enhanced generative model for tabular data imputation and two plug-ins that can be chosen to improve the performance based on data characteristics and downstream tasks.
- We adopt the denoising diffusion GAN for missing value imputation. To the best of our knowledge, we are the first effort to model the complex data distribution implied in missing data through denoising diffusion steps to obtain high-quality imputation results (Section 4).
- We design a data utility plug-in to boost the performance of our imputation framework via a meta-learning paradigm, which guarantees the applicability of our imputed data for training downstream models for a target task (Section 5).
- We propose a data dependency plug-in to ensure the preservation of inherent data dependencies in the imputed relational table, which iteratively imputes the missing values and refines the dependency models based on the confidence of mined FDs and imputed tuples. (Section 6).

- We conducted extensive experiments on multiple datasets with both categorical and numerical attributes. The experimental results demonstrate that our approach outperforms a variety of state-of-the-art missing value imputation methods. (Section 8).

## 2 RELATED WORK

We categorize related work as follows. Note that this work extends our short paper [18] by including 1) a new choice on the missing value imputation model, *i.e.*, we employ the denoising diffusion GAN instead of VAE, which stems from the fact that the multi-step diffusion process tends to generate data of higher quality; 2) a more comprehensive attention module, which takes into account the correlations between different attributes instead of just at the tuple level; 3) a data dependency plug-in to ensure the preservation of inherent data dependencies in the relational table; 4) more comprehensive empirical evaluations, mainly including more datasets and baselines, the effectiveness of newly added plug-ins, the impact of hyper-parameters, [the capability to handle rare values](#), and [the runtime analysis](#).

**Statistical Imputation Approaches.** Statistical-based methods typically calculate statistics such as mean, median, and mode of observed data to estimate missing values. Several statistical methods can also be combined to impute missing data [7, 45]. However, most of these methods only work well for imputing numerical data, which is unsuitable for some real-world scenarios and far from enough to be high-quality.

**Heuristic Based Imputation Approaches.** Heuristic methods have been proposed to effectively estimate the missing values without violating the integrity of the original data. For example, Derand [39], RENUVER [6], Graf [34] and Holistic [13] respectively utilize differential dependencies [38], relaxed functional dependencies [9], [conditional functional dependencies](#) [5] and denial constraints [12] to ensure the imputed data satisfy the given constraints. However, they cannot impute the attributes not appearing in the dependencies. There exist some methods directly mining constraints from missing data for the subsequent imputation [3, 26] and estimating the confidence (*i.e.*, genuineness score) of discovered constraints based on the tuple frequency while ignoring each tuple’s confidence. Graf [34] [iteratively repairs the dirty data and inaccurate CFDs based on their confidence, but it requires lots of complete tuples for constraint learning](#).

**Machine Learning Based Imputation Approaches.** ML-based imputation methods generally train ML models to predict the missing values. The k-nearest neighbors (KNN) based method [44] uses a weighted average of values from the k nearest neighbors to estimate the missing value. MissForest [41] predict the missing values utilizing observed data through the random forest. Other ML-based methods exist, such as MICE [36] and PCAI [29]. Although these ML-based imputation methods outperform statistical ones [31], they may have strict distributional assumptions or struggle with categorical data.

**Deep Learning Based Imputation Approaches.** DL-based methods generally outperform the aforementioned approaches as they can capture the correlations between data through multiple layers of non-linear computations, enabling effective imputation of missing data. [DataWig \[4\] builds a DL model for each attribute to predict](#)

[the missing values](#). GAIN [48] and RandomGAN [17] adapt GAN architecture to reconstruct the whole tabular data. HI-VAE [33] and VAEI [30] extend VAE for missing value imputation. GINN [40] and EEG-GAE [43] utilize graph neural networks. [Baran \[27\] employs pre-trained DL models. It iteratively asks the user to label the data and leverages the provided user corrections to update the pre-trained models incrementally.](#) Pre-trained language models have recently been used for data imputation tasks, including TURL [15] and RPT [42]. [Some studies such as AimNet \[46\] impute the missing values based on DL models while incorporating data dependencies.](#) Our work falls into the scope of DL-based methods, and the experimental results show that our method outperforms the above baselines.

## 3 PROBLEM FORMULATION

**Data Model.** Let  $\mathbf{X}$  denote a relational table with  $n$  tuples  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  and  $k$  attributes  $U = \{A_1, \dots, A_k\}$ . Each attribute is either numerical or categorical. We use  $x_{ij}$  to denote the value of attribute  $A_j$  in tuple  $\mathbf{x}_i$ . Each tuple  $\mathbf{x}_i$  has a label  $y_i$  for classification.

Note that for ease of illustration in the examples, we sometimes refer to each tuple by its ID  $t_i$  in Table 1 and use  $t_i[A_j]$  to denote the value of attribute  $A_j$  in tuple  $t_i$ .

**Missing Data.**  $\mathbf{X}$  contains *missing data*: each tuple  $\mathbf{x}$  may have some attribute values that are not observed. To formalize these missing values, we follow the existing work [25] to introduce a *mask matrix*  $\mathbf{M} = \{\mathbf{m}_1, \dots, \mathbf{m}_n\}$  to indicate which value in  $\mathbf{X}$  is observed. Each  $\mathbf{m}_i = \{m_{i1}, \dots, m_{ik}\}$  represents the missing status of tuple  $\mathbf{x}_i$ :  $x_{ij}$  is observed if  $m_{ij} = 1$ , otherwise  $x_{ij}$  is missing.

*Example 3.1.* Given a relational table in Table 1 where a gray cell means it is missing. For tuple  $t_4$  in Table 1, it has missing values in  $t_4[\text{State}]$  and  $t_4[\text{Rate}]$  while other attribute values are observed. Therefore, its corresponding mask is  $\mathbf{m}_4 = \{1, 1, 1, 0, 1, 1, 1, 0\}$ .

**Data Dependency.** A dependency is a constraint that applies to or defines the relationship between attributes, where knowing the value of one attribute (or a set of attributes) is enough to infer the value or range of another attribute. Let’s take the functional dependency (FD) as an example. We say  $\mathbf{X}$  satisfies the FD  $\varphi : Z \rightarrow A$  ( $Z, A$  are two disjoint subsets of the attributes, and the right-hand side  $A$  of  $\varphi$  has only one attribute) if, for any tuple, the attribute values in  $Z$  uniquely determine the attribute value in  $A$ .

*Example 3.2.* Consider the following two functional dependencies in Table 1:

- $\varphi_1 : \text{AreaCode} \rightarrow \text{State}$
- $\varphi_2 : \text{State}, \text{MaritalStatus} \rightarrow \text{SingleExemp}$

where  $\varphi_1$  means that AreaCode uniquely determines State and there should not exist two tuples that have the same value in AreaCode but different values in State.

**Missing Data Imputation.** The missing data imputation problem is to replace missing data in  $\mathbf{X}$  with an estimated value based on other observed values. If  $\tilde{\mathbf{X}} = \{\tilde{x}_1, \dots, \tilde{x}_n\}$  denotes the imputed table and  $\mathbf{X}^\checkmark = \{x_1^\checkmark, \dots, x_n^\checkmark\}$  denotes the complete ground-truth table, which is not available in practical situations, our goal is to minimize the difference between  $\tilde{\mathbf{X}}$  and  $\mathbf{X}^\checkmark$ .

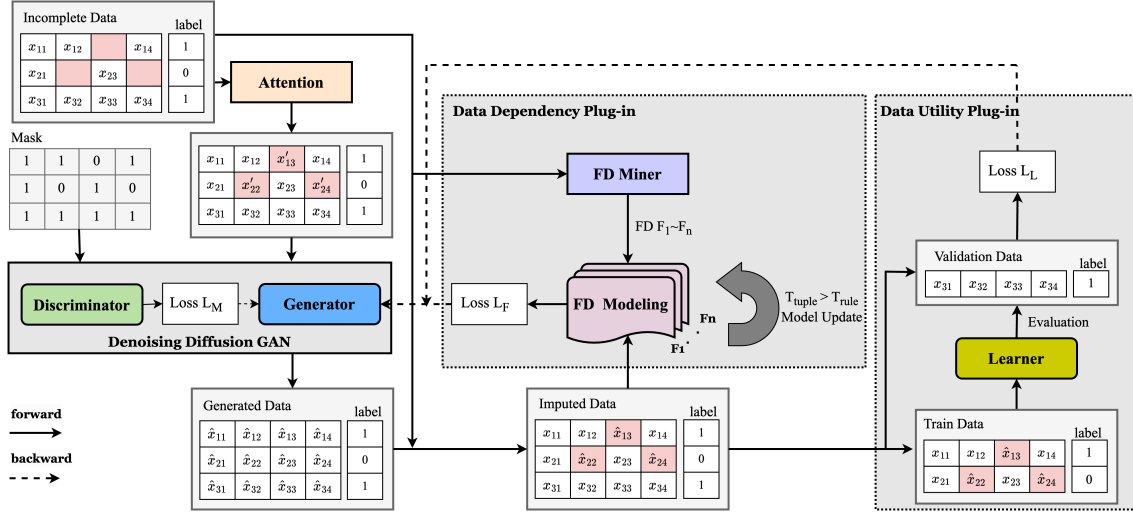


Figure 1: The Architecture of DDGANI

## 4 DENOISING DIFFUSION GAN BASED IMPUTATION NETWORK

In this section, we will start with an introduction to denoising diffusion GAN and then discuss how to utilize it to impute the missing values in a relational table.

### 4.1 Denoising Diffusion GAN

**Why Diffusion Models?** Diffusion models have been shown to outperform generative models in a one-shot manner in many cases despite the complexity of their training as a tradeoff. Kotelnikov et al. [24] have demonstrated that the diffusion models could better learn the distribution characteristics of tabular data. This is because the one-shot generative processes such as GANs and VAEs face the difficulty of directly generating samples from a complex distribution in one shot, which might not capture all data details. In contrast, diffusion models first establish a coarse distribution and then incrementally refine it, ensuring high-quality data generation while obtaining good mode coverage.

**Denoising Diffusion Probabilistic Model.** The denoising diffusion probabilistic model (DDPM) defines a forward diffusion process that adds varying degrees of Gaussian noise to the original data  $X_0 \sim q(X_0)$  over  $T$  steps where  $q(X_0)$  is the original data distribution. The forward process at step  $t$  can be represented as  $q(X_t|X_{t-1})$  where  $X_t$  denotes the data obtained by adding noise to  $X_{t-1}$ , i.e., to  $X_0$  for  $t$  steps. The reverse diffusion process is to guess the Gaussian noise added in the forward process and remove it iteratively (within  $T$  steps), using a neural network model  $\Phi_\theta$ . The reverse process can be represented as  $p_\theta(X_{t-1}|X_t)$  where  $\theta$  represents the parameters of the model. Thus, the training objective of the diffusion-based generative model amounts to maximizing the log-likelihood of the data generated at the end of the reverse process belonging to the original data distribution, which is formulated as maximizing the evidence lower bound (ELBO) [23].

However, due to the typically unknown nature of  $q(X_{t-1}|X_t)$  in the diffusion model, the denoising step size  $T$  is assumed to be

on the order of thousands of steps, and the diffusion rate  $\beta_t$  in each step should be small, which makes its training time-consuming.

**Denoising Diffusion GAN.** The large denoising diffusion steps  $T$  leads to significant computational costs during model training. To address this issue, Xiao et al. [47] employed conditional GANs to approximate  $q(X_{t-1}|X_t)$ . Specifically, the generator  $G$  takes a random noise as input and outputs  $\hat{X}_0$  with  $X_t$  and step  $t$  as conditions. Then,  $\hat{X}_{t-1}$  is obtained through a posterior sampling process  $q(X_{t-1}|X_t, X_0)$  from  $\hat{X}_0$ . The discriminator  $D$  takes  $(\hat{X}_{t-1}, X_t, t)$  or  $(X_{t-1}, X_t, t)$  as inputs and determines the probability that the input data is  $X_{t-1}$  (real) rather than  $\hat{X}_{t-1}$  (fake) given  $X_t$ . After achieving stability in training, the generated  $\hat{X}_{t-1}$  from the generator will closely resemble  $X_{t-1}$ , resulting in a consistent distribution between  $p_\theta(X_{t-1}|X_t)$  and  $q(X_{t-1}|X_t)$ . Therefore, denoising diffusion GANs could adopt a relatively large diffusion rate  $\beta_t$  and a small denoising diffusion step  $T$  ( $T \leq 8$ ).

### 4.2 Denoising Diffusion GAN For Data Imputation

Now it's time to introduce how to impute the missing values with denoising diffusion GAN.

**Attention Mechanisms.** Instead of feeding the data  $X$  with missing values directly into the generative model, we pre-process  $X$  with the self-attention mechanism, which encourages the subsequent generation model to capture the internal dependencies within the data and pay more attention to the observed data associated with the missing values.

Here, we build an attention map  $\mathcal{A}^j$  for each attribute  $A_j$ , of size  $n \times n$  where  $n$  is the number of tuples. Then, considering a missing value  $x_{ij}$  (the value of attribute  $A_j$  in tuple  $\mathbf{x}_i$  is missing), we utilize the attention mechanisms to implement the action of selectively concentrating on fewer tuples for the value on attribute  $A_j$  while ignoring the others when imputing  $x_{ij}$ . Intuitively, tuples that are more similar to  $\mathbf{x}_i$  on attribute  $A_j$  deserve more attention, thus we



have

$$x'_{ij} = \sum_h s(\mathbf{x}_i, \mathbf{x}_h)_j \times x_{hj}. \quad (1)$$

Here,  $s(\mathbf{x}_i, \mathbf{x}_h)_j$  is the value in row  $i$  and column  $h$  of  $\mathcal{A}^j$ , representing the similarity between  $x_{ij}$  and  $x_{hj}$ . Notice that  $x_{ij}$  is missing, so the question here is how to calculate the similarity to a missing value. Actually, we can estimate the similarity of two tuples on attribute  $A_j$  by comparing their similarity on other attributes, and the more correlated an attribute is to  $A_j$ , the more effective it is in estimating the similarity on attribute  $A_j$ .

