

---

# SUPPLEMENTARY MATERIAL: MASKED LANGUAGE MODELING BECOMES CONDITIONAL DENSITY ESTIMATION FOR TABULAR DATA SYNTHESIS

---

## A Appendix / supplemental material

### A.1 Proof of Proposition 1

*Proof.* Without loss of generality, let  $\sigma$  be an identity permutation. We equally space grid points on the interval  $[0, 1]$  such that bins are represented as  $[b_{l-1}, b_l], l = 1, \dots, L$ , where  $h$  denotes the bin width such that  $Lh = 1$ . We consider arbitrary  $j \in I_C$  and  $\mathbf{x} \in \mathbb{R}^p$ , where  $\mathbf{x}$  is given and fixed. For notational simplicity, let  $\mathbf{x}_{-j} := (\mathbf{x}_1, \dots, \mathbf{x}_{j-1})$ .

Suppose that  $\hat{F}_j(\mathbf{x}_j) \in [b_{l-1}, b_l]$ . By the mean value theorem, there exists  $u^* \in (b_{l-1}, b_l)$  such that

$$h \cdot c_j^*(u^* | \mathbf{x}_{-j}) = \int_{b_{l-1}}^{b_l} c_j^*(v | \mathbf{x}_{-j}) dv = \pi_{jl}^*(\mathbf{x}_{-j}).$$

And we have

$$\begin{aligned} & \left| c_j^*(\hat{F}(\mathbf{x}_j) | \mathbf{x}_{-j}) - \hat{c}_j(\hat{F}(\mathbf{x}_j) | \mathbf{x}_{-j}; \theta) \right| \\ = & \left| c_j^*(\hat{F}(\mathbf{x}_j) | \mathbf{x}_{-j}) - \frac{\pi_{jl}(g(\mathbf{x}; F) \odot \mathbf{m}^{(j)}; \theta)}{h} \right| \\ \leq & \left| c_j^*(\hat{F}(\mathbf{x}_j) | \mathbf{x}_{-j}) - c_j^*(u^* | \mathbf{x}_{-j}) \right| + \left| c_j^*(u^* | \mathbf{x}_{-j}) - \frac{\pi_{jl}(g(\mathbf{x}; F) \odot \mathbf{m}^{(j)}; \theta)}{h} \right| \\ \leq & |\hat{F}(\mathbf{x}_j) - u^*| \left| \frac{c_j^*(\hat{F}(\mathbf{x}_j) | \mathbf{x}_{-j}) - c_j^*(u^* | \mathbf{x}_{-j})}{\hat{F}(\mathbf{x}_j) - u^*} \right| + \left| \frac{\pi_{jl}^*(\mathbf{x}_{-j})}{h} - \frac{\pi_{jl}(g(\mathbf{x}; F) \odot \mathbf{m}^{(j)}; \theta)}{h} \right| \\ \leq & hK_j + \frac{1}{h} \left| \pi_{jl}^*(\mathbf{x}_{-j}) - \pi_{jl}(g(\mathbf{x}; F) \odot \mathbf{m}^{(j)}; \theta) \right| \quad (\text{by Assumption 2}) \\ \leq & hK_j + \frac{1}{h} \sum_{l=1}^L \left| \pi_{jl}^*(\mathbf{x}_{-j}) - \pi_{jl}(g(\mathbf{x}; F) \odot \mathbf{m}^{(j)}; \theta) \right|. \end{aligned}$$

where the first equality follows from the definition of  $\hat{c}_j$ .

By the Pinsker's inequality,

$$\begin{aligned} & \frac{1}{2} \sum_{l=1}^L \left| \pi_{jl}^*(\mathbf{x}_{-j}) - \pi_{jl}(g(\mathbf{x}; F) \odot \mathbf{m}^{(j)}; \theta) \right| \\ \leq & \frac{1}{\sqrt{2}} \left( \sum_{l=1}^L \pi_{jl}^*(\mathbf{x}_{-j}) \log \pi_{jl}^*(\mathbf{x}_{-j}) - \sum_{l=1}^L \pi_{jl}^*(\mathbf{x}_{-j}) \log \pi_{jl}(g(\mathbf{x}; F) \odot \mathbf{m}^{(j)}; \theta) \right)^{1/2} \\ \leq & \frac{1}{\sqrt{2}} \left( \sum_{l=1}^L \pi_{jl}^*(\mathbf{x}_{-j}) \log \pi_{jl}^*(\mathbf{x}_{-j}) - \mathbb{E}_{\mathbf{y}_j | \mathbf{x}_{-j}} \left[ \sum_{l=1}^L \mathbb{I}(\mathbf{y}_j = l) \log \pi_{jl}(g(\mathbf{x}; F) \odot \mathbf{m}^{(j)}; \theta) \right] \right)^{1/2}, \end{aligned}$$

where  $\mathbf{y}_j|\mathbf{x}_{-j}$  is a random variable having a categorical distribution such that  $\Pr(\mathbf{y}_j = l|\mathbf{x}_{-j}) = \Pr(g(\mathbf{x}; \hat{F})_j = l|\mathbf{x}_{-j}) = \pi_{jl}^*(\mathbf{x}_{-j})$  for all  $l \in [L]$ .

Then,

$$\begin{aligned} & \left| c_j^*(\hat{F}(\mathbf{x}_j)|\mathbf{x}_{-j}) - \hat{c}_j(\hat{F}(\mathbf{x}_j)|\mathbf{x}_{-j}; \theta) \right| \\ & \leq hK_j + \frac{\sqrt{2}}{h} \left( \sum_{l=1}^L \pi_{jl}^*(\mathbf{x}_{-j}) \log \pi_{jl}^*(\mathbf{x}_{-j}) - \mathbb{E}_{\mathbf{y}_j|\mathbf{x}_{-j}} \left[ \sum_{l=1}^L \mathbb{I}(\mathbf{y}_j = l) \log \pi_{jl}(g(\mathbf{x}; F) \odot \mathbf{m}^{(j)}; \theta) \right] \right)^{1/2}. \quad (1) \end{aligned}$$

The total variation distance between  $p_j^*(\cdot|\mathbf{x}_{-j})$  and  $\hat{p}_j(\cdot|\mathbf{x}_{-j}; \theta)$  is written as

$$\begin{aligned} & \text{TV}\left(p_j^*(\cdot|\mathbf{x}_{-j}), \hat{p}_j(\cdot|\mathbf{x}_{-j}; \theta)\right) \\ & = \frac{1}{2} \int_{\mathbb{R}} \left| p_j^*(\mathbf{x}_j|\mathbf{x}_{-j}) - \hat{p}_j(\mathbf{x}_j|\mathbf{x}_{-j}; \theta) \right| d\mathbf{x}_j \\ & = \frac{1}{2} \int_{\mathbb{R}} \left| c_j^*(\hat{F}(\mathbf{x}_j)|\mathbf{x}_{-j}) - \hat{c}_j(\hat{F}(\mathbf{x}_j)|\mathbf{x}_{-j}; \theta) \right| \cdot \hat{p}_j(\mathbf{x}_j) d\mathbf{x}_j \quad (\text{by Assumption 1}) \\ & \leq \frac{hK_j}{2} + \frac{1}{\sqrt{2}h} \left( \sum_{l=1}^L \pi_{jl}^*(\mathbf{x}_{-j}) \log \pi_{jl}^*(\mathbf{x}_{-j}) - \mathbb{E}_{\mathbf{y}_j|\mathbf{x}_{-j}} \left[ \sum_{l=1}^L \mathbb{I}(\mathbf{y}_j = l) \log \pi_{jl}(g(\mathbf{x}; F) \odot \mathbf{m}^{(j)}; \theta) \right] \right)^{1/2} \\ & = \frac{hK_j}{2} + \frac{\sqrt{\text{Bias}(\theta)}}{\sqrt{2}h} \\ & = \frac{K_j}{2L} + \frac{\sqrt{\text{Bias}(\theta)}}{\sqrt{2}/L}, \end{aligned}$$

where the inequality holds for (1).

The proof is complete. □

## A.2 Proof of Proposition 2

*Proof.* Since the data is MAR,  $p(\mathbf{r}|\mathbf{x}) = p(\mathbf{r}|\mathbf{x}_{obs})$  and

$$\begin{aligned} p(\mathbf{m}, \mathbf{r}|\mathbf{x}) &= p(\mathbf{r}|\mathbf{x}) \cdot p(\mathbf{m}|\mathbf{x}) = p(\mathbf{r}|\mathbf{x}) \cdot p(\mathbf{m}) \\ &= p(\mathbf{r}|\mathbf{x}_{obs}) \cdot p(\mathbf{m}) = p(\mathbf{m}, \mathbf{r}|\mathbf{x}_{obs}). \end{aligned}$$

The proof is complete. □

### A.3 Dataset Descriptions

#### Download links.

- abalone [23]: <https://archive.ics.uci.edu/dataset/1/abalone>
- banknote [18]: <https://archive.ics.uci.edu/dataset/267/banknote+authentication>
- breast [32]: <https://archive.ics.uci.edu/dataset/17/breast+cancer+wisconsin+diagnostic>
- concrete [34]: <https://archive.ics.uci.edu/dataset/165/concrete+compressive+strength>
- covertype [2]: <https://www.kaggle.com/datasets/uciml/forest-cover-type-dataset>
- kings (CC0: Public Domain): <https://www.kaggle.com/datasets/harlfoxem/housesalesprediction>
- letter [28]: <https://archive.ics.uci.edu/dataset/59/letter+recognition>
- loan (CC0: Public Domain): <https://www.kaggle.com/datasets/teertha/personal-loan-modeling>
- redwine [6]: <https://archive.ics.uci.edu/dataset/186/wine+quality>
- whitewine [6]: <https://archive.ics.uci.edu/dataset/186/wine+quality>

Dataset	Train/Test Split	#continuous	#categorical	Classification Target
abalone	3.3K/0.8K	7	2	Rings
banknote	1.1K/0.3K	4	1	class
breast	0.5K/0.1K	30	1	Diagnosis
concrete	0.8K/0.2K	8	1	Age
covtype	0.46M/0.12M	10	1	Cover_Type
kings	17.3K/4.3K	11	7	grade
letter	16K/4K	16	1	lettr
loan	4K/1K	5	6	Personal Loan
redwine	1.3K/0.3K	11	1	quality
whitewine	3.9K/1K	11	1	quality

Table 1: **Description of datasets.** #continuous represents the number of continuous and ordinal variables. #categorical denotes the number of categorical variables. The ‘Classification Target’ refers to the variable used as the response variable in a classification task to evaluate machine learning utility.

### A.4 Experimental Settings for Reproduction

- We run experiments using NVIDIA A10 GPU, and our experimental codes are available with pytorch.
- In practice, for all  $j \in I_C$ , we estimate  $\hat{F}_j$  using empirical measure. In the presence of missing data,  $\hat{F}_j$  is estimated solely using the observed dataset (refer to [4]).

**Hyper-parameters of MaCoDE:** As shown in Section A.9, MaCoDE consistently generates high-quality synthetic data without requiring an extensive hyperparameter tuning process, unlike methods such as [14, 15, 13]. This demonstrates the generalizability of our proposed model to various tabular datasets. **Our implementation codes for the proposed model, MaCoDE, are provided in the supplementary material. For all tabular datasets, we applied the following hyperparameters uniformly without any additional tuning:**

- epochs: 500
- batch size: 1024
- learning rate: 0.001 (with AdamW optimizer [19] with 0.001 weight decay parameter)
- the number of bins:  $L = 50$
- Transformer encoder dimension: 128
- Transformer encoder #heads: 4
- Transformer encoder #layer: 2
- Transformer encoder dropout ratio: 0.0

#### A.4.1 Details of Implementing Baseline Models

To compare our proposed method with baseline models, we performed experiments in a manner similar to [36]. We adjusted the hidden or latent dimensions across various methods, ensuring that the number of trainable parameters is comparable. Since each model employs distinct neural network architectures, it would be misleading to compare their performances using the same network configuration. It’s worth noting that model size typically reflects the performance of the baseline models. Under these conditions, we reproduced the baseline methods using their official codes and default configurations (except for the hidden or latent dimensions). **Our reproduced codes for the baseline models are provided in the supplementary material.** Below are the detailed implementations of the baseline methods:

- CTGAN and TVAE [33]: We follow the implementations provided in the official repository\* for CTGAN and TVAE. We adopt the default hyperparameters specified in the module. To ensure fairness in comparison by aligning model sizes, we adjust the latent dimension of CTGAN and TVAE to 100.
- CTAB-GAN and CTAB-GAN+ [40, 41]: We follow the implementations provided in the official repository†. The latent dimension is set to 100, and the maximum number of clusters is set to 10. For the specification of feature types, we use three categories: continuous, integer, and categorical variables.
- DistVAE [1]: We follow the implementations provided in the official repository‡. We changed the latent dimension to 100 to increase the model size.
- TabDDPM [14]: We utilized the TabDDPM module in `synthcity.Plugins`§ for synthetic data generation. This module follows the implementations provided in the official repository.
- TabMT [8]: Rather than using K-means clustering, we utilize the Gaussian Mixture Model (GMM) to discretize continuous columns to preserve the original continuous domain. We determine the optimal number of clusters for the GMM, ranging from 2 to 10, based on the Bayesian Information Criterion (BIC). During the generation of continuous columns, we initially predict the component label and then sample from the selected Gaussian component.
- MICE [31]: We employed the `IterativeImputer` package from `Scikit-learn` for multiple imputation using chained equations. Following the authors’ experiments, a `max_iter` range of 10 to 20 was considered sufficient for convergence, and we adopted this setting for our experiments. Additionally, to introduce randomness, we set `imputation_order` to `random`, which randomly selects a variable for imputation in each iteration. The remaining parameters were left at their default values to maintain the integrity of the MICE implementation.
- GAIN [35]: We follow the implementations provided in the official repository¶. As the paper does not explicitly discuss the separate handling of categorical and continuous variables, the code treats them simultaneously. Consequently, a rounding process is employed to handle categorical variables afterward.
- missMDA [12]: We utilized the `missMDA` package¶ in R for multiple imputation.
- VAEAC [10]: We follow the implementations provided in the official repository\*\*. The authors provided hyperparameters that adequately address both continuous and categorical variables, so we used these without further modification during model fitting.
- MIWAE [20]: The implemented MIWAE code in the official repository†† was designed for continuous variables only. To accommodate heterogeneous tabular datasets, we treated the conditional distribution of categorical columns as categorical distributions and employed cross-entropy loss for reconstruction. For comparison with not-MIWAE, we set the latent dimension to  $p - 1$ .
- not-MIWAE [9]: The implemented not-MIWAE code in the official repository‡‡ also focused on continuous variables exclusively. To handle categorical variables, we made the same modifications as in MIWAE. For comparison with MIWAE, we set the latent dimension to  $p - 1$ . Training was conducted for 100K steps, consistent with the official implementations.

---

\*<https://github.com/sdv-dev/CTGAN>

†<https://github.com/Team-TUD/CTAB-GAN>, <https://github.com/Team-TUD/CTAB-GAN-Plus>

‡<https://github.com/an-seunghwan/DistVAE>

§<https://github.com/vanderschaarlab/synthcity>

¶<https://github.com/jsyoon0823/GAIN/tree/master>

¶<https://cran.r-project.org/web/packages/missMDA/index.html>

\*\*<https://github.com/tigvarts/vaeac>

††<https://github.com/pamattei/miwae>

‡‡<https://github.com/nbip/notMIWAE>

- EGC [39]: We utilized the `gcimpute` package<sup>§§</sup> to implement EGC. The model is free of hyperparameters.

Tasks	Model	Description
Regression	Random Forest	Package: <code>sklearn.ensemble.RandomForestRegressor</code> , setting: <code>random_state=0</code> , defaulted values
	Logistic Regression	Package: <code>sklearn.linear_model.LogisticRegression</code> , setting: <code>random_state=0</code> , <code>max_iter=1000</code> , defaulted values
Classification	Gaussian Navie Bayes	Package: <code>sklearn.naive_bayes.GaussianNB</code> , setting: defaulted values
	K-Nearest Neighbors	Package: <code>sklearn.neighbors.KNeighborsClassifier</code> , setting: defaulted values
	Decision Tree	Package: <code>sklearn.tree.DecisionTreeClassifier</code> , setting: <code>random_state=0</code> , defaulted values
	Random Forest	Package: <code>sklearn.ensemble.RandomForestClassifier</code> , setting: <code>random_state=0</code> , defaulted values

Table 2: **Regressor and classifier used to evaluate synthetic data quality in machine learning utility.** The names of all parameters used in the description are consistent with those defined in corresponding packages.

## A.5 Evaluation Procedure of Q1

### Regression performance (SMAPE).

1. Train a synthesizer using the real training dataset.
2. Generate a synthetic dataset with the same size as the real training dataset.
3. Train a machine learning model (Random Forest regressor) using the synthetic dataset, where each continuous column serves as the regression target variable.
4. Assess regression prediction performance by averaging the SMAPE values from the test dataset for each Random Forest regressor trained on the continuous columns.

### Classification performance ( $F_1$ ).

1. Train a synthesizer using the real training dataset.
2. Generate a synthetic dataset with the same size as the real training dataset.
3. Train machine learning models (Logistic Regression, Gaussian Naive Bayes, K-Nearest Neighbors classifier, Decision Tree classifier, and Random Forest classifier) using the synthetic dataset (see Table 1 for the classification target variable and refer to Table 2 for detailed configuration).
4. Assess classification prediction performance by averaging the  $F_1$  values from the test dataset from five different classifiers.

### Model selection performance (Model).

1. Train a synthesizer using the real training dataset.
2. Generate a synthetic dataset with the same size as the real training dataset.
3. Train machine learning models (Logistic Regression, Gaussian Naive Bayes, K-Nearest Neighbors classifier, Decision Tree classifier, and Random Forest classifier) using both the real training dataset and the synthetic dataset (see Table 1 for the classification target variable and refer to Table 2 for detailed configuration).
4. Evaluate the classification performance (AUROC) of all trained classifiers on the test dataset.
5. Assess model selection performance by comparing the AUROC rank orderings of classifiers trained on the real training dataset and those trained on the synthetic dataset using Spearman’s Rank Correlation.

<sup>§§</sup><https://github.com/udellgroup/gcimpute/tree/master>

### Feature selection performance (Feature).

1. Train a synthesizer using the real training dataset.
2. Generate a synthetic dataset with the same size as the real training dataset.
3. Train a Random Forest classifier using both the real training dataset and the synthetic dataset (see Table 1 for the classification target variable and refer to Table 2 for detailed configuration).
4. Determine the rank-ordering of important features for both classifiers.
5. Assess feature selection performance by comparing the feature importance rank orderings of classifiers trained on the real training dataset and those trained on the synthetic dataset using Spearman’s Rank Correlation.

## A.6 Evaluation Procedure of Q2

### Regression performance (SMAPE).

1. For each random seed, we generate the mask and train MaCoDE using a masked training dataset (i.e., incomplete dataset).
2. Generate a synthetic dataset with the same size as the real training dataset.
3. Train a machine learning model (Random Forest regressor) using the synthetic dataset, where each continuous column serves as the regression target variable.
4. Assess regression prediction performance by averaging the SMAPE values from the test dataset for each Random Forest regressor trained on the continuous columns.

### Classification performance ( $F_1$ ).

1. For each random seed, we generate the mask and train MaCoDE using a masked training dataset (i.e., incomplete dataset).
2. Generate a synthetic dataset with the same size as the real training dataset.
3. Train machine learning models (Logistic Regression, Gaussian Naive Bayes, K-Nearest Neighbors classifier, Decision Tree classifier, and Random Forest classifier) using the synthetic dataset (see Table 1 for the classification target variable and refer to Table 2 for detailed configuration).
4. Assess classification prediction performance by averaging the  $F_1$  values from the test dataset from five different classifiers.

## A.7 Evaluation of Q3

### A.7.1 Related Works

Data missingness is a common challenge in research and practical analysis, categorized into three primary missing mechanisms: (1) Missing Completely at Random (MCAR), (2) Missing at Random (MAR), and (3) Missing Not at Random (MNAR).

Under the **MCAR** mechanism, the reason for missingness has no relationship with any data, neither observed nor unobserved. In other words, the likelihood of data being missing is equal across all observations. The primary advantage of MCAR is that it does not introduce bias into the data analysis. However, despite this advantage, data missingness can still reduce the statistical power of the study because of the reduced sample size.

**MAR** occurs when the probability of missingness is related to the observed data but not the unobserved missing data. Essentially, even though the data is missing, the mechanism assumes that the missingness is explainable by other variables in the dataset. That is, the missingness can be modeled and imputed using the information available in the data, allowing for more accurate analyses despite the missingness.

Lastly, if the missingness is not specified by either MCAR or MAR, it becomes **MNAR**. The MNAR is the most challenging mechanism, as it implies that the missingness is related to the unobserved data itself. In this case, the missing data is systematically different from the observed data, which introduces bias if not properly accounted for. For example, patients with severe symptoms may be less likely to report their health status, making their data missing. MNAR requires sophisticated statistical methods to address, as ignoring or improperly handling it can lead to biased and unreliable results.