The next question is how to calculate the correlation between two attributes. We build a symmetric correlation matrix  $\mathbf{C}$  while each element  $c_{ij}$  represents the correlation between attribute  $A_i$  and  $A_j$ . For two attributes  $A_i$  and  $A_j$ , we set  $c_{ij} = 0$  if  $i = j$ . Otherwise, we pick all the tuples with no missing values on both  $A_i$  and  $A_j$  and quantify the strength of correlation between them (denoted by  $c_{ij}$  and also  $c_{ji}$ ) by the absolute value of *Cramér's V coefficient* [14].

Then, given two tuples  $\mathbf{x}_i$  and  $\mathbf{x}_h$ , the similarity between  $x_{ij}$  and  $x_{hj}$  is measured with the weighted sum of the similarity of these two tuples on the other attributes and the weights are determined by the strength of correlation with  $A_j$ . That is,

$$s(\mathbf{x}_i, \mathbf{x}_h)_j = \sum_z c_{zj} \times \text{sim}(x_{iz}, x_{hz}). \quad (2)$$

The similarity  $\text{sim}(x_{iz}, x_{hz})$  is defined as follows. If  $A_z$  is a categorical attribute,  $\text{sim}(x_{iz}, x_{hz}) = 1$  when  $x_{iz}$  and  $x_{hz}$  have the same value, otherwise  $\text{sim}(x_{iz}, x_{hz}) = 0$ . If  $A_z$  is a numerical attribute,  $\text{sim}(x_{iz}, x_{hz}) = 1 - \frac{|x_{iz} - x_{hz}|}{\max_z - \min_z}$  where  $\max_z$  and  $\min_z$  are the maximum and minimum values of attribute  $A_z$  respectively. Then, each row  $\mathcal{A}^j_i$  of the attention map  $\mathcal{A}^j$  will be

$$\mathcal{A}^j_i = \text{Softmax}(s(\mathbf{x}_i, \mathbf{x}_1)_j, \dots, s(\mathbf{x}_i, \mathbf{x}_{n-1})_j, s(\mathbf{x}_i, \mathbf{x}_n)_j). \quad (3)$$

*Example 4.1.* Suppose we attempt to estimate the missing value  $t_2[\text{Salary}]$  in Table 1. We first compute the correlations between Salary (i.e.,  $A_7$ ) and other attributes as  $\mathbf{c}_7 = (0, 0.03, 0, 0.05, 0.03, 0.04, 0, 0.62)$ . Then, the similarity between  $t_2$  and  $t_3$  on Salary is computed as  $s(\mathbf{x}_2, \mathbf{x}_3)_7 = 0 \times 0 + 0.03 \times 1 + 0 \times 1 + 0.05 \times 1 + 0.03 \times 0 + 0.62 \times (1 - |(2.50 - 2.34)/(2.58 - 1.38)|) \approx 0.62$ . Similarly, we can obtain  $s(\mathbf{x}_2, \mathbf{x}_1)_7 \approx 0.21$ . Thus,  $t_3$  has a higher attention score than  $t_1$  when imputing  $t_2[\text{Salary}]$ , since it has more similar values on the attributes (e.g., Rate) which are more relevant to Salary.

*Remark.* (1) Equation 1 only applies to the case where  $x_{ij}$  is missing but  $x_{hj}$  is not. When  $x_{hj}$  is also missing,  $s(\mathbf{x}_i, \mathbf{x}_h)_j = 0$ . When  $x_{ij}$  is not missing, we will utilize  $\mathbf{M}$  to mask out  $x'_{ij}$  as Equation 4 (there is no need to utilize the attention mechanisms for an observed value). (2) When traversing the other attributes  $A_z$  of tuples  $\mathbf{x}_i$  and  $\mathbf{x}_h$  in Equation 2, if  $x_{iz}$  or  $x_{hz}$  is missing, we skip their similarity computation on attribute  $A_z$  as the above example shows. (3) We utilize a threshold  $\kappa$  to set the values less than  $\kappa$  in  $\mathcal{A}^j_i$  to 0, which could help to prevent the decision bias caused by data imbalance.

Let  $\mathbf{X}_j$  be the projection of  $\mathbf{X}$  on attribute  $A_j$ , and  $\mathbf{X}^{\mathcal{A}}$  is the concatenation of the product  $\mathcal{A}^j \times \mathbf{X}_j$  across all attributes. We can then reconstruct  $\mathbf{X}$  by

$$\mathbf{X}' = (\mathbf{1} - \mathbf{M}) \odot \mathbf{X}^{\mathcal{A}} + \mathbf{M} \odot \mathbf{X} \quad (4)$$

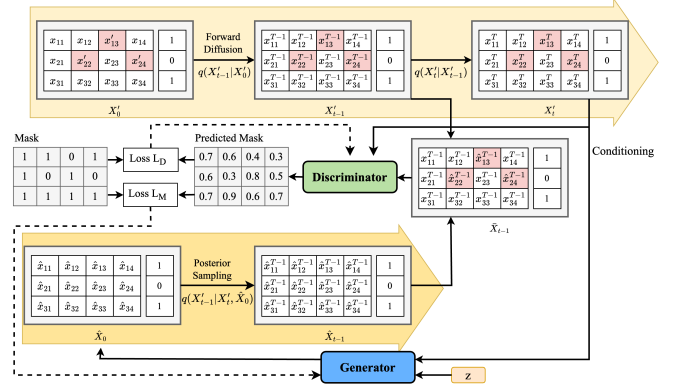


Figure 2: The Training Process of Data Imputation

where  $\odot$  denotes element-wise multiplication and  $\mathbf{1}$  is all-ones matrix. Note that we only use attention mechanisms during the data pre-processing phase and not after each iteration. This is because the cost of calculating similarities per iteration increases significantly with large datasets but does not significantly enhance the final results.

**Generator.** After obtaining  $\mathbf{X}'$ , we further impute the missing values with the use of the denoising diffusion GAN, which consists of a forward diffusion process and a posterior sampling process, as presented in Figure 2. To ease the following description, we denote  $\mathbf{X}'$  by  $\mathbf{X}'_0$ , which is the output of the attention mechanism and the input of the following diffusion model.

In the forward diffusion process, we iteratively add Gaussian noise to  $\mathbf{X}'_0$  over  $T$  steps. In each step, we have:

$$q(\mathbf{X}'_t | \mathbf{X}'_{t-1}) = \mathcal{N}(\mathbf{X}'_t; \sqrt{1 - \beta_t} \mathbf{X}'_{t-1}, \beta_t \mathbf{I}). \quad (5)$$

Here,  $\mathcal{N}$  represents the Gaussian distribution;  $t$  denotes any step between 1 and  $T$ ;  $\beta_t$  is the diffusion rate, representing the degree of noise added at step  $t$ ;  $\mathbf{X}'_t$  denotes the data achieved by introducing noise into  $\mathbf{X}'_0$  over  $t$  steps;  $\mathbf{I}$  denotes the identity matrix of the same dimension as  $\mathbf{X}'_t$ .

In the posterior sampling process, we employ a generator  $G$  to perform denoising, i.e., approximating  $q(\mathbf{X}'_{t-1} | \mathbf{X}'_t)$ . However, instead of taking  $\mathbf{X}'_t$  as input and outputting  $\mathbf{X}'_{t-1}$ , the generator will try to estimate  $\mathbf{X}'_0$  and then adds noise to obtain  $\mathbf{X}'_{t-1}$ . In this way,  $G$  has the singular goal of generating  $\mathbf{X}'_0$ , rather than a series of  $\mathbf{X}'_{t-1}$  for different step  $t$ , which is beneficial to the effective training of  $G$ .

Concretely, let  $\hat{\mathbf{X}}_0$  denote the actual output of  $G$  to be as similar as  $\mathbf{X}'_0$  in the observed data. We inject the noise into  $\hat{\mathbf{X}}_0$  again to get  $\hat{\mathbf{X}}_{t-1}$ :

$$q(\hat{\mathbf{X}}_{t-1} | \mathbf{X}'_t, \hat{\mathbf{X}}_0) = \mathcal{N}(\hat{\mathbf{X}}_{t-1}; \tilde{\mu}_t(\mathbf{X}'_t, \hat{\mathbf{X}}_0), \tilde{\beta}_t \mathbf{I}) \quad (6)$$

where  $\alpha_t = 1 - \beta_t$ ,  $\tilde{\alpha}_t = \prod_{s=1}^t \alpha_s$ ,  $\tilde{\mu}_t(\mathbf{X}'_t, \hat{\mathbf{X}}_0) = \frac{\sqrt{\tilde{\alpha}_{t-1}} \beta_t}{1 - \tilde{\alpha}_t} \hat{\mathbf{X}}_0 + \frac{\sqrt{\alpha_t}(1 - \tilde{\alpha}_{t-1})}{1 - \tilde{\alpha}_t} \mathbf{X}'_t$  and  $\tilde{\beta}_t = \frac{1 - \tilde{\alpha}_{t-1}}{1 - \tilde{\alpha}_t} \beta_t$ . And we define the imputed table  $\tilde{\mathbf{X}}_{t-1}$  in step  $t-1$  to be the combination of the partial observations from  $\mathbf{X}'_{t-1}$  and imputed data from  $\hat{\mathbf{X}}_{t-1}$ :

$$\tilde{\mathbf{X}}_{t-1} = (\mathbf{1} - \mathbf{M}) \odot \mathbf{X}'_{t-1} + \mathbf{M} \odot \hat{\mathbf{X}}_{t-1} \quad (7)$$

Subsequently, we use  $\tilde{\mathbf{X}}_{t-1}$  as the input to  $G$  and obtain a new  $\hat{\mathbf{X}}_0$ . Then we add noise to  $\hat{\mathbf{X}}_0$  to generate  $\hat{\mathbf{X}}_{t-2}$  by Equation 6, and derive

$\tilde{\mathbf{X}}_{t-2}$  using Equation 7. This process continues until we generate the final imputed table  $\tilde{\mathbf{X}}_0$  (*i.e.*,  $\tilde{\mathbf{X}}$ ).

**Discriminator.** To make the fake (imputed) values in  $\tilde{\mathbf{X}}$  to be as compatible as possible with the real (observed) values, some methods [18, 48] introduce a discriminator  $D$  to distinguish which components of  $\tilde{\mathbf{X}}$  are real or fake.  $D$  serves as an adversary to the generator  $G$ , ensuring that the generated data closely resembles the real data. While these approaches can effectively enhance the performance of  $G$ , the potential pitfall is that  $D$  will most likely collapse to a few states since the real value it sees is fixed each time.

In order to make  $G$  better approximate  $q(\mathbf{X}'_{t-1}|\mathbf{X}'_t)$  while increasing the sample diversity, in our framework,  $D$  takes  $\mathbf{X}'_t$ ,  $\tilde{\mathbf{X}}_{t-1}$  and  $t$  as inputs, identify whether each attribute value  $\tilde{x}_{ij}$  come from  $\mathbf{X}'_{t-1}$  or  $\tilde{\mathbf{X}}_{t-1}$ , equivalently predicting  $\mathbf{M}$ . In this way,  $D$  does not merely consider individual positive samples but their noisy versions, making it easier to grasp the data distribution and counteract model overfitting.

**Objective.** Let's summarize the objective functions of the generator  $G$  and discriminator  $D$ . The output of  $G$  at each step is  $\tilde{\mathbf{X}}_0$ , to be as similar as  $\mathbf{X}'_0$  in terms of the observed data. Therefore,  $G$  needs to minimize the reconstruction loss  $\mathcal{L}_R$ :

$$\mathcal{L}_R = \mathbf{M} \odot R(\tilde{\mathbf{X}}_0, \mathbf{X}'_0) \quad (8)$$

where  $R(\cdot)$  represents the calculation of mean squared error (MSE) loss for numerical data and cross-entropy loss for categorical data.

As adversarial training,  $G$  is to make  $D$  unable to distinguish which components of  $\tilde{\mathbf{X}}_{t-1}$  are real in each step, which ensures the consistency of the distributions of  $p_\theta(\mathbf{X}_{t-1}|\mathbf{X}_t)$  and  $q(\mathbf{X}_{t-1}|\mathbf{X}_t)$ . Therefore,  $G$  should minimize the loss  $\mathcal{L}_M$ :

$$\mathcal{L}_M = -\mathbb{E}_{\tilde{\mathbf{X}}_t, \mathbf{M}} \left[ (1 - \mathbf{M}) \log(D(\tilde{\mathbf{X}}_t, t)) \right]. \quad (9)$$

For  $m_{ij} = 1$ , the corresponding loss is 0. And for  $m_{ij} = 0$ , minimizing  $\mathcal{L}_M$  is to incentivize the discriminator to predict it as 1, *i.e.*, the quality of the imputed data is sufficient to pass off as real.

In each step  $t$ , we train  $D$  to maximize the probability of correctly predicting  $\mathbf{M}$ . Consequently, we define the loss  $\mathcal{L}_D$  for the discriminator as:

$$\mathcal{L}_D = -\mathbb{E}_{\tilde{\mathbf{X}}_t, \mathbf{M}} \left[ \mathbf{M} \log D(\tilde{\mathbf{X}}_t, t) + (1 - \mathbf{M}) \log(1 - D(\tilde{\mathbf{X}}_t, t)) \right]. \quad (10)$$

Minimizing  $\mathcal{L}_D$  means minimizing the difference between the predicted mask  $D(\tilde{\mathbf{X}}_t, t)$  and the ground truth  $\mathbf{M}$ , which is pre-determined by  $\mathbf{X}$ .

## 5 DATA UTILITY PLUG-IN

If the imputed relational table  $\tilde{\mathbf{X}}$  is taken as training data, the data utility plug-in can be utilized to ensure the applicability of the imputed data for subsequent classification tasks. We introduce a *Learner*  $L$ , which acts as the classifier used in the subsequent tasks. Thus, we expect the imputed data  $\tilde{\mathbf{X}}$  to minimize the learner's classification loss on the given validation set. This validation loss is denoted by  $\mathcal{L}_L$  in Figure 1, which affects the update of the generator's parameters via backward propagation to generate the expected imputed data.

Suppose the downstream classifier and the validation data are available. We can directly employ them in this plug-in: training the downstream classifier while imputing the missing values in

the training data. Otherwise, we construct the learner  $L$  using a two-layer fully connected neural network, and the validation data will come from  $\tilde{\mathbf{X}}$ . Specifically, we randomly take a portion of the tuples from which all of its attribute values are observed in  $\mathbf{X}$  as the validation set  $\mathbf{X}_v$ , and the rest of the tuples in  $\tilde{\mathbf{X}}$  as the training set  $\mathbf{X}_\eta$  to train  $L$ . The validation loss is determined by the difference between the labels obtained by applying  $L$  to the validation set and the true labels:

$$\mathcal{L}_L = \frac{1}{|\mathbf{X}_v|} \sum_{i=1}^{|\mathbf{X}_v|} H(y_i, p_i) \quad (11)$$

where  $H(\cdot)$  denotes the cross-entropy loss, and  $y_i, p_i$  is the ground truth and predicted label of  $\tilde{x}_i$  respectively.

$\mathcal{L}_L$  is derivable with respect to the parameters  $\theta_G$  of the generator  $G$ . According to the chain rule, we have

$$\frac{\partial \mathcal{L}_L}{\partial \theta_G} = \frac{\partial \mathcal{L}_L}{\partial \theta_L} \cdot \frac{\partial \theta_L}{\partial \mathbf{X}_\eta} \cdot \frac{\partial \mathbf{X}_\eta}{\partial \theta_G}. \quad (12)$$

Here  $\theta_L$  denote the parameter of  $L$ . That is, we can teach  $G$  to generate the training data (*i.e.*, impute the missing values in the training data) that enables the learner to perform well on the validation set (the observed data). If the distribution of  $\mathbf{X}$  and test data in the target task is similar (an assumption often made when applying machine learning models), we can guarantee the utility of our imputed data for training the downstream model. Meanwhile, the learner can contribute to keeping the same data distribution between the imputed and observed data.

Several considerations arise in the design of this plug-in. First, it will be costly to continuously train a new classifier for each generated data instance. As a result, we incorporate  $\mathcal{L}_L$  only after the diffusion process is completely finished. Second, in the case of a high rate of missing data, it becomes difficult to find tuples that do not contain missing values for the validation set. In this situation, we utilize the discriminator  $D$  to construct the most appropriate validation set. Concretely, we assign a confidence score to each tuple in  $\tilde{\mathbf{X}}$ : for each observed attribute value, it has a confidence score of 1, whereas each missing attribute value receives a score derived from  $D$ , *i.e.*, the probability that this value is enough to pass off as real; the confidence score for a tuple is the sum of scores across all its attributes. We then select some top of the tuples after sorting in descending order of their confidence score to compose the validation set.

## 6 DATA DEPENDENCY PLUG-IN

In this section, we introduce our data dependency plug-in, which extracts potential FDs from  $\mathbf{X}$  with missing values and encodes these constraints into the training process of the data imputation module through one extra loss  $\mathcal{L}_F$ . Moreover, we will simultaneously impute missing values and repair inaccurate constraints based on the relative confidence between data and FDs.

### 6.1 Discovery of Functional Dependencies

Berti-Équille et al. [3] have summarized previous work for discovering FDs from data with missing values and mentioned their respective drawbacks. For example, using NULL semantics or directly skipping all tuples with missing values will make many correct FDs undiscoverable. To avoid this, given an FD  $\varphi : Z \rightarrow A$ , we expect

to verify the soundness of  $\varphi$  using all tuples that do not contain missing values on  $Z \cup \{A\}$ , rather than skipping all tuples with missing values on any of the attributes  $U$ .

We adapt TANE [19] to achieve this goal. TANE introduces the notation of a *partition*, denoted by  $\pi$ . For a set of attributes  $Z$ ,  $\pi_Z$  is a collection of *equivalence classes* of tuples, such that each class has a unique value for the attribute set  $Z$ . We redefine the partition. Specifically, given a set of attributes  $Z$ , we omit the missing data to calculate  $\pi_Z$ . Thus, the union of these equivalence classes in  $\pi_Z$  equals a set of tuples with no missing value in  $Z$ . For example,  $\pi_{\text{AreaCode}} = \{\{t_1, t_2, t_3, t_4\}, \{t_6, t_7\}, \{t_8\}\}$ . Then, we redefine the formulation for deriving  $\pi_{Z \cup \{A\}}$  from  $\pi_Z$  and  $\pi_{\{A\}}$  under the new definition of partition:

$$\pi_{Z \cup \{A\}} = (\cup_{z_i \in \pi_Z} (z_i \cap (\cup_{a_j \in \pi_{\{A\}}} a_j))) \setminus \{\emptyset\} \quad (13)$$

This remove tuples from  $\pi_{Z \cup \{A\}}$  that contain missing values on attribute  $Z$  or  $A$ . Under the above setting, an FD  $Z \rightarrow A$  holds if and only if for each  $z_i \in \pi_Z$ , there exists at most one subset of  $z_i$  in  $\pi_{Z \cup \{A\}}$ .

We set a threshold  $\tau$  to limit the maximum number of attributes probed on the left-hand side of an FD. The reason is that when considering an excessive number of attribute combinations, the FD holds mostly because there are no two tuples with the same value in attributes on the left-hand side, which does not work for reasoning about missing values. **Moreover, we will prune some of the candidate FDs during the mining process. For example, if  $Z \rightarrow A$  holds, we will not consider  $Z \cup Z' \rightarrow A$  unless  $Z \rightarrow A$  is found to be invalid in the following phase of refinement. We also filter out FDs that contain the primary key on the left-hand side, which does not contribute to missing value imputation. Thus, the number of potential FDs mined from the incomplete tabular data is not particularly large.** This approach might inadvertently lead to the discovery of some erroneous FDs. The set  $F$  of potential FDs mined from  $\mathbf{X}$  may differ from those  $F^\checkmark$  held in  $\mathbf{X}^\checkmark$ . Thus, they will be refined based on the imputed data in the subsequent process, illustrated in Section 6.3, **where most of the erroneous FDs can be rectified in a few iterations.**

## 6.2 Modeling of Functional Dependencies

Given a set  $F$  of potential FDs mined from  $\mathbf{X}$ , we build a dependency model  $\Phi_\varphi$  for each FD  $\varphi \in F : Z \rightarrow A$ , which takes the attribute values of  $Z$  as input and predicts the attribute value of  $A$ . Thus, we utilize all tuples with no missing value in  $Z \cup \{A\}$  as training data for  $\Phi_\varphi$ , and for each tuple, their attribute values of  $Z$  and  $A$  are taken as features and the label, respectively. Actually, these dependency models could also be pre-trained using the master data or knowledge bases after alignment. This allows more external information to guide the imputation of missing values.

It is worth noting that the dependency model  $\Phi_\varphi$  actually accepts values in all attributes  $U \setminus \{A\}$  as input, not just in  $Z$ . But all attribute values except  $Z$  are set to zero for input to the model. That means we manually perform a “dropout” operation that ignores the inputs in the other attributes  $U \setminus (Z \cup \{A\})$ , but still retains the nodes on the input layer of  $\Phi_\varphi$ . The reason behind this is that as the imputation becomes more accurate,  $\varphi$  might be proven to be incorrect. Thus we need to add more attributes to the left-hand side of  $\varphi$  to make it valid (more details can be found in Section 6.3). In such a case, we

need to “activate” the corresponding input nodes and fine-tune the model instead of training a new dependency model, which will save a lot of training costs without sacrificing the model’s performance.

Then, given the imputed table  $\mathbf{X}$ , we expect it to satisfy all the FDs in  $F$ . That is, the values imputed by the generator  $G$  are desired to be consistent with the predictions of the dependency models. Let  $\mathbf{R}_F = \{A_{\varphi_1}, \dots, A_{\varphi_{|F|}}\}$  where  $A_{\varphi_i}$  be the attribute in the right-hand side of the  $i$ th FD in  $F$ . Then, we define the following loss for the generator:

$$\mathcal{L}_F = \sum_{i=1}^n \sum_{j=1}^{|F|} \|\bar{\mathbf{v}}_{ij} - \mathbf{v}_{ij}\|_1 \quad (14)$$

where  $\|\cdot\|_1$  denotes the L1-norm as the distance measure between two vectors;  $\bar{\mathbf{v}}_{ij}$  and  $\mathbf{v}_{ij}$  are representations of the attribute value of the  $i$ th tuple in attribute  $A_{\varphi_j}$  given by the generator  $G$  and corresponding dependency model  $\Phi_{\varphi_j}$ . In fact, since we take one-hot encoding for the categorical attributes in the data imputation module,  $\bar{\mathbf{v}}_{ij}$  and  $\mathbf{v}_{ij}$  are predicted probability vectors given by the model  $G$  and  $\Phi_{\varphi_j}$  respectively.

## 6.3 Inaccurate FDs Refinement

When mining the potential functional dependencies, we ignore the missing data, assuming that tuples containing missing values always satisfy the mined FDs, which is clearly flawed. As the imputation becomes more accurate, we must go back to determine whether these FDs are valid. In this section, we introduce a framework to refine inaccurate FDs based on the confidence of data and constraints.

**Confidence of Tuple.** The confidence of a tuple is determined by the confidence of its attribute values. Given an FD  $\varphi : Z \rightarrow A$ , the confidence of a tuple  $\bar{\mathbf{x}}_i$  under  $\varphi$  is the minimum value of the confidence of attribute values in  $Z \cup \{A\}$ :

$$c_\varphi(\bar{\mathbf{x}}_i) = \min(c(\bar{x}_{ij}), A_j \in Z \cup \{A\}) \quad (15)$$

**Confidence of attribute value.** It goes without saying that the confidence of observed data should be 1. As for the missing data, we derive the confidence of imputed value  $\bar{x}_{ij}$  based on the predictions of generator  $G$  and discriminator  $D$ . As stated earlier, we adopt one-hot encoding for the categorical attribute. Therefore, for the  $j$ th attribute of the  $i$ th tuple, the output of  $G$  is a predicted probability vector, and  $\bar{x}_{ij}$  is the value with the highest probability, which is denoted by  $p_G(x_{ij})$  and can also be taken as the confidence of imputed value  $\bar{x}_{ij}$  given by  $G$ .  $D$  gives the probability  $p_D(x_{ij})$  that  $\bar{x}_{ij}$  turns out to be observed. The larger  $p_D(x_{ij})$ , the more  $\bar{x}_{ij}$  looks like being real. We have also found from the experiment that  $p_D(x_{ij})$  being larger means  $x_{ij}$  is imputed correctly. Thus, we take it as the confidence of imputed value  $\bar{x}_{ij}$  provided by  $D$ . In brief, the confidence of  $\bar{x}_{ij}$  is the weighted sum of  $p_G(x_{ij})$  and  $p_D(x_{ij})$ :

$$c(\bar{x}_{ij}) = \begin{cases} 1, & m_{ij} = 1 \\ \frac{w_G \times p_G(x_{ij}) + w_D \times p_D(x_{ij})}{w_G + w_D}, & m_{ij} = 0 \end{cases} \quad (16)$$

Here  $w_G$  and  $w_D$  are the accuracy of  $G$  and  $D$  measured from the observed data. Concretely,  $w_G$  is the proportion of observed values  $G$  correctly imputed among all observed values. The more values  $G$  can correctly impute, the more reliable the probability  $p_G(x_{ij})$  it gives.  $w_D$  is the sum of the predicted probabilities over

the observed values divided by the total number of the observed values. The larger  $w_D$  is, the better D is at recognizing observed values, and the more reliable the probability  $p_D(x_{ij})$  it gives.

*Example 6.1.* Given the FD  $\varphi_3 : \text{AreaCode} \rightarrow \text{SingleExemp}$ , let us compute the confidence of  $t_5$  under  $\varphi_3$ . Suppose the predicted probability vector given by G for  $t_5[\text{AreaCode}]$  is (0.05, 0.8, 0.15), and the output from D is 0.9. There are totally 42 observed data in all categorical attributes. Let us say G correctly predicts 32 of them, and the sum of predicted probabilities by D is 38.64. We have  $w_G = 32/42 \approx 0.76$  and  $w_D = 38.64/42 \approx 0.92$ . The confidence of  $t_5[\text{AreaCode}]$  is  $(0.76 \times 0.8 + 0.92 \times 0.9)/(0.76 + 0.92) \approx 0.85$ . Subsequently, we can determine the confidence of  $t_5$  under  $\varphi_3$  as  $\min(0.85, 1) = 0.85$ .

**Confidence of FD.** Following the previous work [3], we construct a subset of tuples  $\bar{X}_\varphi$  from  $\bar{X}$  that satisfy  $\varphi$  to measure the confidence  $c(\varphi)$  of an FD  $\varphi : Z \rightarrow A$ . Berti-Equille et al. [3] continue the frequency-based strategy (the more frequently a value occurs, the more likely it is to be correct) while ignoring each tuple's confidence. To avoid using erroneous tuples to measure the confidence of an FD, we choose tuples with higher confidence, which is consistent with the intuition that the higher the confidence of tuples that satisfy a constraint, the more likely the constraint is correct.

$\bar{X}_\varphi$  should initially contain all tuples with no missing value in attributes  $Z \cup \{A\}$  (The FD  $\varphi$  is extracted from these tuples, so they are naturally supposed to serve as evidence that  $\varphi$  is valid). For the other tuples, we sort them in descending order of their confidence and add them in turn to  $\bar{X}_\varphi$ . Note that if a tuple has already been in  $\bar{X}_\varphi$ , tuples that have violation with it *w.r.t.*  $\varphi$  should be excluded. Finally, we have

$$c(\varphi) = \frac{\sum_{\bar{x}_i \in \bar{X}_\varphi} c_\varphi(\bar{x}_i)}{\sum_{\bar{x}_j \in \bar{X}} c_\varphi(\bar{x}_j)}. \quad (17)$$

It is worth noting that tuples with unique values in attributes  $Z$  are not used to derive the confidence of  $\varphi$  since they do not play any role in determining whether the FD holds. Thus, we remove them from  $\bar{X}_\varphi$  and  $\bar{X}$  before calculating the confidence of rule  $\varphi$ .

*Example 6.2.* Continue to consider  $\varphi_3$  in Example 6.1. Suppose  $t_5[\text{AreaCode}]$  and  $t_3[\text{SingleExemp}]$  are correctly imputed with "501" and "1000" respectively, and the confidence of  $t_3$  is 0.67. We can find from the imputed table that  $t_3$  have the same value with  $\{t_1, t_2, t_4\}$  in AreaCode but different in SingleExemp. Given that  $\{t_1, t_2, t_4\}$  under  $\varphi_3$  are all observed, we put  $\{t_1, t_2, t_4\}$  in  $\bar{X}_{\varphi_3}$ . Similarly,  $\{t_6, t_7\}$  are also added to  $\bar{X}_{\varphi_3}$ . We exclude  $t_8$  since it has a distinct value of AreaCode. Then, the confidence of  $\varphi_3$  is calculated as  $c(\varphi_3) = (1 + 1 + 1 + 1 + 1)/(1 + 1 + 0.67 + 1 + 0.85 + 1 + 1) \approx 0.77$ .