Recent methods employ the deep generative model to estimate a joint distribution and generate samples for imputations. GAN-based imputers such as GAIN [35] and MisGAN [17] adopt an adversarial learning approach to generate both missing entries and masking vectors. Other approaches like VAEAC [10], HI-VAE [24], and ReMasker [7] employ strategies that allow them to learn conditional distributions on arbitrary conditioning sets using the uniform masking strategy. MIWAE [20], based on the Importance Weighted Autoencoder [3], demonstrates that under mild conditions and the MAR assumption, the target likelihood can be approximated regardless of the imputation function. Additionally, not-MIWAE [9] extends MIWAE to handle cases where the data is under MNAR assumption by modeling missing entries as a latent variable. In parallel, the Optimal Transport (OT)-based method employs distributional matching, utilizing the 2-Wasserstein distance to compare distributions in both data and latent spaces [22, 38]. To handle mixed-type tabular datasets, [39] introduced EGC (extended Gaussian copula), which relies on a latent Gaussian distribution to support single and multiple imputations.

## A.7.2 Evaluation Metrics

---

**Algorithm 1** Evaluation procedure for multiple imputation [30, 27]

---

**Input:** Complete dataset  $D = \{x_i\}_{i=1}^n$

**Output:** Bias, Coverage, and Confidence interval length

- 1: Target estimand:  $Q^* = 1/n \sum_{i=1}^n \mathbb{I}(x_i > \bar{x})$ , where  $\bar{x} = 1/n \sum_{i=1}^n x_i$
  - 2: **for** (random seeds)  $s = 1, \dots, S$  **do**
  - 3:   Generate missing dataset  $D^{(s)}$
  - 4:   Training the imputation model using the dataset  $D^{(s)}$
  - 5:   Perform multiple imputation:  $\hat{D}_m^{(s)}, m = 1, \dots, M$
  - 6:   **for**  $m = 1, 2, \dots, M$  **do**
  - 7:      $\hat{Q}_m^{(s)} = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(x_i > \bar{x})$ , where  $\hat{D}_m^{(s)} = \{x_i\}_{i=1}^n$  and  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
  - 8:      $\hat{U}_m^{(s)} = \hat{Q}_m^{(s)} (1 - \hat{Q}_m^{(s)})/n$
  - 9:      $\hat{\bar{Q}}^{(s)} = \frac{1}{M} \sum_{m=1}^M \hat{Q}_m^{(s)}$
  - 10:     $\hat{U}^{(s)} = \frac{1}{M} \sum_{m=1}^M \hat{U}_m^{(s)} + (M+1)/M \cdot \frac{1}{M-1} \sum_{m=1}^M (\hat{U}_m^{(s)} - \bar{U}^{(s)})^2$ ,
  - 11: where  $\bar{U}^{(s)} = \frac{1}{M} \sum_{m=1}^M \hat{U}_m^{(s)}$
  - 12: Bias:  $\frac{1}{S} \sum_{s=1}^S |\hat{\bar{Q}}^{(s)} - Q^*|$
  - 13: Coverage:  $\frac{1}{S} \sum_{s=1}^S \mathbb{I}(Q^* \in (\hat{\bar{Q}}^{(s)} \pm 1.96 \cdot \sqrt{\hat{U}^{(s)}}))$
  - 14: Confidence interval length:  $\frac{1}{S} \sum_{s=1}^S 2 \cdot 1.96 \cdot \sqrt{\hat{U}^{(s)}}$
- 

For Q3, we selected the following multiple imputation models that can handle heterogeneous tabular datasets: MICE [31], GAIN [35], missMDA [12], VAEAC [10], MIWAE [20], not-MIWAE [9], and EGC [39].

We assess the effectiveness of multiple imputations by employing interval inference for the population mean, which was proposed by Rubin [27]. However, since obtaining the population mean is not feasible, we utilize the sample column-wise mean of the complete dataset as a parameter of interest [16, 37]. The evaluation procedure for multiple imputations is outlined in Algorithm 1, and we utilize `torch.random.manual_seed()` to set seeds in performing multiple imputations. We report the mean and standard error of bias, coverage, and confidence interval length across 10 different random seeds, all continuous columns, and datasets [30]. Note that we do not split the data into training and test sets [9].

**Remark 1.** *We acknowledge that existing imputation methods have been evaluated using RMSE (root mean square error), a metric for assessing single imputation methods [9, 20, 24, 22, 38, 11]. However, since our objective focuses on distributional learning rather than recovering missing values, we adopt the evaluation procedure proposed in [27, 30], better suited for assessing multiple imputation methods.*

Following [22, 11, 38], we generate the missing value mask for each dataset with three mechanisms in four settings. (MCAR) In the MCAR setting, each value is masked according to the realization of a Bernoulli random variable with a fixed parameter. (MAR) In the MAR setting, for each experiment, a fixed subset of variables that cannot have missing values is sampled. Then, the remaining variables have missing values according to a logistic model with random weights, which takes the non-missing variables as inputs. A bias term is fitted using line search to attain the desired proportion of missing values. (MNAR) Finally, two different mechanisms are implemented in the MNAR setting. The first, MNARL, is identical to the previously described MAR mechanism, but the inputs of the logistic model are then masked by an MCAR mechanism. Hence, the logistic model’s outcome depends on missing values. The second

mechanism, MNARQ, samples a subset of variables whose values in the lower and upper  $p$ th percentiles are masked according to a Bernoulli random variable, and the values in-between are left not missing.

Model	Bias ↓	Coverage	Width ↓
MICE	.010 $\pm$ .001	.845 $\pm$ .019	<b>.040</b> $\pm$ .002
GAIN	.019 $\pm$ .002	.633 $\pm$ .033	<b>.040</b> $\pm$ .002
missMDA	.015 $\pm$ .001	.700 $\pm$ .022	.043 $\pm$ .002
VAEAC	<u>.008</u> $\pm$ .001	.905 $\pm$ .016	<b>.040</b> $\pm$ .002
MIWAE	<b>.006</b> $\pm$ .000	<u>.952</u> $\pm$ .012	.043 $\pm$ .002
not-MIWAE	<b>.006</b> $\pm$ .000	<b>.949</b> $\pm$ .012	<u>.042</u> $\pm$ .002
EGC	<b>.006</b> $\pm$ .000	.996 $\pm$ .004	.058 $\pm$ .002
MaCoDE(MAR)	<b>.006</b> $\pm$ .000	.963 $\pm$ .009	.051 $\pm$ .003

Table 3: **Q3**: Multiple imputation under MAR at 0.3 missingness. The means and standard errors of the mean across 5 datasets and 10 repeated experiments are reported. ↓ denotes lower is better. Coverage close to 0.95 indicates better performance. The best value is bolded, and the second best is underlined.

### A.7.3 Results

The missing data mechanism within the parentheses refers to the mechanism applied to the dataset on which MaCoDE was trained. Table 3 indicates that MaCoDE consistently exhibits competitive performance against all baseline models across metrics assessing multiple imputation performances, including bias, coverage, and confidence interval length. This suggests our proposed approach can support multiple imputations for deriving statistically valid inferences from missing data with the MAR mechanism.



## A.8 Additional Experiments

### A.8.1 Privacy Preservability

- **Evaluation metrics:** The  $k$ -anonymity property [29] is a measure used to assess the level of privacy protection in synthetic data. It ensures that each individual’s information in the dataset cannot be distinguished from that of at least  $k - 1$  other individuals. In other words, each record in the dataset has at least  $k - 1$  similar records in terms of quasi-identifiers, attributes that could potentially identify a subject. A higher  $k$  value implies a higher level of anonymity and better privacy preservation.

*DCR (Distance to Closest Record)* [25, 40] is defined as the distances between all real training and synthetic samples. A higher DCR value indicates more effective privacy preservation, indicating a lack of overlap between the real training data and the synthetic samples. Conversely, an excessively large DCR score suggests a lower quality of the generated synthetic dataset. Therefore, the DCR metric provides insights into both the privacy-preserving capability and the quality of the synthetic dataset.

*Attribute disclosure* [5, 21] refers to the situation where attackers can uncover additional covariates of a record by leveraging a subset of covariates they already possess, along with similar records from the synthetic dataset. To quantify the extent to which attackers can accurately identify these additional covariates, we employ classification metrics. Higher attribute disclosure metrics indicate an increased risk of privacy leakage, implying that attackers can precisely infer unknown variables. In terms of privacy concerns, attribute disclosure can be considered a more significant issue than membership inference attacks, as attackers are assumed to have access to only a subset of covariates for a given record.

- **Evaluation procedure:** We evaluate the  $k$ -anonymity following the approach described in [26]<sup>¶</sup>. Additionally, similar to [40], we define DCR as the 5<sup>th</sup> percentile of the  $L_2$  distances between all real training samples and synthetic samples. Since DCR relies on  $L_2$  distance and continuous variables, it is computed solely using continuous variables. We assess attribute disclosure using the methodology outlined in [5].
- **Result:** The right panel of Figure 1 demonstrates MaCoDE’s capability to regulate the privacy level by adjusting the temperature parameter  $\tau$ . Simultaneously, the left panel illustrates that the quality of synthetic data, measured by feature selection performance, remains notable even with increasing privacy levels from  $\tau = 1$  to  $\tau = 3$ . However, increasing  $\tau$  beyond 3 leads to declining feature selection performance compared to other models despite DCR remaining competitive. For additional results on other metrics related to the trade-off between privacy level and synthetic data quality as  $\tau$  varies, please refer to Table 4 and Table 5.

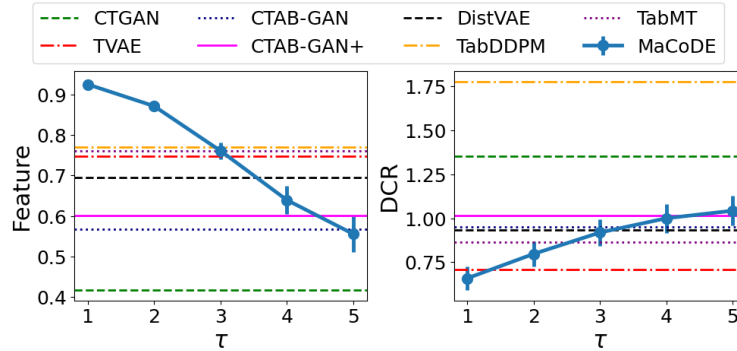


Figure 1: **Trade-off between privacy and quality.** Left: feature selection performance (synthetic data quality). Right: DCR (privacy preservability). The means and standard errors of the mean across 10 datasets and 10 repeated experiments are reported. Error bars represent the standard errors of the mean.