**FD Refinement.** We compare the confidence of a constraint and tuples violating it to determine whether an FD  $\varphi$  is valid, i.e.,  $\varphi \in F^\vee$ . Concretely, the confidence of  $\varphi$  can be calculated by Equation 17. In the meantime, we choose the tuple  $\bar{x}_\varphi^*$  from  $\bar{X} \setminus \bar{X}_\varphi$  with the highest confidence as an indicator that  $\varphi$  may be invalid. It stems from the notion that higher confidence indicates the tuple  $\bar{x}_\varphi^*$  is closer to its true value, and  $\varphi$  obviously cannot hold on  $\bar{X}_\varphi \cup \{\bar{x}_\varphi^*\}$ . When the FD  $\varphi$  has higher confidence than  $\bar{x}_\varphi^*$ , we continue to employ it to guide the imputation of missing values. Otherwise,  $\varphi$  is considered inaccurate and necessitates a refinement.

---

#### Algorithm 1: Training Process of DDGANI

---

**Input:** relational table  $X$  with missing values, mask  $M$ , the number of diffusion steps  $T$ , input noise  $z$ , hyper-parameters such as loss weights  $\alpha, \beta$

**Output:** imputed relational table  $\bar{X}$

- 1  $X' \leftarrow$  pre-processing  $X$  with the self-attention mechanism;
- 2 initialize the parameters of generator G and discriminator D;
- 3 **extract potential FDs  $F$  from  $X$  and build their corresponding dependency models;**
- 4 **repeat**
- 5    $\{X'_1, \dots, X'_T\} \leftarrow$  add noise to  $X'$  over  $T$  steps;
- 6   **for** step  $t \in [T, \dots, 1]$  **do**
- 7      $\hat{X}_0 \leftarrow G(X'_t, z, t)$ ;
- 8     compute the reconstruction loss  $\mathcal{L}_R$  with  $\hat{X}_0$  and  $X'$ ;
- 9      $\bar{X}_0 \leftarrow (1 - M) \odot \hat{X}_0 + M \odot X'$ ;
- 10    **compute the data dependency loss  $\mathcal{L}_F$  with  $\bar{X}_0$ ;**
- 11     $\hat{X}_{t-1} \leftarrow$  posterior sampling with  $\hat{X}_0$  and  $X'_t$ ;
- 12     $\bar{X}_{t-1} \leftarrow (1 - M) \odot \hat{X}_{t-1} + M \odot X'_{t-1}$ ;
- 13     $\bar{M} \leftarrow D(X'_t, \bar{X}_{t-1}, t)$ ;
- 14    compute the discrimination loss  $\mathcal{L}_D$  and  $\mathcal{L}_M$ ;
- 15     $\mathcal{L}_G \leftarrow \mathcal{L}_R + \alpha \mathcal{L}_M + \beta \mathcal{L}_F$ ;
- 16    update the parameters of G, D by  $\mathcal{L}_G, \mathcal{L}_D$  respectively;
- 17   **refine  $F$  based on  $\bar{X}_0$  and retrain the dependency models;**
- 18     $\bar{X}_v \leftarrow$  randomly pick some tuples from  $\bar{X}_0$  whose attribute values are all observable;
- 19    train the learner L with the rest tuples  $\bar{X}_0 \setminus \bar{X}_v$ ;
- 20    evaluate L on  $\bar{X}_v$  and compute  $\mathcal{L}_L$ ;
- 21    update the parameters of G by  $\mathcal{L}_L$ ;
- 22 **until convergence of G and D;**
- 23 **return  $\bar{X}_0$**

---

*Remark.* Here, we can find the calculation of FD confidence in Equation 17 indirectly improves the confidence of the violated tuples required for the refutation of FD, thus ensuring that only a correct pair of tuples conflicting on the FD can give the conclusion that it is inaccurate.

Before discussing specific strategies of FD refinement, we first analyze what changes  $F$  has to make to become  $F^\vee$ . There are four optional steps: (1) add a new FD  $\varphi'$  that  $\varphi' \in F^\vee$ ; (2) remove an FD  $\varphi$  from  $F$  that  $\varphi \notin F^\vee$ ; (3) add some new attributes to the left-hand side of  $\varphi \in F$ ; (4) remove some attributes from the left-hand side of  $\varphi \in F$ . Choice 1 is only required if the missing rate is very high (experimentally found to be >85%). This is uncommon and what we can do is re-mining the FDs from  $\bar{X}$ . Thus, we will not discuss this scenario in detail. Choice 4 will never be adopted, and that is because if  $\varphi$  does not hold, then it cannot be valid after removing an attribute from the left-hand side of  $\varphi$ . Hence, we only consider choice 2 and 3, i.e., when the confidence of an FD is less than the violated tuple, we add some attributes to its left-hand side. If the number of attributes in the left-hand side after appending exceeds the threshold  $\tau$ , we discard this FD when it still does not hold.



How to repair an FD by adding more attributes to the left-hand side is not a new problem and has been studied in previous work such as [11]. The difference is that we will try to find all possible solutions for appending attributes. For example, upon discovering that `Areacode`, `MaritalStatus`  $\rightarrow$  `SingleExemp` is one way to fix the erroneous FD `Areacode`  $\rightarrow$  `SingleExemp`, we will not stop there and continue to find another way to refine it such as `Areacode`, `MarriedExemp`  $\rightarrow$  `SingleExemp`. This allows us to dig into more dependencies between attributes, which can help with the missing value imputation.

Here, we could adopt some simple sorting of the candidate attributes appended to an FD  $\varphi : Z \rightarrow A$ . We need the power set of  $U \setminus (Z \cup \{A\})$  (excluding  $\emptyset$  and sets larger than  $\tau - |Z|$ ). Each item in the power set is taken in ascending order of the size into consideration whether it can be appended to the left-hand side of  $\varphi$ . Here, if appending an attribute set  $Z'$  yields an FD  $\varphi' : Z \cup Z' \rightarrow A$ , and  $c(\varphi') \geq c(\bar{x}_{\varphi'}^*)$ , then no more super sets of  $Z'$  are considered. We can also choose to combine the correlation of attributes but without associated theoretical guarantees: if appending  $Z'$  does not result in a valid FD  $\varphi' : Z \cup Z' \rightarrow A$ , we will not consider other candidates that have a weaker correlation with  $A$  than  $Z'$  referring to our correlation matrix **C**. We discard  $\varphi$  if it still cannot hold after trying all the possible appends.

**Retraining of Dependency Model.** If more attributes are added to the left-hand side of FD  $\varphi$  ( $\varphi$  is refined to  $\varphi' : Z \cup Z' \rightarrow A$ ), we need to retrain the dependency model  $\Phi_\varphi$  after enabling the input nodes of attributes  $Z'$  to obtain a new model  $\Phi_{\varphi'}$ . That is, we utilize all tuples in  $\bar{X}_{\varphi'}$  to retrain  $\Phi_\varphi$ : for each tuple, their attribute values of  $Z \cup Z'$  and  $A$  are taken as features and the label, respectively, and all attribute values except  $Z \cup Z'$  still need to be set to zero for input to the model. Since  $\Phi_\varphi$  has already captured some inner connection between  $Z$  and  $A$ , the retraining mainly goes to determine the effect of the newly added attributes on  $A$ , which will not take long. If the FD  $\varphi$  is entirely discarded, the corresponding dependency model  $\Phi_\varphi$  is also eliminated with it.

## 7 ENTIRE TRAINING PROCESS OF DDGANI

The overall workflow of DDGANI is presented in Algorithm 1. Let us first focus on the data imputation module, which is in black. The missing data is initially pre-processed with the attention mechanism to obtain  $X'$  (line 1). We then initialize the parameters of generator  $G$  and discriminator  $D$  (line 2) and train the models (lines 4-22). A forward diffusion process is applied on  $X'$  in the beginning of each epoch to produce  $\{X'_1, \dots, X'_T\}$  (line 5). Then in each step  $t$  of denoising,  $G$  takes  $X'_t$ ,  $t$  and a random Gaussian noise  $z$  as input and outputs  $\hat{X}_0$  (line 7). The loss  $\mathcal{L}_R$  is computed based on the differences between  $\hat{X}_0$  and  $X'$  (line 8). Then, a posterior sampling process derives  $\hat{X}_{t-1}$  from  $\hat{X}_0$  and  $X'_t$  (line 11). We integrate  $\hat{X}_{t-1}$  and  $X'_{t-1}$  to obtain the final imputed table  $\bar{X}_{t-1}$  at step  $t-1$ , which is fed into  $D$  with  $X'_t$  and  $t-1$  to predict the mask  $M$  (lines 12-13). The differences between the predicted result and  $M$  will be used to compute the loss  $\mathcal{L}_D$  and  $\mathcal{L}_M$ , which are used to update the parameters of  $D$  and  $G$  along with  $\mathcal{L}_R$ , respectively (lines 14-16).

Then, let us look at the usage of the data utility plug-in, which is in gray. After the data imputation module has finished its work, we randomly draw some tuples without missing values from  $\bar{X}_0$  as the

Table 2: Dataset Description

Dataset	#Instance	#Numerical	#Categorical	#FDs
Wine	178	13	0	0
Wireless	2000	7	0	0
Spam	4601	57	0	0
Tic-Tac-Toe	958	0	9	0
Flare	1066	0	12	0
Balance	625	0	4	0
Adult	30148	5	9	2
Hospital	1000	2	7	8
Tax	100000	2	7	9
Thoracic	470	3	13	0
Contraceptive	1473	2	7	0
Australian	690	6	8	0
Beers	2410	3	3	0

validation set  $\bar{X}_0$ , and a learner  $L$  is trained using the remaining data (lines 18-19).  $\mathcal{L}_L$  is the validation loss of  $L$ , which will be used to update the parameters of  $G$  via meta-learning paradigm (lines 20-21).

The data dependency plugin is required to be integrated for use in the data imputation module, which is in brown. We will first mine potential FDs at the beginning and build their corresponding dependency models (line 3). After obtaining  $\hat{X}_0$  at each step  $t$ , we combine  $\hat{X}_0$  and  $X'$ , yielding the imputed data  $\bar{X}_0$  and calculating the loss  $\mathcal{L}_F$  (lines 9-10). This loss will be added into  $\mathcal{L}_G$  to update the parameters of  $G$  (line 15). These FDs may be refined according to the confidence of tuples and constraints after the missing value imputation (line 17).

## 8 EXPERIMENTS

### 8.1 Experimental Settings

**Dataset.** We evaluate our approach on **thirteen** publicly available datasets, including three datasets with all numerical attributes *Wine*, *Wireless*, *Spam*, three datasets with all categorical attributes *Tic-Tac-Toe*, *Flare*, *Balance*, and seven datasets with both numerical and categorical attributes *Adult*, *Hospital* [35], *Tax* [5], *Thoracic*, *Contraceptive*, *Australian*, *Beers*<sup>1</sup>. Datasets not specifically attributed are from the UCI Machine Learning Repository (<https://archive.ics.uci.edu/>). Table 2 presents an overview of these datasets, including the number of tuples, numerical attributes, categorical attributes, and functional dependencies. Here, the FDs are mined by TANE [19] from the complete dataset and treated as the true FDs in Section 8.6.

Each dataset is randomly divided into three parts: 70% of the training set with missing values injected to evaluate the accuracy of imputation methods, 20% of the test set which remains complete and is used to report the accuracy of the downstream classifier, and 10% of the validation set for the hyper-parameter searching.

**Missing Value Injection.** We remove some attribute values from the datasets completely at random with probability  $p$  (*Missing Completely at Random*, *MCAR*). The missing rate  $p$  is set to 20% by default, and it is changed from 10% to 50% in Section 8.3 to show the effectiveness of our proposed method under different missing rates. Section 8.4 shows the evaluation of our method under *MAR*