<sup>¶</sup><https://github.com/vanderschaarlab/synthcity>

Model	Statistical fidelity				Machine learning utility			
	KL ↓	GoF ↓	MMD ↓	WD ↓	SMAPE ↓	$F_1$ ↑	Model ↑	Feature ↑
Baseline	.016 $\pm$ .002	.029 $\pm$ .002	.002 $\pm$ .000	1.019 $\pm$ .156	.107 $\pm$ .008	.686 $\pm$ .023	.887 $\pm$ .018	.956 $\pm$ .005
CTGAN	.221 $\pm$ .014	.561 $\pm$ .046	.094 $\pm$ .007	6.435 $\pm$ 1.011	.256 $\pm$ .016	.411 $\pm$ .027	.208 $\pm$ .048	.417 $\pm$ .043
TVAE	.066 $\pm$ .003	.119 $\pm$ .005	.016 $\pm$ .001	1.631 $\pm$ .173	.192 $\pm$ .011	.608 $\pm$ .021	.486 $\pm$ .041	.747 $\pm$ .027
CTAB-GAN	.116 $\pm$ .008	.196 $\pm$ .025	.044 $\pm$ .004	3.327 $\pm$ .460	.218 $\pm$ .012	.524 $\pm$ .026	.263 $\pm$ .042	.568 $\pm$ .041
CTAB-GAN+	.136 $\pm$ .018	.144 $\pm$ .010	.054 $\pm$ .007	3.971 $\pm$ .772	.226 $\pm$ .017	.530 $\pm$ .020	.227 $\pm$ .048	.601 $\pm$ .041
DistVAE	.059 $\pm$ .007	.070 $\pm$ .004	.016 $\pm$ .001	2.272 $\pm$ .282	.226 $\pm$ .017	.588 $\pm$ .021	.194 $\pm$ .048	.695 $\pm$ .030
TabDDPM	.696 $\pm$ .117	.374 $\pm$ .087	.057 $\pm$ .011	42.916 $\pm$ 8.127	.161 $\pm$ .011	.576 $\pm$ .022	.507 $\pm$ .039	.770 $\pm$ .027
TabMT	.011 $\pm$ .001	.035 $\pm$ .003	.012 $\pm$ .001	2.299 $\pm$ .346	.188 $\pm$ .013	.622 $\pm$ .024	.528 $\pm$ .039	.761 $\pm$ .028
MaCoDE( $\tau = 1$ )	.034 $\pm$ .004	.072 $\pm$ .004	.007 $\pm$ .001	1.630 $\pm$ .245	.158 $\pm$ .010	.635 $\pm$ .023	.599 $\pm$ .035	.925 $\pm$ .007
MaCoDE( $\tau = 2$ )	.047 $\pm$ .004	.117 $\pm$ .009	.021 $\pm$ .002	2.450 $\pm$ .302	.177 $\pm$ .012	.601 $\pm$ .024	.432 $\pm$ .038	.871 $\pm$ .011
MaCoDE( $\tau = 3$ )	.071 $\pm$ .006	.290 $\pm$ .046	.035 $\pm$ .003	3.191 $\pm$ .354	.196 $\pm$ .013	.562 $\pm$ .024	.305 $\pm$ .041	.761 $\pm$ .020
MaCoDE( $\tau = 4$ )	.088 $\pm$ .008	.555 $\pm$ .109	.043 $\pm$ .003	3.601 $\pm$ .383	.210 $\pm$ .014	.528 $\pm$ .024	.128 $\pm$ .041	.639 $\pm$ .035
MaCoDE( $\tau = 5$ )	.099 $\pm$ .010	.847 $\pm$ .181	.045 $\pm$ .003	3.822 $\pm$ .402	.219 $\pm$ .014	.501 $\pm$ .025	.065 $\pm$ .047	.555 $\pm$ .044

Table 4: **Trade-off between privacy and quality** (statistical fidelity and machine learning utility). The means and standard errors of the mean across 10 datasets and 10 repeated experiments are reported. ‘Baseline’ refers to the result obtained using half of the real training dataset.  $\uparrow$  ( $\downarrow$ ) denotes higher (lower) is better.

Model	$k$ -anonymity(%) $\uparrow$	DCR $\uparrow$	AD $\downarrow$
Baseline	1.151 $\pm$ 0.089	0.501 $\pm$ 0.051	0.631 $\pm$ 0.023
CTGAN	2.413 $\pm$ 0.139	1.355 $\pm$ 0.152	0.398 $\pm$ 0.024
TVAE	1.648 $\pm$ 0.116	0.709 $\pm$ 0.056	0.583 $\pm$ 0.023
CTAB-GAN	2.048 $\pm$ 0.126	0.951 $\pm$ 0.076	0.514 $\pm$ 0.025
CTAB-GAN+	2.215 $\pm$ 0.121	1.016 $\pm$ 0.113	0.471 $\pm$ 0.018
DistVAE	2.345 $\pm$ 0.136	0.932 $\pm$ 0.071	0.540 $\pm$ 0.019
TabDDPM	1.249 $\pm$ 0.108	1.779 $\pm$ 0.391	0.567 $\pm$ 0.024
TabMT	1.656 $\pm$ 0.135	0.862 $\pm$ 0.072	0.566 $\pm$ 0.021
MaCoDE( $\tau = 1$ )	1.405 $\pm$ 0.104	0.658 $\pm$ 0.067	0.589 $\pm$ 0.022
MaCoDE( $\tau = 2$ )	1.656 $\pm$ 0.109	0.797 $\pm$ 0.070	0.547 $\pm$ 0.022
MaCoDE( $\tau = 3$ )	1.959 $\pm$ 0.147	0.919 $\pm$ 0.077	0.508 $\pm$ 0.022
MaCoDE( $\tau = 4$ )	2.075 $\pm$ 0.159	1.000 $\pm$ 0.082	0.476 $\pm$ 0.022
MaCoDE( $\tau = 5$ )	2.262 $\pm$ 0.166	1.043 $\pm$ 0.085	0.457 $\pm$ 0.022

Table 5: **Privacy preservability**. The means and standard errors of the mean across 10 datasets and 10 repeated experiments are reported. ‘Baseline’ refers to the result obtained using half of the real training dataset.  $\uparrow$  ( $\downarrow$ ) denotes higher (lower) is better.