<sup>1</sup><https://www.kaggle.com/datasets/nickhould/craft-cans>

**Table 3: Overall Evaluation of the Model (Average  $\pm$  Std)**

Dataset	Metrics	Mean	MissFI	GAIN	DataWig	Baran	AimNet	NOT-MIWAE	Graf	VAIM	EGG-GAE	DDGANI
Wine	ARMSE	0.296 $\pm$ 0.000	0.151 $\pm$ 0.003	0.171 $\pm$ 0.007	0.163 $\pm$ 0.005	0.225 $\pm$ 0.013	0.241 $\pm$ 0.014	0.169 $\pm$ 0.004	0.297 $\pm$ 0.014	0.182 $\pm$ 0.013	0.150 $\pm$ 0.004	0.145 $\pm$ 0.001
	AMAE	0.250 $\pm$ 0.000	0.113 $\pm$ 0.002	0.135 $\pm$ 0.007	0.129 $\pm$ 0.004	0.153 $\pm$ 0.010	0.185 $\pm$ 0.010	0.128 $\pm$ 0.004	0.251 $\pm$ 0.012	0.143 $\pm$ 0.010	0.120 $\pm$ 0.003	0.106 $\pm$ 0.001
Wireless	ARMSE	0.352 $\pm$ 0.000	0.120 $\pm$ 0.000	0.140 $\pm$ 0.007	0.119 $\pm$ 0.002	0.231 $\pm$ 0.018	0.224 $\pm$ 0.012	0.121 $\pm$ 0.001	0.355 $\pm$ 0.001	0.148 $\pm$ 0.006	0.164 $\pm$ 0.008	0.118 $\pm$ 0.001
	AMAE	0.318 $\pm$ 0.000	0.097 $\pm$ 0.000	0.109 $\pm$ 0.006	0.094 $\pm$ 0.001	0.187 $\pm$ 0.017	0.165 $\pm$ 0.007	0.095 $\pm$ 0.001	0.322 $\pm$ 0.001	0.117 $\pm$ 0.005	0.134 $\pm$ 0.007	0.093 $\pm$ 0.001
Spam	ARMSE	0.053 $\pm$ 0.000	0.048 $\pm$ 0.000	0.051 $\pm$ 0.001	0.052 $\pm$ 0.001	0.057 $\pm$ 0.002	0.051 $\pm$ 0.007	0.056 $\pm$ 0.001	0.054 $\pm$ 0.002	0.054 $\pm$ 0.001	0.049 $\pm$ 0.001	0.047 $\pm$ 0.000
	AMAE	0.022 $\pm$ 0.000	0.023 $\pm$ 0.000	0.022 $\pm$ 0.000	0.025 $\pm$ 0.000	0.016 $\pm$ 0.001	0.022 $\pm$ 0.001	0.030 $\pm$ 0.001	0.022 $\pm$ 0.000	0.020 $\pm$ 0.003	0.019 $\pm$ 0.000	0.021 $\pm$ 0.000
Tic-Tac-Toe	ARMSE	0.590 $\pm$ 0.009	0.529 $\pm$ 0.008	0.537 $\pm$ 0.013	0.514 $\pm$ 0.008	0.607 $\pm$ 0.033	0.525 $\pm$ 0.007	0.529 $\pm$ 0.019	0.523 $\pm$ 0.006	0.582 $\pm$ 0.006	0.527 $\pm$ 0.003	0.512 $\pm$ 0.002
	AMAE	0.590 $\pm$ 0.009	0.529 $\pm$ 0.008	0.537 $\pm$ 0.013	0.514 $\pm$ 0.008	0.607 $\pm$ 0.033	0.525 $\pm$ 0.007	0.529 $\pm$ 0.019	0.523 $\pm$ 0.006	0.582 $\pm$ 0.006	0.527 $\pm$ 0.003	0.512 $\pm$ 0.002
Flare	ARMSE	0.260 $\pm$ 0.000	0.240 $\pm$ 0.000	0.200 $\pm$ 0.009	0.205 $\pm$ 0.007	0.241 $\pm$ 0.009	0.241 $\pm$ 0.002	0.226 $\pm$ 0.004	0.262 $\pm$ 0.004	0.240 $\pm$ 0.014	0.183 $\pm$ 0.001	0.173 $\pm$ 0.003
	AMAE	0.260 $\pm$ 0.000	0.240 $\pm$ 0.000	0.200 $\pm$ 0.009	0.205 $\pm$ 0.007	0.241 $\pm$ 0.009	0.241 $\pm$ 0.002	0.226 $\pm$ 0.004	0.262 $\pm$ 0.004	0.240 $\pm$ 0.014	0.183 $\pm$ 0.001	0.173 $\pm$ 0.003
Balance	ARMSE	0.831 $\pm$ 0.004	0.818 $\pm$ 0.023	0.823 $\pm$ 0.004	0.818 $\pm$ 0.020	0.828 $\pm$ 0.019	0.783 $\pm$ 0.013	0.822 $\pm$ 0.008	0.802 $\pm$ 0.008	0.814 $\pm$ 0.016	0.776 $\pm$ 0.001	0.772 $\pm$ 0.007
	AMAE	0.831 $\pm$ 0.004	0.818 $\pm$ 0.023	0.823 $\pm$ 0.004	0.818 $\pm$ 0.020	0.828 $\pm$ 0.019	0.783 $\pm$ 0.013	0.822 $\pm$ 0.008	0.802 $\pm$ 0.008	0.814 $\pm$ 0.016	0.776 $\pm$ 0.001	0.772 $\pm$ 0.007
Adult	ARMSE	0.341 $\pm$ 0.000	0.258 $\pm$ 0.000	0.250 $\pm$ 0.007	0.243 $\pm$ 0.004	0.292 $\pm$ 0.000	0.293 $\pm$ 0.001	0.254 $\pm$ 0.005	0.338 $\pm$ 0.000	0.245 $\pm$ 0.012	0.217 $\pm$ 0.003	0.195 $\pm$ 0.000
	AMAE	0.325 $\pm$ 0.000	0.244 $\pm$ 0.000	0.234 $\pm$ 0.008	0.230 $\pm$ 0.004	0.273 $\pm$ 0.000	0.272 $\pm$ 0.001	0.238 $\pm$ 0.007	0.323 $\pm$ 0.000	0.229 $\pm$ 0.012	0.203 $\pm$ 0.003	0.180 $\pm$ 0.000
Hospital	ARMSE	0.618 $\pm$ 0.000	0.254 $\pm$ 0.003	0.365 $\pm$ 0.026	0.398 $\pm$ 0.005	0.276 $\pm$ 0.062	0.169 $\pm$ 0.003	0.568 $\pm$ 0.012	0.267 $\pm$ 0.003	0.600 $\pm$ 0.002	0.147 $\pm$ 0.008	0.133 $\pm$ 0.008
	AMAE	0.610 $\pm$ 0.000	0.239 $\pm$ 0.003	0.351 $\pm$ 0.026	0.387 $\pm$ 0.005	0.255 $\pm$ 0.061	0.139 $\pm$ 0.001	0.554 $\pm$ 0.012	0.258 $\pm$ 0.003	0.580 $\pm$ 0.003	0.137 $\pm$ 0.008	0.115 $\pm$ 0.002
Tax	ARMSE	0.615 $\pm$ 0.000	0.380 $\pm$ 0.011	0.558 $\pm$ 0.035	0.367 $\pm$ 0.019	0.534 $\pm$ 0.011	0.422 $\pm$ 0.026	0.517 $\pm$ 0.027	0.369 $\pm$ 0.082	0.586 $\pm$ 0.042	0.353 $\pm$ 0.012	0.349 $\pm$ 0.004
	AMAE	0.600 $\pm$ 0.000	0.360 $\pm$ 0.009	0.526 $\pm$ 0.039	0.352 $\pm$ 0.020	0.497 $\pm$ 0.008	0.391 $\pm$ 0.026	0.505 $\pm$ 0.029	0.355 $\pm$ 0.076	0.569 $\pm$ 0.047	0.337 $\pm$ 0.008	0.325 $\pm$ 0.003
Thoracic	ARMSE	0.170 $\pm$ 0.005	0.161 $\pm$ 0.010	0.175 $\pm$ 0.004	0.154 $\pm$ 0.001	0.142 $\pm$ 0.014	0.198 $\pm$ 0.001	0.154 $\pm$ 0.006	0.166 $\pm$ 0.008	0.153 $\pm$ 0.005	0.151 $\pm$ 0.001	0.139 $\pm$ 0.007
	AMAE	0.158 $\pm$ 0.005	0.151 $\pm$ 0.009	0.163 $\pm$ 0.003	0.146 $\pm$ 0.000	0.129 $\pm$ 0.014	0.171 $\pm$ 0.001	0.145 $\pm$ 0.006	0.155 $\pm$ 0.009	0.143 $\pm$ 0.004	0.148 $\pm$ 0.001	0.132 $\pm$ 0.007
Contraceptive	ARMSE	0.351 $\pm$ 0.003	0.328 $\pm$ 0.003	0.341 $\pm$ 0.006	0.301 $\pm$ 0.007	0.343 $\pm$ 0.017	0.337 $\pm$ 0.003	0.321 $\pm$ 0.006	0.313 $\pm$ 0.004	0.326 $\pm$ 0.011	0.309 $\pm$ 0.001	0.290 $\pm$ 0.001
	AMAE	0.339 $\pm$ 0.003	0.318 $\pm$ 0.002	0.331 $\pm$ 0.006	0.292 $\pm$ 0.007	0.331 $\pm$ 0.016	0.328 $\pm$ 0.003	0.312 $\pm$ 0.006	0.302 $\pm$ 0.004	0.317 $\pm$ 0.010	0.303 $\pm$ 0.001	0.281 $\pm$ 0.001
Australian	ARMSE	0.295 $\pm$ 0.004	0.269 $\pm$ 0.013	0.264 $\pm$ 0.007	0.264 $\pm$ 0.001	0.297 $\pm$ 0.008	0.255 $\pm$ 0.006	0.260 $\pm$ 0.004	0.265 $\pm$ 0.008	0.272 $\pm$ 0.006	0.235 $\pm$ 0.000	0.237 $\pm$ 0.003
	AMAE	0.273 $\pm$ 0.004	0.254 $\pm$ 0.014	0.247 $\pm$ 0.008	0.246 $\pm$ 0.002	0.278 $\pm$ 0.008	0.232 $\pm$ 0.004	0.243 $\pm$ 0.004	0.247 $\pm$ 0.005	0.255 $\pm$ 0.004	0.230 $\pm$ 0.000	0.226 $\pm$ 0.002
Beers	ARMSE	0.152 $\pm$ 0.000	0.112 $\pm$ 0.001	0.147 $\pm$ 0.000	0.113 $\pm$ 0.000	0.131 $\pm$ 0.028	0.145 $\pm$ 0.000	0.120 $\pm$ 0.001	0.152 $\pm$ 0.000	0.117 $\pm$ 0.000	0.111 $\pm$ 0.000	0.104 $\pm$ 0.001
	AMAE	0.083 $\pm$ 0.000	0.049 $\pm$ 0.000	0.080 $\pm$ 0.000	0.050 $\pm$ 0.000	0.072 $\pm$ 0.004	0.080 $\pm$ 0.000	0.055 $\pm$ 0.002	0.083 $\pm$ 0.000	0.052 $\pm$ 0.000	0.050 $\pm$ 0.000	0.047 $\pm$ 0.001

(Missing at Random), MNAR (Missing Not at Random) and Regional Missing. More details can be found in Section 8.4. Note that Beers is a real-world dataset with 1067 missing values. Thus, we do not inject more missing values into Beers, and the ground truth is obtained through crowdsourcing techniques.

**Evaluation Metrics.** We employ two commonly used metrics to measure the accuracy of imputation methods: *average root mean square error* (ARMSE) [41] and *average mean absolute error* (AMAE) [33]. ARMSE calculates the average of the errors for each attribute. Specifically, it utilizes the *root mean square error* (RMSE) [21] for numerical attributes and the *accuracy error* (AR) [33] for categorical ones. The smaller ARMSE corresponds to the better imputation result. Similar to ARMSE, AMAE employs *mean absolute error* (MAE) [21] for numerical attributes and AR for categorical ones. The smaller AMAE corresponds to the better imputation result. We also measure the *accuracy* of the downstream classifier in some experiments to verify the utility of the imputed dataset.

**Baseline.** We compare DDGANI against the following data imputation methods. (1) *Mean* [16], which imputes the missing values by the mean of each numerical attribute and the mode of each categorical attribute; (2) *MissFI* [41], which utilizes the missing forest to impute the missing data; (3) *GAIN* [48], which employs a GAN-based model and utilizes a hinting mechanism to optimize the model for imputing missing values; (4) *DataWig* [4], which trains an imputation model for each attribute that imputes the missing values based on the observed values in all other attributes. (5) *Baran* [27], which generates potential imputations using pre-trained error corrector models and selects representative tuples for manual correction to fine-tune the corrector models. (6) *AimNet* [46], which utilizes a transformer to reconstruct missing data, also based on FDs and the attention scores between attributes. (7) *NOT-MIWAE* [20], which

applies deep latent variable models (DLVMs) to the problem of missing data imputation; (8) *Graf* [34], which uses sequence GANs to learn conditional functional dependencies from the dirty data and imputes the missing values based on these learned CFDs. (9) *VAIM* [18], which utilizes a VAE-based model with attention mechanisms and learner to optimize the generative model for estimating missing values; (10) *EGG-GAE* [43], which imputes the missing values by adopting graph neural networks; We also compare with pre-trained language model based methods including *ChatGPT* and *RPT* [42], and traditional rule-based approach *Holistic* [13] for data imputation.

**Hyperparameter Settings.** We set the diffusion timestep  $T$  to 6, the threshold  $\kappa$  for the attention mechanism to 0.05, and the limit  $\tau$  for the number of attributes on the left-hand side of FDs to 3. The diffusion rate  $\beta_t$  at step  $t$  is defined to be linearly increasing constants between  $\beta_1 = 10^{-4}$  and  $\beta_T = 2 \times 10^{-2}$ . Other hyperparameters in DDGANI and each competing method are decided via cross-validation.

**Experimental Environment.** All experiments were conducted on an Intel(R) Xeon(R) Silver 4210R 2.40GHz server with NVIDIA GeForce RTX 3090 24GiB (GPU) and 94GB of RAM, running Ubuntu 20.04 system.

## 8.2 Overall Evaluation of the Model

We first evaluate the performance of DDGANI on *twelve* datasets with 20% missing rate under the missing completely at random (MCAR) model, and *Beers* with existing missing values. The results are presented in Table 3. We conduct five runs for each method and record the average and standard deviation of the experimental results. The best scores are in bold, and the second-best scores are underlined. The results of our method are on a gray background.

Table 3 shows that DDGANI outperforms other methods in most datasets. This can be attributed to the way DDGANI approaches the target distribution through a series of incremental approximations. This iterative approach focuses on refining the current distribution rather than learning a distribution from scratch, allowing the model to better capture subtle nuances and intricate dependencies within the data. Additionally, the utilization of the two plug-ins also helps the model better learn the data distribution and the data dependencies.

Furthermore, it’s worth noting that GAIN, NOT-MIWAE, and VAIM all exhibit subpar performance when dealing with datasets containing a mix of numerical and categorical attributes. This is due to the following reasons: GAIN suffers from mode collapse when generating categorical data, which may result in overfitting certain observed values in the categorical attribute while ignoring the true data distribution. NOT-MIWAE and VAIM have the assumption that the latent space obeys a Gaussian distribution and generates samples from a complex distribution in one shot, which is not optimal for the dataset with multiple categorical attributes. DataWig makes an initial guess for the missing values using Mean imputation, and the prediction of missing values in one attribute needs to be based on the imputation results of other attributes. Thus, if it imputes one attribute incorrectly, it can greatly affect the accuracy of its imputation on other attributes. Baran achieves lower AMAE than DDGANI in two datasets, Spam and Thoracic, but it incorporates user supervision in the form of a number of manual imputations for examples with missing values, which leads to human costs. AimNet is also based on the attention mechanism and FDs for missing value imputation. But it only utilizes the attention scores between different attributes, while DDGANI not only considers the correlation between attributes but also different tuples, thereby fully utilizing the information provided from similar tuples to impute the missing values. Moreover, AimNet relies on the given FDs for missing value imputation without refinement, which can easily go wrong if the FDs are inaccurate. Graf utilizes CFDs for missing value imputation and co-cleans the inaccurate CFDs and missing values simultaneously. However, it employs sequence GANs to generate CFDs, which requires a lot of clean training data. And it can only impute the attribute values under the learned CFDs (In our experiments, we use Mean imputation for other attributes Graf could not touch). EGG-GAE works better with mixed datasets, and has lower ARMSE in the Australian dataset than DDGANI. It adopts graph neural networks for missing value imputation and reconstructs a tuple with missing value by leveraging neighboring nodes in the graph. Thus, it leads to difficulty in correctly imputing the missing values of “outlier”, the tuple that can not find similar ones in the entire tabular data. Instead, DDGANI utilizes a generative model to rebuild the entire relational table step by step based on the distribution of all observed data, unlike EGG-GAE which strongly relies on similar tuples.

It is worth noting that all methods have high imputation errors in the Balance dataset. This is because most existing methods, including DDGANI, rely strongly on attribute correlations, making it difficult to impute the missing values appearing in the attribute that is not strongly correlated with other attributes. The correlations between different attributes are all zero in the Balance dataset, causing all methods to fail.

**Table 4: Comparison with PLMs**

Dataset	Metrics	ChatGPT	RPT	DDGANI
Wine	ARMSE	0.211	0.193	<b>0.145</b>
	AMAE	0.174	0.159	<b>0.106</b>
	Accuracy	94.44	97.22	<b>99.78</b>
Flare	ARMSE	0.258	0.250	<b>0.173</b>
	AMAE	0.258	0.250	<b>0.173</b>
	Accuracy	78.66	78.60	<b>80.37</b>
Adult	ARMSE	0.312	0.301	<b>0.195</b>
	AMAE	0.301	0.270	<b>0.180</b>
	Accuracy	82.13	83.33	<b>85.12</b>

**Comparison with Pre-trained Language Models.** Since the tuples in tabular data are often viewed as sequences [42], which is similar to a natural language processing task, we compare DDGANI additionally with pre-trained language models (PLMs) including ChatGPT and RPT [42] on the missing value imputation task. In this part of the experiments, we only show the results of Wine, Flare and Adult in Table 4 as the representatives of datasets with all numerical attributes, all categorical attributes and mixed types of attributes. Actually, our experiments on thirteen datasets all reveal that DDGANI has the best imputation performance, particularly on datasets consisting entirely of numerical data. The reason is that PLMs lack sufficient prior knowledge when confronted with purely numerical data. They are also less capable of capturing the underlying data distribution than deep generative models.

### 8.3 Comparison under Different Missing Rates

To better showcase the effectiveness of DDGANI, we vary the missing rate from 10% to 50% to test the performance of each method. In this set of experiments, we still choose Wine, Flare and Adult as the representatives of the datasets. The results are given in Figure 3.

Obviously, as the missing rate increases, the performance of each imputation method declines. This is because the decrease in observed data makes it difficult to estimate the missing values, especially when the data near the missing values are gradually unavailable. Nevertheless, DDGANI consistently performs better than other baselines. As the missing rate increases, the imputation error of our method increases at a slower rate than that of GAIN and VAIM, which also employ generative models. Specifically, with the missing rate changing from 0.1 to 0.5 in Adult, the ARMSE of DDGANI increases from 0.179 to 0.247, while GAIN and VAIM increase from 0.213 to 0.298 and from 0.214 to 0.345, respectively. The reason behind this is that DDGANI adopts the diffusion model, which is better at predicting missing values from the context through several conditional denoising diffusion steps. Meanwhile, the usage of self-attention mechanisms and data dependency plug-in allows the model to better capture the correlation between missing values and the information of observed data.

### 8.4 Comparison under Different Missing Types

In this set of experiments, we evaluate the impact of different missing mechanisms on the performance of imputation methods. According to a recent study [20], we inject different types of missing values. 1) MCAR: removing some data from the dataset completely at random; 2) MAR: randomly selecting a numerical attribute to sort

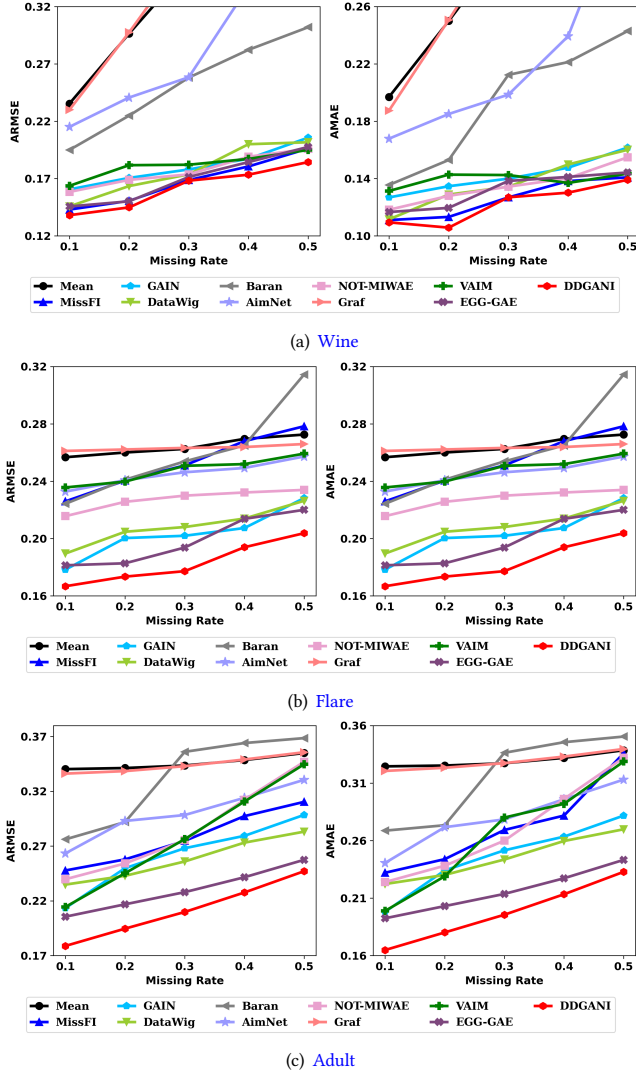


Figure 3: Comparison under Different Missing Rates

the tuples and then select some consecutive tuples to inject missing values into its other attributes; 3) MNAR: randomly removing two types of attribute values; one is the data in a numerical attribute that is less than the median, and the other is from a categorical attribute; Particularly, we inject 4) Regional Missing: designating a rectangular area by randomly selecting several consecutive attributes from some consecutive tuples and subsequently removing all data within this area. This kind of missing data often appears as the result of an outer join of relational tables.

The results shown in Figure 4 demonstrate the superiority of DDGANI. The advantages of our method are more evident under MAR, MNAR, and regional missing. Specifically, under the MAR mechanism, DDGANI achieves an ARMSE of 0.251 in Adult, notably superior to GAIN and VAIM with 0.322 and 0.357, respectively. This implies that our method can accurately capture the dependencies between attribute values. Regional missing leads to the loss of substantial sequential data. In this situation, our attention mechanism

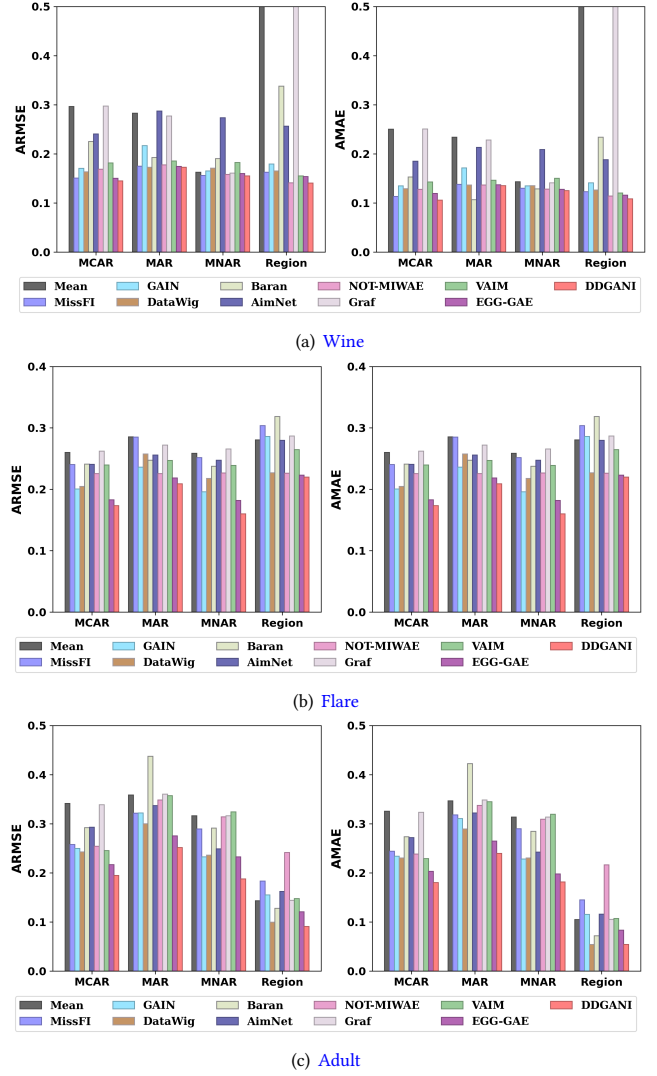


Figure 4: Comparison under Different Missing Types

and data dependency plug-in help the model to place its eyes on more distant tuples even if it cannot find clues from nearby data to impute the missing values.

### 8.5 Utility of the Imputed Data

In this part of the experiments, we train a classifier with the imputed data and measure its accuracy on the test set to demonstrate the applicability of imputed training data. Here, we fix the missing rate at 20% under the MCAR mechanism and employ Random Forest as the downstream classifier. It is worth noting that Hospital and Tax are label-free. In order to conduct this experiment, we take the values in attribute EmergencyService and HasChild as labels for these two datasets, respectively, since these two attributes do not directly depend on other attributes but rather exhibit some correlations with most other attributes.

The experimental results are shown in Table 5, the first twelve columns. *Intuitively, the smaller the imputation error, the higher*



Table 5: Accuracy of Downstream Classifier

Dataset	Mean	MissFI	GAIN	DataWig	Baran	AimNet	NOT-MIWAE	Graf	VAIM	EGG-GAE	DDGANI	Removal	Original
Wine	94.44±0.00	99.44±1.11	97.78±1.11	98.89±1.36	94.44±1.76	93.52±1.31	98.89±1.36	95.00±3.24	98.89±1.36	98.89±1.36	99.78±0.12	88.89±0.00	100.00±0.00
Wireless	97.00±0.00	97.35±0.00	97.35±0.46	97.25±0.43	96.95±0.71	96.67±2.36	97.25±0.32	96.30±0.10	96.70±0.40	96.55±0.21	97.70±0.25	96.75±0.00	98.25±0.00
Spam	93.38±0.00	93.25±0.00	93.05±0.08	93.22±0.33	92.70±0.52	92.67±2.36	92.68±0.43	93.21±0.63	93.05±0.48	92.94±0.52	93.55±0.31	-	95.55±0.00
Tic-Tac-Toe	80.21±0.00	83.44±0.88	86.98±1.00	85.83±0.53	86.77±1.14	83.68±0.49	87.90±0.91	86.67±0.63	85.94±1.57	86.29±1.07	88.72±1.30	77.92±0.00	98.85±0.00
Flare	78.60±0.00	80.06±0.79	79.91±1.34	79.93±1.08	80.09±1.34	77.57±1.01	80.06±1.10	80.22±1.34	79.60±0.22	79.13±0.96	80.37±1.10	78.60±0.00	81.79±0.00
Balance	69.20±0.00	78.13±1.60	76.00±1.36	65.73±1.56	64.00±0.00	73.60±0.73	73.33±2.10	75.73±0.75	77.07±0.99	80.27±0.75	80.33±0.31	77.60±0.00	82.72±0.00
Adult	84.92±0.00	85.00±0.05	84.80±0.27	84.41±0.22	81.24±0.00	82.62±0.43	84.59±0.20	84.21±0.32	84.92±0.21	84.85±0.25	85.12±0.09	81.43±0.00	85.34±0.00
Hospital	96.50±0.00	96.50±0.00	96.60±0.20	95.80±0.87	96.70±0.24	96.83±0.24	96.60±0.20	96.70±0.24	96.50±0.0	96.70±0.27	97.00±0.00	95.00±0.00	97.00±0.00
Tax	86.48±0.00	87.04±0.82	85.75±0.56	86.30±0.42	84.93±0.72	84.91±0.37	86.82±0.40	86.77±0.61	86.58±0.61	86.43±0.27	88.73±0.66	86.29±0.00	89.41±0.00
Thoracic	78.72±0.00	80.14±1.00	79.57±0.43	78.72±0.00	78.94±1.04	78.01±1.33	78.72±0.95	78.09±0.85	79.15±0.52	78.72±0.21	81.79±0.00	77.45±0.00	82.36±0.00
Contraceptive	53.15±0.00	52.09±0.32	51.53±0.61	51.32±1.40	49.36±1.78	49.27±1.76	54.10±0.92	49.22±0.54	51.59±1.36	49.49±0.96	55.59±0.36	49.56±0.00	55.98±0.00
Australian	86.26±0.00	86.71±0.34	87.25±1.18	80.00±1.45	81.01±1.18	87.20±0.49	87.54±1.55	84.06±0.65	87.54±1.80	88.05±0.68	88.13±0.29	82.90±0.00	89.86±0.00
Beers	85.48±0.00	88.66±1.02	86.47±1.89	90.80±0.84	86.17±1.09	86.69±0.10	93.19±0.61	85.48±0.00	93.85±1.27	93.65±0.03	96.26±0.54	68.88±0.00	100.00±0.00

Table 6: An Ablation Study

Dataset	Metrics	DDGANI	w/o $\mathcal{L}_L$	w/o $\mathcal{L}_F$
Wine	ARMSE	0.145±0.001	0.157±0.002	-
	AMAE	0.106±0.001	0.121±0.002	-
	Accuracy	99.78±0.12	98.33±0.22	-
Flare	ARMSE	0.173±0.003	0.194±0.002	-
	AMAE	0.173±0.003	0.194±0.002	-
	Accuracy	80.37±1.10	78.88±1.16	-
Adult	ARMSE	0.195±0.000	0.206±0.000	0.201±0.002
	AMAE	0.180±0.000	0.192±0.000	0.190±0.001
	Accuracy	85.12±0.09	84.58±0.09	84.96±0.13
Tax	ARMSE	0.349±0.004	0.360±0.006	0.366±0.004
	AMAE	0.325±0.003	0.340±0.004	0.345±0.001
	Accuracy	88.73±0.66	86.70±0.30	87.62±0.62

the accuracy of the downstream classifier. This has indeed been validated in many datasets. For example, in the Wine dataset, the Mean, AimNet, Baran, and Graf result in higher ARMSE, and the downstream classifier training by their imputed data is relatively poorer. However, as mentioned in [37], many factors, not just the imputation error of the training set, affect the performance of the downstream classifier. The impact of data quality on downstream classification tasks varies across different datasets. EGG-GAE, for example, does well in the missing value imputation task but cannot result in the superior performance of the downstream classifier. This is because of its difficulty in dealing with the outliers, which misleads the classifier. The deviation of the distribution of imputed outliers from the true distribution greatly affects the training of the classifier. DDGANI outperforms other methods in all datasets. The superior performance of DDGANI is attributed to the accurate imputation and the data utility plug-in, which could [keep the imputed data as close to the original distribution as possible](#), allowing missing value imputation to be done in a direction that helps improve the performance of the downstream classifier.