abalone			
Dataset	$k$ -anonymity(%) $\uparrow$	DCR $\uparrow$	AD $\downarrow$
CTGAN	3.634 $\pm$ 0.318	0.423 $\pm$ 0.012	0.346 $\pm$ 0.009
TVAE	2.029 $\pm$ 0.234	0.333 $\pm$ 0.010	0.385 $\pm$ 0.018
CTAB-GAN	2.769 $\pm$ 0.315	0.449 $\pm$ 0.027	0.312 $\pm$ 0.023
CTAB-GAN+	2.718 $\pm$ 0.285	0.311 $\pm$ 0.010	0.323 $\pm$ 0.010
DistVAE	2.074 $\pm$ 0.149	0.316 $\pm$ 0.006	0.349 $\pm$ 0.009
TabDDPM	0.233 $\pm$ 0.047	0.145 $\pm$ 0.003	0.371 $\pm$ 0.016
TabMT	1.595 $\pm$ 0.617	0.423 $\pm$ 0.007	0.337 $\pm$ 0.010
MaCoDE	1.182 $\pm$ 0.292	0.148 $\pm$ 0.002	0.368 $\pm$ 0.014
breast			
Dataset	$k$ -anonymity(%) $\uparrow$	DCR $\uparrow$	AD $\downarrow$
CTGAN	1.582 $\pm$ 0.383	5.536 $\pm$ 0.167	0.446 $\pm$ 0.050
TVAE	0.813 $\pm$ 0.239	2.148 $\pm$ 0.026	0.900 $\pm$ 0.042
CTAB-GAN	0.857 $\pm$ 0.198	2.807 $\pm$ 0.078	0.850 $\pm$ 0.037
CTAB-GAN+	1.473 $\pm$ 0.289	4.025 $\pm$ 0.315	0.596 $\pm$ 0.078
DistVAE	2.264 $\pm$ 0.390	2.666 $\pm$ 0.060	0.810 $\pm$ 0.036
TabDDPM	2.835 $\pm$ 0.301	12.733 $\pm$ 0.896	0.740 $\pm$ 0.089
TabMT	0.967 $\pm$ 0.164	2.718 $\pm$ 0.054	0.825 $\pm$ 0.041
MaCoDE	0.527 $\pm$ 0.140	2.264 $\pm$ 0.058	0.906 $\pm$ 0.032
covtype			
Dataset	$k$ -anonymity(%) $\uparrow$	DCR $\uparrow$	AD $\downarrow$
CTGAN	0.215 $\pm$ 0.000	0.451 $\pm$ 0.003	0.490 $\pm$ 0.023
TVAE	0.215 $\pm$ 0.000	0.415 $\pm$ 0.002	0.639 $\pm$ 0.009
CTAB-GAN	0.215 $\pm$ 0.000	0.418 $\pm$ 0.004	0.579 $\pm$ 0.004
CTAB-GAN+	0.215 $\pm$ 0.000	0.353 $\pm$ 0.001	0.645 $\pm$ 0.001
DistVAE	0.215 $\pm$ 0.000	0.566 $\pm$ 0.002	0.612 $\pm$ 0.002
TabDDPM	0.211 $\pm$ 0.002	0.313 $\pm$ 0.001	0.653 $\pm$ 0.002
TabMT	0.215 $\pm$ 0.000	0.409 $\pm$ 0.001	0.671 $\pm$ 0.003
MaCoDE	0.215 $\pm$ 0.000	0.266 $\pm$ 0.000	0.696 $\pm$ 0.001
letter			
Dataset	$k$ -anonymity(%) $\uparrow$	DCR $\uparrow$	AD $\downarrow$
CTGAN	3.520 $\pm$ 0.180	1.583 $\pm$ 0.042	0.089 $\pm$ 0.009
TVAE	2.657 $\pm$ 0.232	0.779 $\pm$ 0.009	0.274 $\pm$ 0.008
CTAB-GAN	3.343 $\pm$ 0.267	1.438 $\pm$ 0.016	0.100 $\pm$ 0.006
CTAB-GAN+	3.112 $\pm$ 0.253	1.218 $\pm$ 0.008	0.207 $\pm$ 0.004
DistVAE	1.959 $\pm$ 0.265	1.378 $\pm$ 0.006	0.273 $\pm$ 0.005
TabDDPM	1.435 $\pm$ 0.173	0.744 $\pm$ 0.005	0.260 $\pm$ 0.008
TabMT	2.112 $\pm$ 0.236	0.822 $\pm$ 0.008	0.402 $\pm$ 0.005
MaCoDE	1.819 $\pm$ 0.089	1.174 $\pm$ 0.011	0.380 $\pm$ 0.007
redwine			
Dataset	$k$ -anonymity(%) $\uparrow$	DCR $\uparrow$	AD $\downarrow$
CTGAN	2.291 $\pm$ 0.411	1.794 $\pm$ 0.032	0.192 $\pm$ 0.034
TVAE	1.009 $\pm$ 0.325	0.976 $\pm$ 0.016	0.500 $\pm$ 0.041
CTAB-GAN	1.345 $\pm$ 0.140	1.191 $\pm$ 0.027	0.454 $\pm$ 0.027
CTAB-GAN+	2.056 $\pm$ 0.208	1.227 $\pm$ 0.027	0.448 $\pm$ 0.024
DistVAE	3.815 $\pm$ 0.220	1.171 $\pm$ 0.016	0.436 $\pm$ 0.026
TabDDPM	2.697 $\pm$ 0.289	1.330 $\pm$ 0.046	0.480 $\pm$ 0.025
TabMT	1.071 $\pm$ 0.205	1.236 $\pm$ 0.024	0.420 $\pm$ 0.021
MaCoDE	1.016 $\pm$ 0.086	0.940 $\pm$ 0.018	0.454 $\pm$ 0.026
banknote			
Dataset	$k$ -anonymity(%) $\uparrow$	DCR $\uparrow$	AD $\downarrow$
CTGAN	2.862 $\pm$ 0.146	0.341 $\pm$ 0.020	0.763 $\pm$ 0.033
TVAE	2.607 $\pm$ 0.305	0.220 $\pm$ 0.007	0.920 $\pm$ 0.021
CTAB-GAN	3.118 $\pm$ 0.202	0.275 $\pm$ 0.008	0.847 $\pm$ 0.025
CTAB-GAN+	2.799 $\pm$ 0.151	0.279 $\pm$ 0.012	0.560 $\pm$ 0.038
DistVAE	3.081 $\pm$ 0.267	0.290 $\pm$ 0.010	0.757 $\pm$ 0.027
TabDDPM	0.985 $\pm$ 0.180	0.168 $\pm$ 0.007	0.997 $\pm$ 0.003
TabMT	2.662 $\pm$ 0.262	0.257 $\pm$ 0.009	0.910 $\pm$ 0.015
MaCoDE	2.552 $\pm$ 0.255	0.076 $\pm$ 0.005	0.930 $\pm$ 0.016
concrete			
Dataset	$k$ -anonymity(%) $\uparrow$	DCR $\uparrow$	AD $\downarrow$
CTGAN	3.750 $\pm$ 0.305	1.236 $\pm$ 0.076	0.156 $\pm$ 0.040
TVAE	2.464 $\pm$ 0.386	0.681 $\pm$ 0.018	0.340 $\pm$ 0.054
CTAB-GAN	3.519 $\pm$ 0.216	0.973 $\pm$ 0.044	0.310 $\pm$ 0.054
CTAB-GAN+	4.066 $\pm$ 0.201	0.972 $\pm$ 0.034	0.310 $\pm$ 0.033
DistVAE	4.235 $\pm$ 0.231	1.012 $\pm$ 0.037	0.358 $\pm$ 0.051
TabDDPM	1.663 $\pm$ 0.253	0.425 $\pm$ 0.011	0.388 $\pm$ 0.044
TabMT	4.029 $\pm$ 0.228	0.949 $\pm$ 0.038	0.308 $\pm$ 0.053
MaCoDE	2.864 $\pm$ 0.293	0.239 $\pm$ 0.019	0.348 $\pm$ 0.043
kings			
Dataset	$k$ -anonymity(%) $\uparrow$	DCR $\uparrow$	AD $\downarrow$
CTGAN	1.243 $\pm$ 0.168	0.502 $\pm$ 0.006	0.601 $\pm$ 0.004
TVAE	0.848 $\pm$ 0.066	0.361 $\pm$ 0.003	0.685 $\pm$ 0.003
CTAB-GAN	1.158 $\pm$ 0.103	0.461 $\pm$ 0.004	0.631 $\pm$ 0.003
CTAB-GAN+	1.028 $\pm$ 0.126	0.411 $\pm$ 0.006	0.655 $\pm$ 0.003
DistVAE	0.699 $\pm$ 0.113	0.488 $\pm$ 0.002	0.658 $\pm$ 0.004
TabDDPM	0.236 $\pm$ 0.032	0.325 $\pm$ 0.003	0.668 $\pm$ 0.005
TabMT	0.396 $\pm$ 0.103	0.427 $\pm$ 0.003	0.652 $\pm$ 0.003
MaCoDE	0.445 $\pm$ 0.151	0.301 $\pm$ 0.005	0.677 $\pm$ 0.004
loan			
Dataset	$k$ -anonymity(%) $\uparrow$	DCR $\uparrow$	AD $\downarrow$
CTGAN	2.450 $\pm$ 0.244	0.241 $\pm$ 0.007	0.660 $\pm$ 0.012
TVAE	2.442 $\pm$ 0.338	0.160 $\pm$ 0.004	0.716 $\pm$ 0.008
CTAB-GAN	2.407 $\pm$ 0.156	0.208 $\pm$ 0.004	0.687 $\pm$ 0.010
CTAB-GAN+	2.570 $\pm$ 0.159	0.216 $\pm$ 0.015	0.577 $\pm$ 0.005
DistVAE	2.348 $\pm$ 0.092	0.238 $\pm$ 0.004	0.683 $\pm$ 0.002
TabDDPM	0.992 $\pm$ 0.107	0.123 $\pm$ 0.002	0.699 $\pm$ 0.006
TabMT	1.895 $\pm$ 0.126	0.190 $\pm$ 0.003	0.687 $\pm$ 0.003
MaCoDE	1.942 $\pm$ 0.212	0.127 $\pm$ 0.003	0.697 $\pm$ 0.007
whitewine			
Dataset	$k$ -anonymity(%) $\uparrow$	DCR $\uparrow$	AD $\downarrow$
CTGAN	2.580 $\pm$ 0.425	1.445 $\pm$ 0.018	0.239 $\pm$ 0.020
TVAE	1.399 $\pm$ 0.222	1.016 $\pm$ 0.005	0.475 $\pm$ 0.032
CTAB-GAN	1.746 $\pm$ 0.287	1.292 $\pm$ 0.027	0.373 $\pm$ 0.022
CTAB-GAN+	2.116 $\pm$ 0.171	1.145 $\pm$ 0.010	0.385 $\pm$ 0.018
DistVAE	2.762 $\pm$ 0.184	1.189 $\pm$ 0.007	0.461 $\pm$ 0.016
TabDDPM	0.896 $\pm$ 0.064	1.043 $\pm$ 0.009	0.441 $\pm$ 0.009
TabMT	1.618 $\pm$ 0.303	1.192 $\pm$ 0.013	0.446 $\pm$ 0.014
MaCoDE	1.539 $\pm$ 0.183	0.988 $\pm$ 0.013	0.446 $\pm$ 0.016

Table 6: **Privacy preservability** for each dataset. The means and standard errors of the mean across 10 repeated experiments are reported.  $\uparrow$  ( $\downarrow$ ) denotes higher (lower) is better.

### A.8.2 Sensitivity Analysis

- **Evaluation procedure:** We generated missing values into the kings dataset at rates of 0.1, 0.3, 0.5, and 0.7 for each missing mechanism and proceeded to train the model. Subsequently, we evaluated the machine learning utility using metrics such as SMAPE and  $F_1$ -score, along with assessing the multiple imputation performance of the model using the same methodology described earlier.
- **Result:** Figure 2 and 3 show the sensitivity analysis conducted by varying the missingness rate of kings dataset across four missing data mechanisms (MCAR, MAR, MNARL, MNARQ). In Figure 2, it's evident that MaCoDE maintains competitive performance in terms of SMAPE, even with increasing missingness rates across all missing data mechanisms. Regarding the  $F_1$  score, MaCoDE outperforms other models at missingness rates of 0.1 and 0.3, but its performance diminishes beyond a missingness rate of 0.5 (except for the MNARQ missing data mechanism). Additionally, as shown in Figure 3, MaCoDE consistently exhibits competitive performance in the multiple imputation, regardless of the increasing missingness rate, without significant performance degradation compared to other imputation models.

Hence, Figures 2 and 3 demonstrate that MaCoDE maintains comparable performance to other models even when trained on a dataset with missing values (i.e., incomplete dataset) without compromising the quality of the synthetic data it generates. This is a notable advantage of our proposed model over other baseline models, which struggle with training on datasets containing missing values.

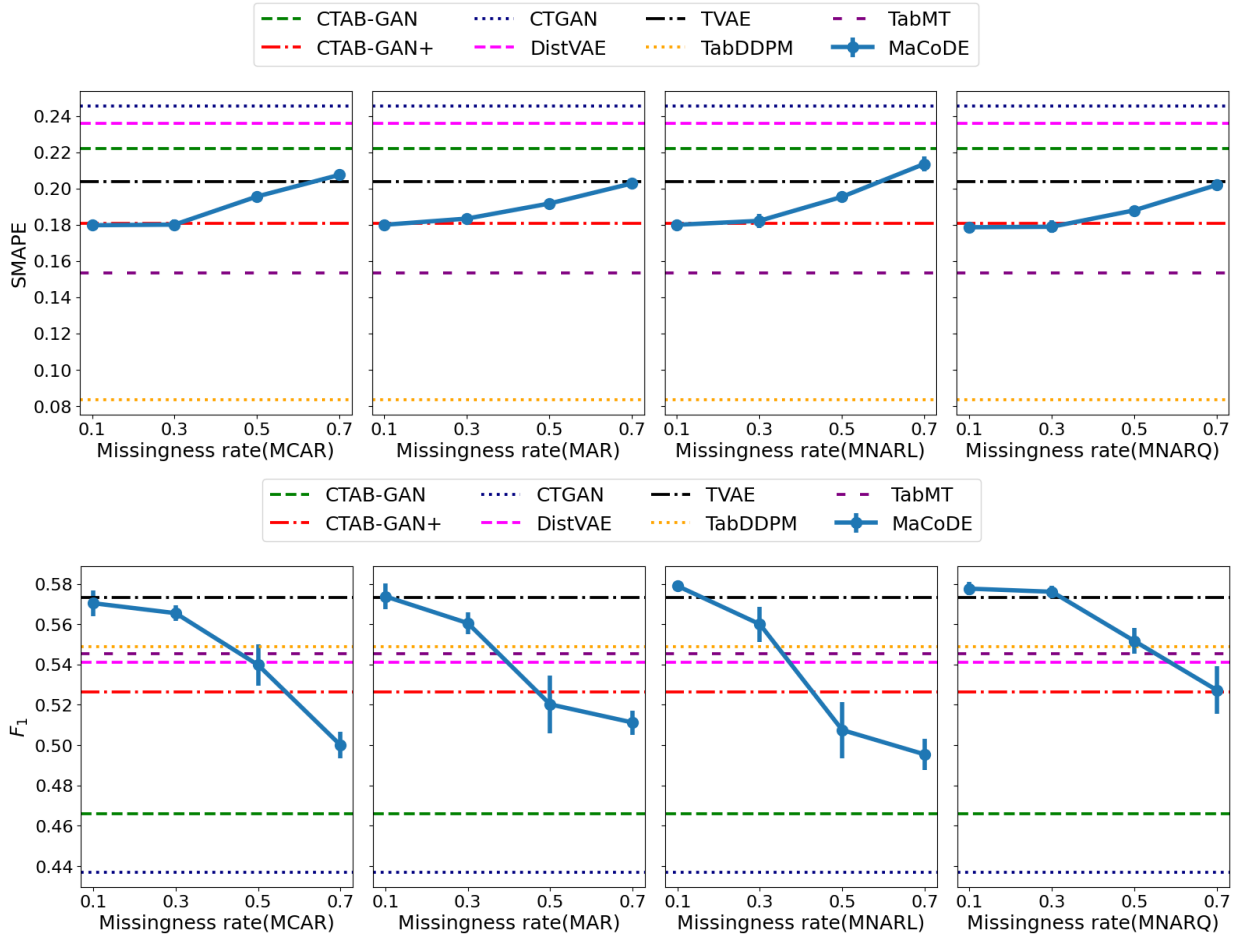


Figure 2: **Q2.** Sensitivity analysis of machine learning utility according to missingness rate. Machine learning utility is evaluated using kings dataset under four missing mechanisms. The means and standard errors of the mean across 10 repeated experiments are reported. Error bars represent standard errors.