**Effectiveness of Data Utility Plug-in.** We conduct an ablation study to delve deeper into the specific impact of data utility plug-in ( $\mathcal{L}_L$ ) on the model performance. Table 6 shows that removing  $\mathcal{L}_L$  does adversely affect the utility of imputed data on downstream tasks more than removing the data dependency plug-in ( $\mathcal{L}_F$ ). To further verify that the optimization of  $\mathcal{L}_L$  can improve data estimation without increasing bias, we also test the performance of the classifier training with the original dataset where no missing values

Table 7: Comparison with Holistic (ARMSE)

Dataset	Rate	Holistic			DDGANI	
		T	P	R	P	R
Hospital	10%	<b>0.082</b>	0.123	0.103	0.114	0.091
	50%	0.396	0.513	0.456	0.401	<b>0.347</b>
Tax	10%	0.186	0.211	0.195	0.179	<b>0.164</b>
	50%	0.365	0.414	0.387	0.355	<b>0.327</b>

are introduced (denoted by Original) and the dataset where tuples with missing values are removed directly (denoted by Removal). We can see from Table 5 that there is not much of a difference between DDGANI and Original [except Tic-Tac-Toe](#). Particularly, we can achieve the same training results as the original Hospital dataset without missing values. Additionally, there is a noticeable improvement compared to the approach of not using any imputation algorithm (There are no complete tuples available in the Spam dataset for training the classifier after injecting 20% missing values).

## 8.6 Effectiveness of Data Dependency Plug-in

**Comparison with Holistic.** We inject 10% and 50% missing values into Hospital and Tax datasets under MCAR mechanisms to compare two ways of utilizing FDs for missing value imputation. We choose these two datasets since they contain more FDs. Since Holistic [13] directly utilizes the FDs to guide the imputation of missing values and cannot deal with the numerical attributes, we compare with it only regarding the attributes involved in FDs. Table 7 shows the ARMSE of imputation results. Here, “T”, “P”, and “R” denote the imputation of missing values with true FDs, potential FDs mined from missing data, and refined FDs, respectively. Table 7 shows that DDGANI with refined FDs performs better than Holistic employing true FDs in most cases. Moreover, the superiority of DDGANI becomes even more obvious as the missing rate increases. This may be attributed to DDGANI’s capability not only to refer to the attributes involved in FDs but also to those outside the FDs when imputing the missing values. Additionally, it is noteworthy that both methods have better performance with refined FDs than initially obtained FDs, which reveals the role of FD refinement.

**An Ablation Study.** We also conduct an ablation study to verify the impact of data dependency plug-in ( $\mathcal{L}_F$ ) on the model performance, which is shown in Table 6. [Note that the removal of  \$\mathcal{L}\_F\$  has no impact on Wine and Flare since they have no FDs. But when dealing with a dataset like Tax that contains a large number of FDs,  \$\mathcal{L}\_F\$](#)

significantly boosts model performance. We can see from Table 6 that the removal of  $\mathcal{L}_F$  makes ARMSE increase from 0.349 to 0.366 and AMAE increase from 0.325 to 0.345 in Tax. It confirms the efficacy of  $\mathcal{L}_F$  in preserving the data dependencies of imputed data.

**Details of FDs Refinement.** Here, we provide the details of the FDs refinement in the three datasets Adult, Hospital and Tax. Table 8 demonstrates the true FDs, the potential FDs mined from the tabular data with missing values, the refined FDs, and the number of refinement iterations needed. The red and gray FDs in Table 8 denote the erroneous FDs and the correct FDs missing after the FDs refinement, respectively.

First of all, we can notice that the number of potential FDs mined from the incomplete tabular data is not particularly large. The total number of true FDs in these three datasets is 19, while the number of candidate FDs mined is 30, which is almost 1.5 times the number of true FDs. The reason for this is we prune some of the candidates during the mining process and set a threshold to limit the maximum number of attributes probed on the left-hand side of an FD.

Table 8 further reveals that it is hard to avoid mining erroneous FDs from the tabular data with missing values, but the vast majority of them can be fixed correctly. We can see that we mine a total of 15 inaccurate FDs from the three datasets, of which 12 are correctly refined. Moreover, the refinement does not require too many iterations, which can prevent the erroneous FDs from having a continuous negative effect on the missing value imputation. However, there are three erroneous FDs that remain after the phase of refinement in Tax. This may be due to the strong correlation between the attributes in these FDs (e.g., State→Areacode) or the values of attributes on the left-hand side being too discrete (e.g., SingleExemp, MarriedExemp→Areacode). These three FDs hold true in most tuples, making it difficult to learn to impute the missing values that do not follow such patterns. It can also be seen that a correct FD in Hospital is mined in the beginning but mistakenly discarded in the refinement phase. The reason is that the imputation is not accurate enough, which makes the confidence of the tuple violating this FD slightly higher than the confidence of FD.

## 8.7 Impact of Hyperparameters

We vary the diffusion timestep  $T$  from 1 to 12 and the threshold  $\kappa$  in the attention mechanism from 0.05 to 1 to show their impact on the performance of missing value imputation. The results are shown in Figure 5.

As we can see from the three figures on the left in Figure 5, the imputation error decreases dramatically as  $T$  grows (The ARMSE decreases from 0.220 at  $T = 1$  to 0.195 at  $T = 6$  in Adult) but remains essentially stable when  $T$  exceeds a certain value, e.g.,  $T = 6$ . The results confirm the effectiveness of decomposing the generation process from a one-shot manner into multiple diffusion steps and the role of denoising diffusion GANs in shortening the diffusion process. As for the attention threshold  $\kappa$ , we observe a decline in imputation performance with increasing  $\kappa$ , which suggests that a smaller value of  $\kappa$  allows the model to focus more on the data associated with the missing values, effectively preventing decision biases caused by data imbalances. Due to the presence of the Learner, the accuracy of the downstream classifier varies within a certain controlled variance, less affected by these two hyperparameters.

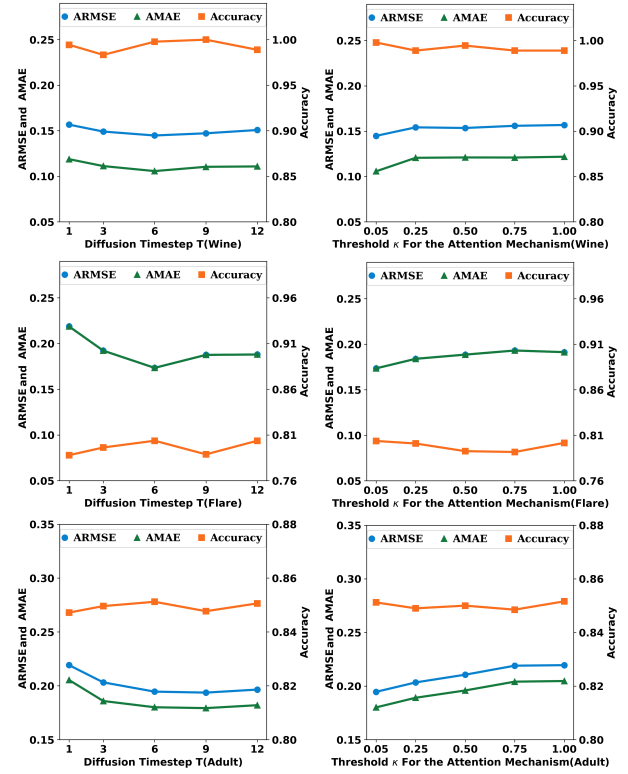


Figure 5: Impact of Hyperparameters

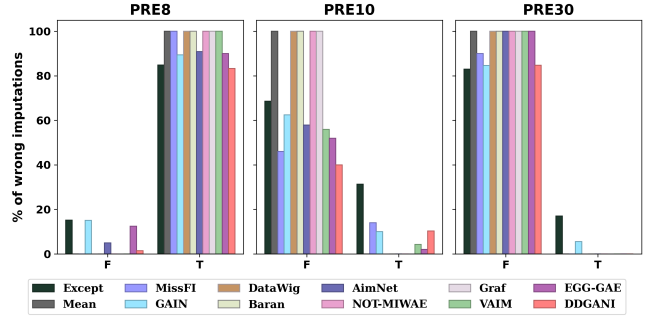


Figure 6: Capability to Handle Rare Values

## 8.8 Capability to Handle Rare Values

In order to evaluate the ability of DDGANI to deal with rare values, we have conducted an experiment on the Thoracic dataset based on the given work [8] with 20% missing rate under the MCAR model. The experimental results are shown in Figure 6. Here, the horizontal axis represents the attribute values, and there are only two attribute values “F” and “T” in the three attributes “PRE8”, “PRE10” and “PRE30”. The vertical axis represents the error rate of imputation for each method (For example, if ten cells need to be imputed with “F” but only one is imputed correctly, the error rate for “F” is 90%). The error rate of Except is  $1 - f$  where  $f$  is the proportion of the corresponding attribute value in the complete ground-truth table. The higher the frequency of a value, the higher the probability it should be correctly imputed.

**Table 8: Details of FDs Refinement**

Dataset	True FDs	Potential FDs	Refined FDs	# Iterations
Adult	education-num→education education→education-num	education-num→education education→education-num <b>marital-status, sex→relationship</b>	education-num→education education→education-num	1
Hospital	HospitalName→City City→State HospitalName→State HospitalName→HospitalOwner MeasureCode→Condition MeasureName→Condition MeasureName→MeasureCode MeasureCode→MeasureName	HospitalName→City <b>MeasureCode→State</b> <b>MeasureName→State</b> City→State HospitalName→State <b>Condition, HospitalOwner→State</b> HospitalName→HospitalOwner <b>City, MeasureCode→HospitalOwner</b> <b>City, MeasureName→HospitalOwner</b> MeasureCode→Condition MeasureName→Condition MeasureName→MeasureCode MeasureCode→MeasureName	HospitalName→City (missing) City→State HospitalName→State HospitalName→HospitalOwner MeasureCode→Condition MeasureName→Condition MeasureName→MeasureCode MeasureCode→MeasureName	4
Tax	Areacode→State State, MaritalStatus→SingleExemp State, MarriedExemp→SingleExemp Areacode, MaritalStatus→SingleExemp Areacode, MarriedExemp→SingleExemp State, MaritalStatus→MarriedExemp State, SingleExemp→MarriedExemp Areacode, MaritalStatus→MarriedExemp Areacode, SingleExemp→MarriedExemp	<b>State→Areacode</b> <b>SingleExemp, MarriedExemp→Areacode</b> Areacode→State <b>City, MarriedExemp→State</b> <b>SingleExemp, MarriedExemp→State</b> <b>Areacode→SingleExemp</b> <b>State→SingleExemp</b> <b>City, MaritalStatus→SingleExemp</b> State, MaritalStatus→MarriedExemp State, SingleExemp→MarriedExemp State, SingleExemp→MarriedExemp Areacode, MaritalStatus→MarriedExemp Areacode, SingleExemp→MarriedExemp <b>City, MaritalStatus→MarriedExemp</b> <b>City, SingleExemp→MarriedExemp</b>	<b>State→Areacode</b> <b>SingleExemp, MarriedExemp→Areacode</b> Areacode→State <b>SingleExemp, MarriedExemp→State</b> Areacode, MaritalStatus→SingleExemp Areacode, MarriedExemp→SingleExemp State, MaritalStatus→SingleExemp State, MarriedExemp→SingleExemp State, MaritalStatus→MarriedExemp State, SingleExemp→MarriedExemp Areacode, MaritalStatus→MarriedExemp Areacode, SingleExemp→MarriedExemp	2

The results reveal that all methods, including DDGANI, have difficulty restoring to rare values. This issue likely arises since imputation models favor imputing missing cells with values occurring more frequently, which can lead to lower loss during the training phase. But DDGANI can achieve higher accuracy than other methods. This improved performance is likely attributed to applying the attention mechanism, which implements the action of selectively concentrating on the tuples with similar rare data in the same attribute for imputation. But if a value appears only once, there is really no chance for DDGANI to recover it, which is one limitation of our proposed method.

## 8.9 Runtime Analysis

In the end, we conduct experiments to compare the runtime of different methods and the results can be found in Figure 7. Since most of the baselines are based on DL models, the time we report is the runtime with a GPU. Without the GPU, many methods would take one to two days to achieve convergence.

We can see from Figure 7 that since we employ a more complex model to achieve a better imputation result, the running time of our method is not dominant. DDGANI significantly takes longer than other methods in the Hospital and Tax datasets. This is due to the need for the data dependency plugin to maintain and fix the dependencies between different attributes. But in other datasets, DDGANI does not need to take too much longer than other methods. For example, in the Spam, Tic-Tac-Toe, Adult, Thoracic, and Australian

datasets, DDGANI’s runtime is lower than that of DataWig, Baran, AimNet, Graf, and VAIM in many cases. This is because DDGANI adopts a diffusion model that approaches the target distribution step by step, with each step requiring only small distributional shifts. This makes the task relatively simple and stable, resulting in faster convergence than other methods.