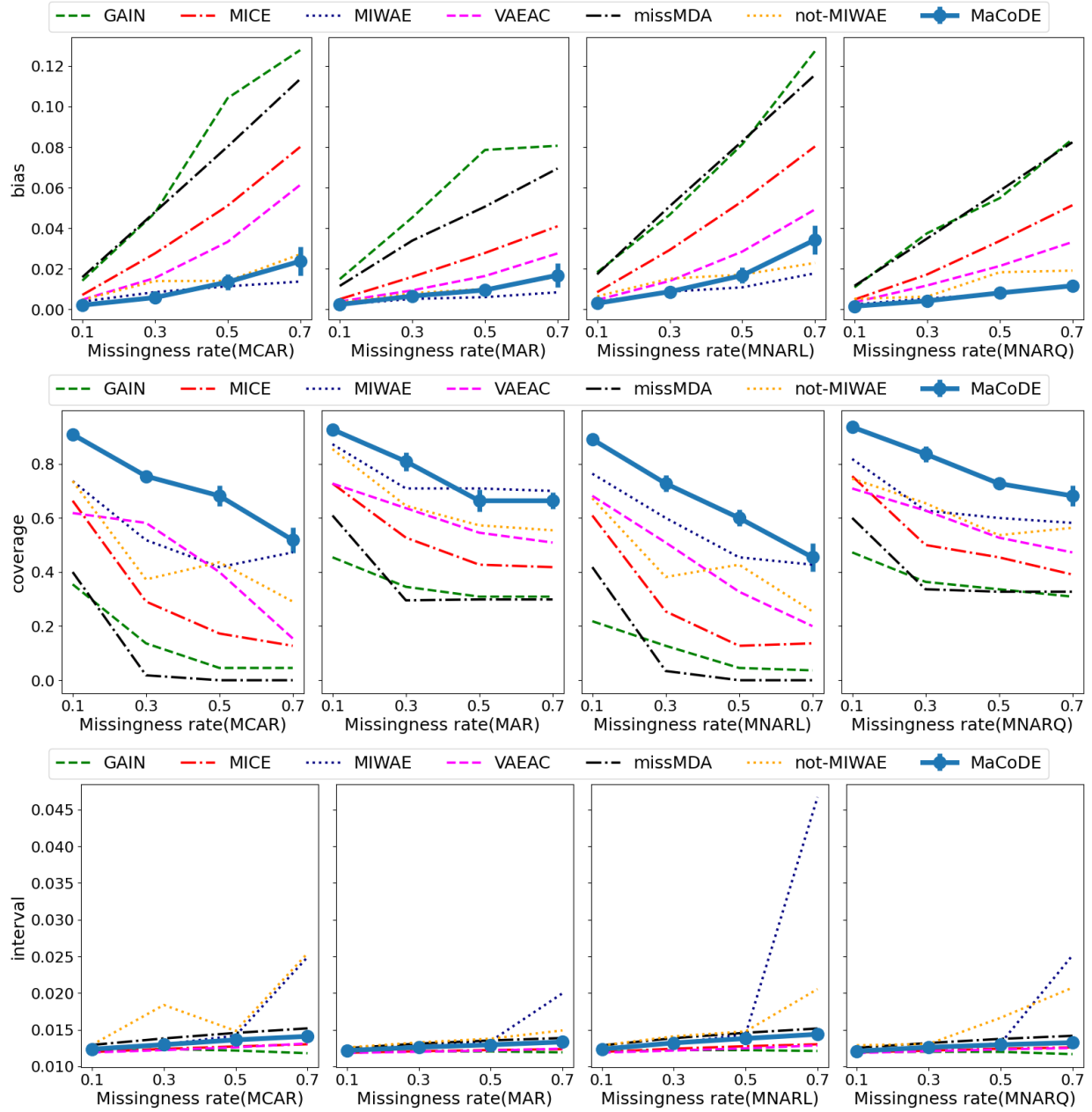


Figure 3: **Q3.** Sensitivity analysis of multiple imputations according to missingness rate. Multiple imputation is evaluated using kings dataset under four missing mechanisms. The means and standard errors of the mean across 10 repeated experiments are reported. Error bars represent standard errors.

## A.9 Detailed Experimental Results

### A.9.1 Q1. Synthetic Data Quality

abalone								
Model	Statistical fidelity				Machine learning utility			
	KL ↓	GoF ↓	MMD ↓	WD ↓	SMAPE ↓	$F_1$ ↑	Model ↑	Feature ↑
Baseline	.006±.000	.023±.001	.001±.000	.152±.010	.031±.000	.229±.003	.790±.057	.981±.006
CTGAN	.114±.011	.289±.058	.109±.010	1.174±.069	.101±.003	.166±.004	.625±.087	.283±.095
TVAE	.064±.003	.121±.004	.016±.001	.387±.007	.120±.005	.247±.003	.152±.146	.417±.105
CTAB-GAN	.208±.012	.403±.054	.090±.013	1.161±.120	.109±.005	.162±.005	.274±.161	.395±.158
CTAB-GAN+	.035±.008	.091±.020	.022±.006	.485±.094	.068±.006	.189±.013	.805±.147	.771±.100
DistVAE	.010±.001	.043±.003	.005±.000	.263±.007	.061±.001	.216±.002	.825±.050	.288±.034
TabDDPM	.119±.030	.905±.577	.005±.002	4.592±1.055	.033±.000	.222±.002	.747±.059	.957±.008
TabMT	.003±.000	.022±.001	.033±.001	.949±.033	.081±.001	.198±.003	.744±.050	.852±.022
MaCoDE	.017±.002	.064±.006	.005±.001	.150±.009	.039±.000	.217±.003	.744±.049	.952±.010
banknote								
Model	Statistical fidelity				Machine learning utility			
	KL ↓	GoF ↓	MMD ↓	WD ↓	SMAPE ↓	$F_1$ ↑	Model ↑	Feature ↑
Baseline	.016±.001	.040±.003	.002±.000	.059±.004	.232±.002	.957±.002	.903±.027	1.000±.000
CTGAN	.146±.013	.187±.013	.095±.009	1.260±.119	.646±.009	.745±.020	−.321±.143	.720±.085
TVAE	.047±.009	.084±.008	.019±.004	.260±.023	.436±.006	.879±.006	.616±.079	.940±.031
CTAB-GAN	.030±.003	.066±.003	.021±.001	.415±.018	.519±.004	.826±.011	.207±.166	.880±.033
CTAB-GAN+	.176±.076	.137±.042	.072±.038	.847±.227	.670±.031	.634±.059	−.388±.274	.280±.605
DistVAE	.075±.002	.068±.001	.031±.001	.563±.011	.651±.005	.763±.009	−.513±.081	1.000±.000
TabDDPM	.052±.009	.054±.004	.006±.001	.364±.059	.315±.004	.942±.003	.839±.033	1.000±.000
TabMT	.012±.001	.033±.001	.012±.001	.301±.008	.495±.005	.883±.007	.468±.126	.960±.027
MaCoDE	.031±.003	.074±.005	.009±.001	.148±.011	.374±.005	.918±.005	.652±.071	1.000±.000
breast								
Model	Statistical fidelity				Machine learning utility			
	KL ↓	GoF ↓	MMD ↓	WD ↓	SMAPE ↓	$F_1$ ↑	Model ↑	Feature ↑
Baseline	.081±.005	.080±.003	.007±.001	5.511±.116	.056±.001	.946±.002	.662±.106	.865±.015
CTGAN	.419±.029	.345±.013	.242±.008	35.684±1.061	.237±.006	.645±.046	.057±.167	−.163±.045
TVAE	.049±.002	.087±.002	.010±.001	6.519±.079	.090±.001	.938±.004	.687±.071	.895±.012
CTAB-GAN	.218±.010	.186±.005	.092±.005	16.515±.534	.151±.005	.877±.010	.350±.129	.558±.034
CTAB-GAN+	.564±.255	.356±.092	.229±.090	25.009±9.867	.242±.046	.696±.138	.265±.525	.187±.275
DistVAE	.071±.002	.090±.001	.033±.001	1.407±.107	.123±.002	.889±.006	.274±.071	.750±.022
TabDDPM	3.110±.025	.587±.006	.306±.001	262.226±6.849	.345±.008	.638±.016	−.260±.151	.364±.041
TabMT	.032±.002	.060±.002	.028±.001	12.405±.123	.128±.002	.911±.003	.503±.131	.756±.017
MaCoDE	.036±.003	.079±.006	.012±.001	8.634±.122	.098±.002	.921±.004	.332±.167	.807±.010
concrete								
Model	Statistical fidelity				Machine learning utility			
	KL ↓	GoF ↓	MMD ↓	WD ↓	SMAPE ↓	$F_1$ ↑	Model ↑	Feature ↑
Baseline	.026±.002	.055±.002	.003±.000	.433±.006	.121±.002	.466±.009	.912±.017	.950±.009
CTGAN	.384±.027	.830±.129	.134±.012	4.443±.289	.334±.004	.091±.020	.070±.151	.114±.138
TVAE	.125±.007	.219±.007	.017±.002	1.057±.031	.276±.003	.398±.010	.030±.161	.290±.045
CTAB-GAN	.204±.015	.226±.010	.051±.006	2.383±.098	.290±.004	.329±.011	.035±.061	−.260±.097
CTAB-GAN+	.263±.054	.278±.027	.066±.023	2.472±.380	.285±.012	.331±.045	−.009±.175	.017±.344
DistVAE	.186±.026	.159±.002	.024±.001	1.808±.011	.289±.003	.374±.007	−.160±.092	.107±.104
TabDDPM	.036±.001	.087±.007	.001±.000	.532±.031	.156±.004	.446±.007	.537±.049	.938±.019
TabMT	.031±.001	.115±.002	.019±.001	1.810±.034	.284±.003	.355±.010	.392±.060	.121±.086
MaCoDE	.089±.019	.118±.008	.006±.001	.684±.025	.229±.002	.405±.009	.439±.081	.919±.021
covtype								
Model	Statistical fidelity				Machine learning utility			
	KL ↓	GoF ↓	MMD ↓	WD ↓	SMAPE ↓	$F_1$ ↑	Model ↑	Feature ↑
Baseline	.000±.000	.002±.000	.000±.000	.454±.002	.067±.000	.804±.000	1.000±.000	1.000±.000
CTGAN	.105±.007	.493±.063	.017±.002	1.338±.053	.154±.001	.495±.023	.320±.098	.878±.011
TVAE	.048±.004	.078±.004	.013±.002	.950±.021	.151±.001	.634±.008	.420±.055	.930±.019
CTAB-GAN	.020±.002	.064±.005	.005±.000	.892±.019	.147±.001	.574±.005	.250±.067	.922±.010
CTAB-GAN+	.002±.000	.020±.003	.001±.000	.601±.007	.122±.001	.652±.003	.300±.000	.958±.006
DistVAE	.007±.000	.028±.001	.005±.000	1.228±.010	.182±.001	.626±.002	−.080±.033	.895±.016
TabDDPM	.003±.000	.011±.000	.001±.000	.556±.006	.115±.000	.665±.001	.300±.000	.964±.000
TabMT	.001±.000	.009±.000	.001±.000	.621±.005	.119±.000	.701±.002	.700±.000	.960±.002
MaCoDE	.001±.000	.026±.001	.001±.000	.509±.002	.098±.000	.742±.001	.900±.000	.956±.002