Despite the slower runtime, DDGANI achieves better imputation results than other methods. Thus, DDGANI is better suited for scenarios where the accuracy of imputation, especially the accuracy of the downstream classifier, is desired.

## 9 CONCLUSION

In this paper, we have proposed DDGANI, a missing value imputation framework based on denoising diffusion GAN. We have incorporated the self-attention mechanism, which enables the generative model to capture more relevant contextual information. We have designed two plugins that can be chosen based on data characteristics and downstream tasks, which ensure the utility of imputed data and the preservation of inherent data dependencies. Extensive experimental results have demonstrated the approach outperforms state-of-the-art approaches.

## ACKNOWLEDGMENTS

This work was supported by the National Natural Science Foundation of China (62372034, 61902017).

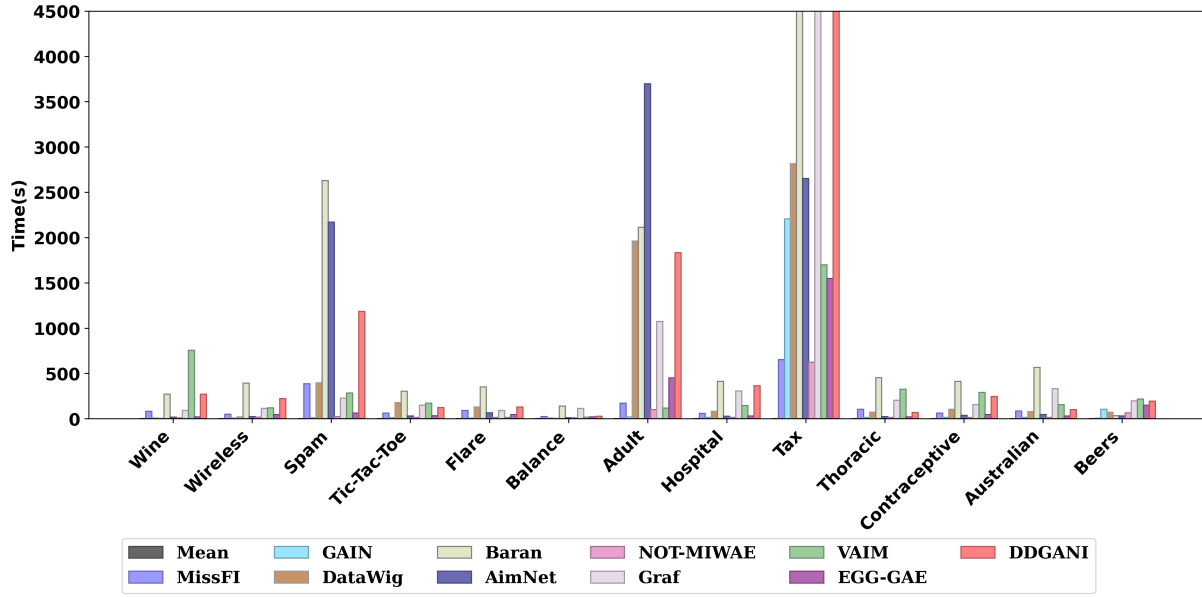


Figure 7: Comparison of Execution Time

## REFERENCES

- [1] Naomi S Altman. 1992. An introduction to kernel and nearest-neighbor non-parametric regression. *The American Statistician* 46, 3 (1992), 175–185.
- [2] Olawale F Ayilara, Lisa Zhang, Tolulope T Sajobi, Richard Sawatzky, Eric Bohm, and Lisa M Lix. 2019. Impact of missing data on bias and precision when estimating change in patient-reported outcomes from a clinical registry. *Health and quality of life outcomes* 17 (2019), 1–9.
- [3] Laure Berti-Équille, Hazar Harmouch, Felix Naumann, Noël Novelli, and Saravan Thirumuruganathan. 2018. Discovery of genuine functional dependencies from relational data with missing values. *Proceedings of the VLDB Endowment* 11, 8 (2018), 880–892.
- [4] Felix Biessmann, Tammo Rukat, Philipp Schmidt, Prathik Naidu, Sebastian Schelter, Andrey Taptunov, Dustin Lange, and David Salinas. 2019. DataWig: Missing value imputation for tables. *Journal of Machine Learning Research* 20, 175 (2019), 1–6.
- [5] Philip Bohannon, Wenfei Fan, Floris Geerts, Xibei Jia, and Anastasios Kementsitsidis. 2006. Conditional functional dependencies for data cleaning. In *2007 IEEE 23rd international conference on data engineering*. IEEE, 746–755.
- [6] Bernardo Breve, Loredana Caruccio, Vincenzo Deufemia, and Giuseppe Polese. 2022. RENUVER: A Missing Value Imputation Algorithm based on Relaxed Functional Dependencies. In *EDBT/ICDT 2024*, 27th International Conference on Extending Database Technology.
- [7] S. Van Buuren, J. P.L. Brand, C.G.M. Groothuis-Oudshoorn, and D. B. Rubin. 2006. Fully conditional specification in multivariate imputation. *Journal of Statistical Computation and Simulation* 76, 12 (2006), 1049–1064.
- [8] Riccardo Cappuzzo, Saravan Thirumuruganathan, and Paolo Papotti. 2024. Relational Data Imputation with Graph Neural Networks. In *EDBT/ICDT 2024*, 27th International Conference on Extending Database Technology.
- [9] Loredana Caruccio, Vincenzo Deufemia, and Giuseppe Polese. 2015. Relaxed functional dependencies—a survey of approaches. *IEEE Transactions on Knowledge and Data Engineering* 28, 1 (2015), 147–165.
- [10] Nadiia Chepurko, Ryan Marcus, Emanuel Zraggen, Raul Castro Fernandez, Tim Kraska, and David Karger. 2020. ARDA: automatic relational data augmentation for machine learning. *Proceedings of the VLDB Endowment* 13, 9 (2020), 1373–1387.
- [11] Fei Chiang and Renée J Miller. 2011. A unified model for data and constraint repair. In *2011 IEEE 27th International Conference on Data Engineering*. IEEE, 446–457.
- [12] Xu Chu, Ihab F Ilyas, and Paolo Papotti. 2013. Discovering denial constraints. *Proceedings of the VLDB Endowment* 6, 13 (2013), 1498–1509.
- [13] Xu Chu, Ihab F Ilyas, and Paolo Papotti. 2013. Holistic data cleaning: Putting violations into context. In *2013 IEEE 29th International Conference on Data Engineering (ICDE)*. IEEE, 458–469.
- [14] Harald Cramér. 1999. *Mathematical methods of statistics*. Vol. 43. Princeton university press.
- [15] Xiang Deng, Huan Sun, Alyssa Lees, You Wu, and Cong Yu. 2022. Turl: Table understanding through representation learning. *ACM SIGMOD Record* 51, 1 (2022), 33–40.
- [16] Alireza Farhangfar, Lukasz A Kurgan, and Witold Pedrycz. 2007. A novel framework for imputation of missing values in databases. *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans* 37, 5 (2007), 692–709.
- [17] Hanan Hammad Alharbi and Masaomi Kimura. 2020. Missing Data Imputation Using Data Generated By GAN. In *2020 the 3rd International Conference on Computing and Big Data*. 73–77.
- [18] Xuerui Hong and Shuang Hao. 2023. Imputation of Missing Values in Training Data using Variational Autoencoder. In *2023 IEEE 39th International Conference on Data Engineering Workshops (ICDEW)*. IEEE, 49–54.
- [19] Yka Huhtala, Juha Kärkkäinen, Pasi Porkka, and Hannu Toivonen. 1999. TANE: An efficient algorithm for discovering functional and approximate dependencies. *The computer journal* 42, 2 (1999), 100–111.
- [20] Niels Bruun Ipsen, Pierre-Alexandre Mattei, and Jes Frellsen. 2021. not-MIWAE: Deep Generative Modelling with Missing not at Random Data. In *ICLR 2021-International Conference on Learning Representations*.
- [21] Shawn R Jeffery, Minos Garofalakis, and Michael J Franklin. 2006. Adaptive cleaning for RFID data streams. In *Vldb*, Vol. 6. Citeseer, 163–174.
- [22] Julie Josse, Jérôme Pages, and François Husson. 2011. Multiple imputation in principal component analysis. *Advances in data analysis and classification* 5 (2011), 231–246.
- [23] Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013).
- [24] Akim Kotelnikov, Dmitry Baranchuk, Ivan Rubachev, and Artem Babenko. 2023. Tabddpm: Modelling tabular data with diffusion models. In *International Conference on Machine Learning*. PMLR, 17564–17579.
- [25] Steven Cheng-Xian Li, Bo Jiang, and Benjamin Marlin. 2018. MisGAN: Learning from Incomplete Data with Generative Adversarial Networks. In *International Conference on Learning Representations*.
- [26] Ester Livshits, Alireza Heidari, Ihab F Ilyas, and Benny Kimelfeld. 2020. Approximate denial constraints. *arXiv preprint arXiv:2005.08540* (2020).
- [27] Mohammad Mahdavi and Ziawasch Abedjan. 2020. Baran: Effective error correction via a unified context representation and transfer learning. *Proceedings of the VLDB Endowment* 13, 12 (2020), 1948–1961.
- [28] Mohammad Mahdavi, Ziawasch Abedjan, Raul Castro Fernandez, Samuel Maden, Mourad Ouzzani, Michael Stonebraker, and Nan Tang. 2019. Raha: A configuration-free error detection system. In *Proceedings of the 2019 International Conference on Management of Data*. 865–882.
- [29] Rahul Mazumder, Trevor Hastie, and Robert Tibshirani. 2010. Spectral regularization algorithms for learning large incomplete matrices. *The Journal of Machine Learning Research* 11 (2010), 2287–2322.
- [30] John T McCoy, Steve Kroon, and Lidia Auret. 2018. Variational autoencoders for missing data imputation with application to a simulated milling circuit. *IFAC-PapersOnLine* 51, 21 (2018), 141–146.
- [31] Xiaoye Miao, Yangyang Wu, Lu Chen, Yunjun Gao, Jun Wang, and Jianwei Yin. 2021. Efficient and effective data imputation with influence functions. *Proceedings of the VLDB Endowment* 15, 3 (2021), 624–632.



- [32] Robin Mitra, Sarah F McGough, Tapabrata Chakraborti, Chris Holmes, Ryan Copping, Niels Hagenbuch, Stefanie Biedermann, Jack Noonan, Brieuc Lehmann, Aditi Shenvi, et al. 2023. Learning from data with structured missingness. *Nature Machine Intelligence* 5, 1 (2023), 13–23.
- [33] Alfredo Nazabal, Pablo M Olmos, Zoubin Ghahramani, and Isabel Valera. 2020. Handling incomplete heterogeneous data using vaes. *Pattern Recognition* 107 (2020), 107501.
- [34] Jinfeng Peng, Derong Shen, Nan Tang, Tieying Liu, Yue Kou, Tiezheng Nie, Hang Cui, and Ge Yu. 2022. Self-supervised and Interpretable Data Cleaning with Sequence Generative Adversarial Networks. *Proceedings of the VLDB Endowment* 16, 3 (2022), 433–446.
- [35] Theodoros Rekatsinas, Xu Chu, Ihab F Ilyas, and Christopher Ré. 2017. Holoclean: Holistic data repairs with probabilistic inference. *arXiv preprint arXiv:1702.00820* (2017).
- [36] Patrick Royston and Ian R White. 2011. Multiple imputation by chained equations (MICE): implementation in Stata. *Journal of statistical software* 45 (2011), 1–20.
- [37] Tolou Shadbahr, Michael Roberts, Jan Stanczuk, Julian Gilbey, Philip Teare, Sören Dittmer, Matthew Thorpe, Ramon Viñas Torné, Evis Sala, Pietro Lió, et al. 2023. The impact of imputation quality on machine learning classifiers for datasets with missing values. *Communications Medicine* 3, 1 (2023), 139.
- [38] Shaoxu Song and Lei Chen. 2011. Differential dependencies: Reasoning and discovery. *ACM Transactions on Database Systems (TODS)* 36, 3 (2011), 1–41.
- [39] Shaoxu Song, Yu Sun, Aoqian Zhang, Lei Chen, and Jianmin Wang. 2018. Enriching data imputation under similarity rule constraints. *IEEE transactions on knowledge and data engineering* 32, 2 (2018), 275–287.
- [40] Indro Spinelli, Simone Scardapane, and Aurelio Uncini. 2020. Missing data imputation with adversarially-trained graph convolutional networks. *Neural Networks* 129 (2020), 249–260.
- [41] Daniel J Stekhoven and Peter Bühlmann. 2012. MissForest—non-parametric missing value imputation for mixed-type data. *Bioinformatics* 28, 1 (2012), 112–118.
- [42] Nan Tang, Ju Fan, Fangyi Li, Jianhong Tu, Xiaoyong Du, Guoliang Li, Sam Madden, and Mourad Ouzzani. 2021. RPT: relational pre-trained transformer is almost all you need towards democratizing data preparation. *Proceedings of the VLDB Endowment* 14, 8 (2021), 1254–1261.
- [43] Lev Telyatnikov and Simone Scardapane. 2023. EGG-GAE: scalable graph neural networks for tabular data imputation. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 2661–2676.
- [44] Olga Troyanskaya, Michael Cantor, Gavin Sherlock, Pat Brown, Trevor Hastie, Robert Tibshirani, David Botstein, and Russ B Altman. 2001. Missing value estimation methods for DNA microarrays. *Bioinformatics* 17, 6 (2001), 520–525.
- [45] Stef van Buuren. 2007. Multiple imputation of discrete and continuous data by fully conditional specification. *Statistical Methods in Medical Research* 16, 3 (2007), 219–242. <https://doi.org/10.1177/0962280206074463>
- [46] Richard Wu, Aoqian Zhang, Ihab Ilyas, and Theodoros Rekatsinas. 2020. Attention-based learning for missing data imputation in HoloClean. *Proceedings of Machine Learning and Systems* 2 (2020), 307–325.
- [47] Zhisheng Xiao, Karsten Kreis, and Arash Vahdat. 2021. Tackling the generative learning trilemma with denoising diffusion gans. *arXiv preprint arXiv:2112.07804* (2021).
- [48] Jinsung Yoon, James Jordon, and Mihaela Schaar. 2018. Gain: Missing data imputation using generative adversarial nets. In *International Conference on Machine Learning*. PMLR, 5689–5698.