Table 7: **Q1**: Statistical fidelity and machine learning utility for each dataset. The means and the standard errors of the mean across 10 repeated experiments are reported. ‘Baseline’ refers to the result obtained using half of the real training dataset. ↑ (↓) denotes higher (lower) is better.

kings								
Model	Statistical fidelity				Machine learning utility			
	KL ↓	GoF ↓	MMD ↓	WD ↓	SMAPE ↓	$F_1$ ↑	Model ↑	Feature ↑
Baseline	.002±.000	.008±.000	.000±.000	.414±.014	.088±.001	.605±.003	1.000±.000	.994±.001
CTGAN	.141±.009	.753±.083	.036±.004	2.702±.347	.245±.003	.437±.010	.880±.025	.881±.008
TVAE	.032±.001	.097±.002	.004±.001	.942±.038	.204±.000	.573±.003	.940±.016	.894±.008
CTAB-GAN	.056±.004	.421±.223	.017±.001	1.388±.114	.222±.002	.466±.010	.947±.016	.898±.008
CTAB-GAN+	.037±.030	.100±.048	.039±.037	2.107±1.594	.181±.007	.526±.022	.937±.055	.956±.022
DistVAE	.159±.017	.085±.002	.005±.000	1.039±.035	.236±.001	.541±.004	.930±.030	.907±.003
TabDDPM	.096±.017	1.621±.684	.001±.000	27.567±6.882	.083±.001	.549±.007	.900±.038	.978±.005
TabMT	.001±.000	.010±.000	.001±.000	.741±.011	.154±.001	.545±.005	.940±.016	.964±.002
MaCoDE	.011±.001	.046±.002	.009±.002	.998±.089	.179±.001	.580±.004	.920±.025	.975±.002
letter								
Model	Statistical fidelity				Machine learning utility			
	KL ↓	GoF ↓	MMD ↓	WD ↓	SMAPE ↓	$F_1$ ↑	Model ↑	Feature ↑
Baseline	.001±.000	.010±.000	.000±.000	.921±.002	.059±.000	.821±.002	1.000±.000	.990±.001
CTGAN	.171±.004	.294±.004	.022±.002	4.743±.182	.169±.003	.169±.010	.120±.135	.555±.055
TVAE	.040±.002	.064±.002	.005±.000	1.604±.013	.106±.001	.496±.011	.550±.056	.944±.005
CTAB-GAN	.113±.005	.121±.008	.013±.001	3.921±.059	.153±.002	.197±.010	.130±.087	.729±.030
CTAB-GAN+	.051±.012	.073±.015	.005±.001	2.782±.112	.123±.001	.471±.021	.440±.126	.899±.023
DistVAE	.021±.001	.044±.001	.004±.000	2.526±.015	.120±.000	.577±.003	.110±.046	.854±.009
TabDDPM	.009±.000	.031±.001	.002±.000	1.371±.023	.075±.000	.448±.005	.500±.001	.531±.008
TabMT	.002±.000	.010±.000	.001±.000	1.579±.009	.083±.000	.745±.003	.900±.001	.976±.003
MaCoDE	.095±.002	.144±.003	.017±.001	2.135±.035	.106±.001	.689±.003	.850±.017	.965±.003
loan								
Model	Statistical fidelity				Machine learning utility			
	KL ↓	GoF ↓	MMD ↓	WD ↓	SMAPE ↓	$F_1$ ↑	Model ↑	Feature ↑
Baseline	.002±.000	.011±.000	.001±.000	.071±.001	.267±.001	.930±.001	.907±.037	.965±.008
CTGAN	.079±.009	.177±.014	.056±.006	.963±.070	.368±.003	.886±.004	.272±.079	.682±.045
TVAE	.121±.004	.170±.003	.023±.000	.751±.028	.339±.002	.850±.007	.685±.102	.916±.017
CTAB-GAN	.037±.004	.087±.004	.023±.003	.502±.032	.333±.004	.904±.003	.440±.097	.802±.014
CTAB-GAN+	.050±.028	.084±.009	.039±.024	.505±.251	.304±.020	.892±.011	.144±.313	.810±.043
DistVAE	.011±.001	.058±.002	.030±.000	.673±.007	.361±.001	.903±.003	.517±.069	.781±.016
TabDDPM	.004±.000	.010±.000	.001±.000	.075±.007	.281±.001	.925±.001	.677±.039	.978±.006
TabMT	.002±.000	.013±.001	.003±.000	.192±.008	.307±.001	.913±.002	.623±.043	.891±.012
MaCoDE	.031±.004	.047±.003	.007±.001	.126±.008	.273±.001	.919±.002	.726±.033	.896±.021
redwine								
Model	Statistical fidelity				Machine learning utility			
	KL ↓	GoF ↓	MMD ↓	WD ↓	SMAPE ↓	$F_1$ ↑	Model ↑	Feature ↑
Baseline	.025±.004	.042±.002	.003±.000	1.180±.027	.080±.001	.576±.005	.802±.072	.905±.014
CTGAN	.463±.017	.914±.178	.164±.007	8.166±.528	.177±.005	.228±.031	.116±.144	.071±.119
TVAE	.057±.002	.115±.002	.012±.001	1.842±.042	.103±.001	.569±.005	.238±.117	.684±.043
CAB-GAN	.111±.012	.143±.008	.042±.006	3.075±.150	.139±.003	.506±.007	.010±.109	.445±.056
CTAB-GAN+	.118±.025	.157±.031	.044±.014	2.969±.312	.152±.012	.488±.049	-.145±.413	.510±.343
DistVAE	.035±.001	.064±.001	.015±.000	2.423±.018	.134±.001	.533±.005	-.108±.091	.776±.030
TabDDPM	2.678±.108	.652±.050	.214±.009	86.984±4.091	.108±.002	.481±.010	.412±.122	.372±.066
TabMT	.023±.002	.052±.001	.015±.000	2.534±.024	.127±.002	.507±.009	.010±.112	.567±.057
MaCoDE	.022±.001	.071±.002	.005±.000	1.526±.016	.105±.001	.524±.007	.197±.092	.870±.032
whitewine								
Model	Statistical fidelity				Machine learning utility			
	KL ↓	GoF ↓	MMD ↓	WD ↓	SMAPE ↓	$F_1$ ↑	Model ↑	Feature ↑
Baseline	.006±.001	.022±.001	.001±.000	.993±.008	.065±.000	.523±.004	.897±.027	.913±.012
CTGAN	.186±.008	1.330±.159	.068±.005	3.881±.153	.123±.002	.248±.019	-.060±.064	.148±.083
TVAE	.078±.005	.152±.005	.040±.004	1.996±.064	.100±.001	.492±.003	.545±.127	.561±.066
CTAB-GAN	.166±.016	.241±.019	.080±.011	3.016±.147	.117±.002	.407±.010	-.010±.089	.311±.090
CTAB-GAN+	.064±.022	.139±.022	.029±.012	1.929±.265	.109±.004	.424±.021	-.080±.326	.620±.189
DistVAE	.018±.001	.059±.001	.009±.001	1.787±.017	.106±.001	.459±.003	.150±.095	.595±.026
TabDDPM	.675±.041	.161±.011	.022±.002	4.292±2.951	.079±.001	.439±.008	.535±.070	.685±.056
TabMT	.005±.001	.031±.000	.007±.000	1.861±.010	.102±.001	.462±.004	-.001±.068	.562±.035
MaCoDE	.010±.001	.049±.002	.002±.000	1.238±.008	.087±.000	.465±.003	.240±.070	.903±.013

Table 8: **Q1**: Statistical fidelity and machine learning utility for each dataset. The means and the standard errors of the mean across 10 repeated experiments are reported. ‘Baseline’ refers to the result obtained using half of the real training dataset. ↑ (↓) denotes higher (lower) is better.

## A.9.2 Q1. Visualization of Marginal Histogram

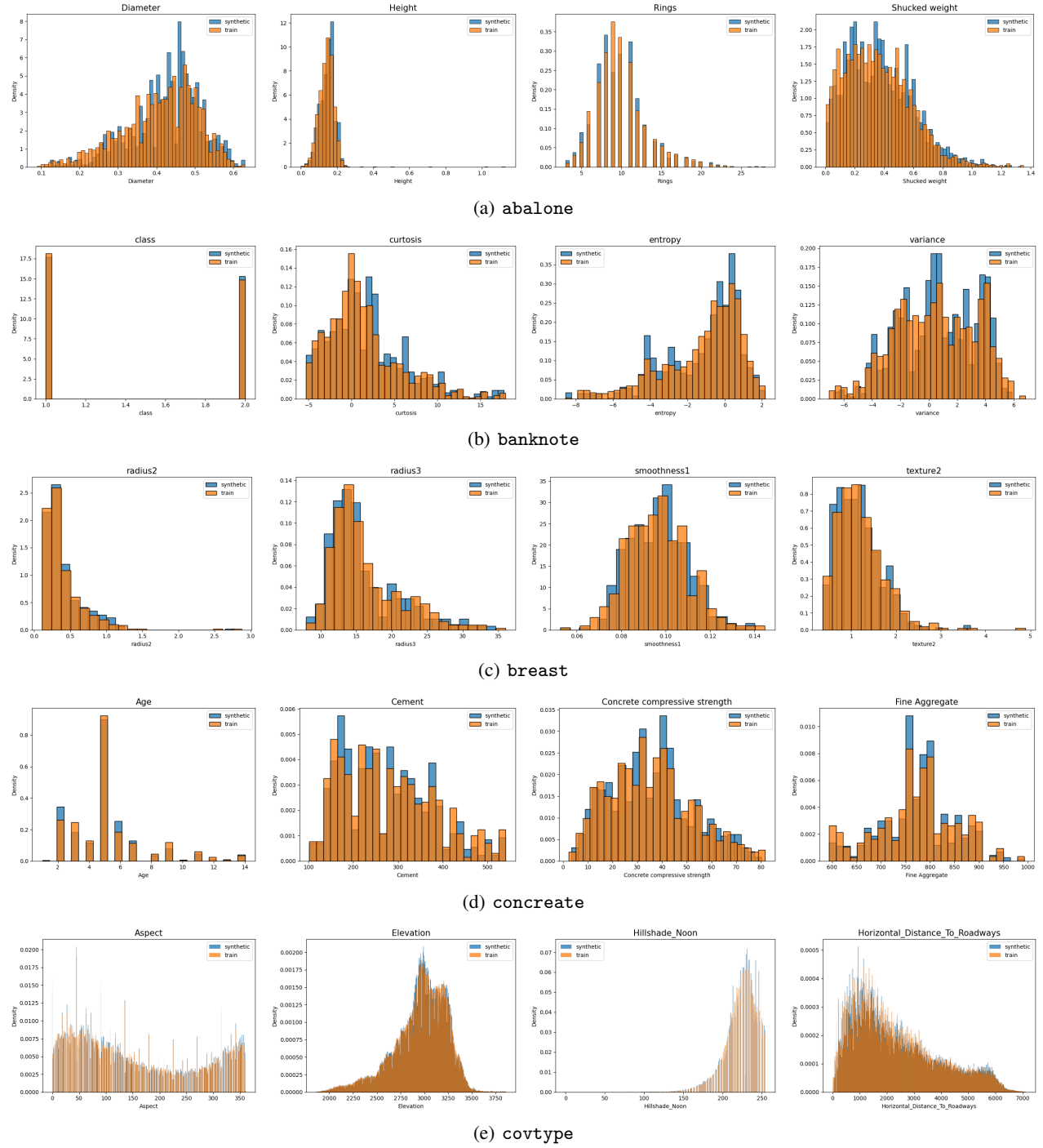


Figure 4: Histograms of observed dataset and synthetic dataset, generated by MaCoDE.



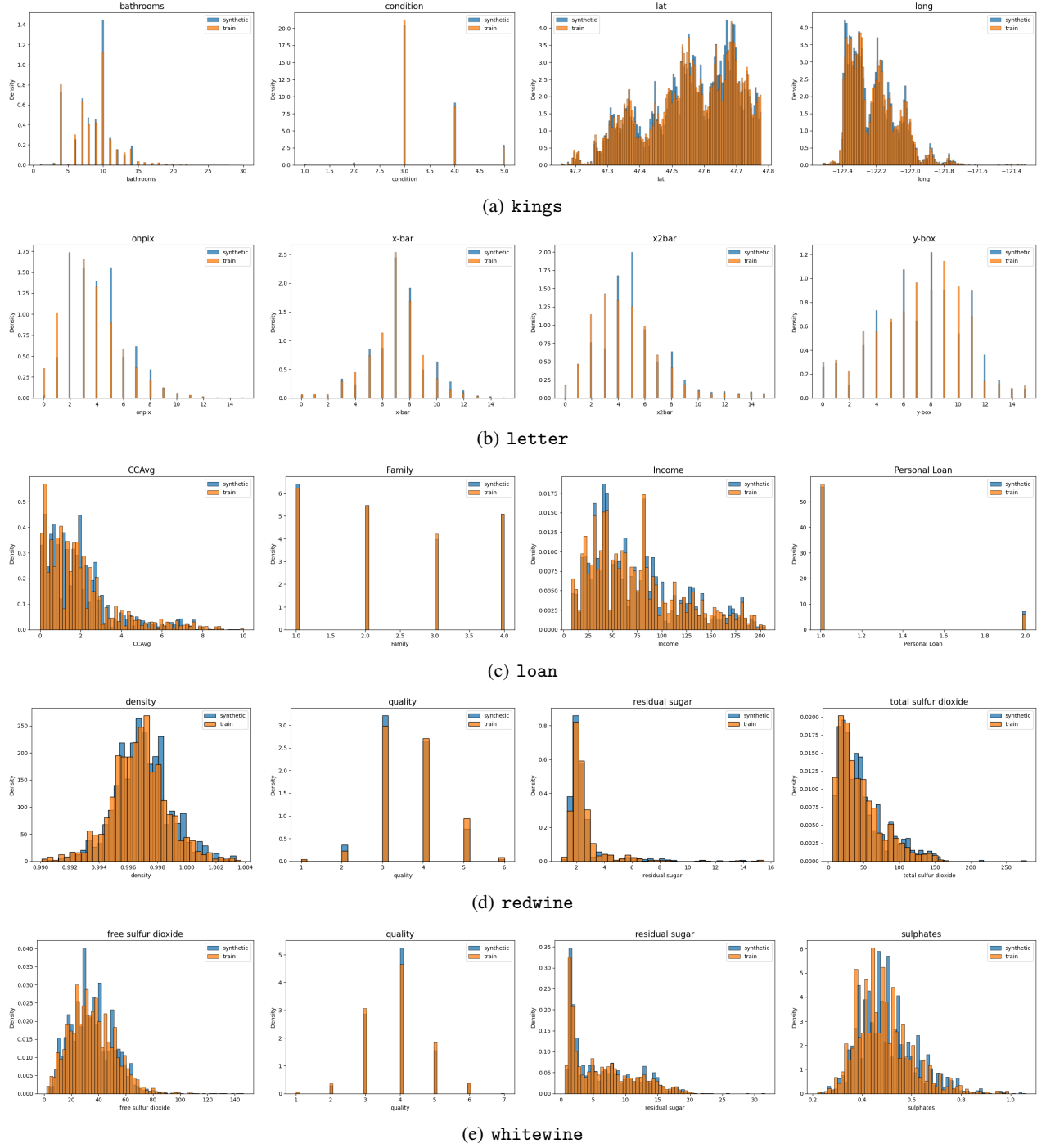


Figure 5: Histograms of observed dataset and synthetic dataset, generated by MaCoDE.

### A.9.3 Q2: Synthetic Data Quality in Scenarios with Incomplete Training Dataset

MCAR			MAR		
Dataset	SMAPE ↓	$F_1$ ↑	Dataset	SMAPE ↓	$F_1$ ↑
abalone	0.042 $\pm$ 0.000	0.216 $\pm$ 0.003	abalone	0.041 $\pm$ 0.000	0.211 $\pm$ 0.003
banknote	0.417 $\pm$ 0.004	0.879 $\pm$ 0.006	banknote	0.410 $\pm$ 0.005	0.891 $\pm$ 0.005
breast	0.111 $\pm$ 0.001	0.909 $\pm$ 0.006	breast	0.107 $\pm$ 0.002	0.916 $\pm$ 0.005
concrete	0.240 $\pm$ 0.002	0.373 $\pm$ 0.010	concrete	0.244 $\pm$ 0.003	0.383 $\pm$ 0.007
covtype	0.103 $\pm$ 0.000	0.721 $\pm$ 0.001	covtype	0.102 $\pm$ 0.000	0.723 $\pm$ 0.001
kings	0.180 $\pm$ 0.003	0.565 $\pm$ 0.004	kings	0.183 $\pm$ 0.001	0.561 $\pm$ 0.005
letter	0.107 $\pm$ 0.000	0.674 $\pm$ 0.004	letter	0.107 $\pm$ 0.000	0.682 $\pm$ 0.004
loan	0.278 $\pm$ 0.003	0.915 $\pm$ 0.002	loan	0.278 $\pm$ 0.001	0.914 $\pm$ 0.002
redwine	0.112 $\pm$ 0.001	0.515 $\pm$ 0.009	redwine	0.110 $\pm$ 0.001	0.518 $\pm$ 0.005
whitewine	0.091 $\pm$ 0.000	0.461 $\pm$ 0.005	whitewine	0.090 $\pm$ 0.001	0.465 $\pm$ 0.006

MNARL			MNARQ		
Dataset	SMAPE ↓	$F_1$ ↑	Dataset	SMAPE ↓	$F_1$ ↑
abalone	0.042 $\pm$ 0.000	0.212 $\pm$ 0.004	abalone	0.041 $\pm$ 0.000	0.211 $\pm$ 0.003
banknote	0.418 $\pm$ 0.007	0.883 $\pm$ 0.006	banknote	0.399 $\pm$ 0.003	0.899 $\pm$ 0.009
breast	0.112 $\pm$ 0.002	0.908 $\pm$ 0.006	breast	0.105 $\pm$ 0.001	0.918 $\pm$ 0.003
concrete	0.243 $\pm$ 0.003	0.374 $\pm$ 0.011	concrete	0.239 $\pm$ 0.003	0.384 $\pm$ 0.009
covtype	0.104 $\pm$ 0.000	0.718 $\pm$ 0.001	covtype	0.101 $\pm$ 0.000	0.730 $\pm$ 0.001
kings	0.182 $\pm$ 0.004	0.560 $\pm$ 0.009	kings	0.179 $\pm$ 0.004	0.576 $\pm$ 0.003
letter	0.108 $\pm$ 0.001	0.676 $\pm$ 0.003	letter	0.107 $\pm$ 0.001	0.676 $\pm$ 0.004
loan	0.277 $\pm$ 0.001	0.915 $\pm$ 0.002	loan	0.275 $\pm$ 0.001	0.916 $\pm$ 0.001
redwine	0.112 $\pm$ 0.001	0.511 $\pm$ 0.009	redwine	0.108 $\pm$ 0.001	0.527 $\pm$ 0.006
whitewine	0.091 $\pm$ 0.001	0.464 $\pm$ 0.006	whitewine	0.089 $\pm$ 0.001	0.460 $\pm$ 0.004

Table 9: **Q2**: Machine learning utility for each dataset under **MCAR**, **MAR**, **MNARL**, and **MNARQ** at 0.3 missingness. The means and standard errors of the mean 10 repeated experiments are reported.  $\uparrow$  ( $\downarrow$ ) denotes higher (lower) is better.

#### A.9.4 Q3: Multiple Imputation Performance

MCAR			
Model	Bias ↓	Coverage	Width ↓
MICE	0.014 $\pm$ 0.001	0.795 $\pm$ 0.022	0.040 $\pm$ 0.002
GAIN	0.024 $\pm$ 0.002	0.591 $\pm$ 0.037	0.040 $\pm$ 0.002
missMDA	0.018 $\pm$ 0.001	0.644 $\pm$ 0.015	0.045 $\pm$ 0.002
VAEAC	0.010 $\pm$ 0.001	0.874 $\pm$ 0.022	0.041 $\pm$ 0.002
MIWAE	0.008 $\pm$ 0.000	0.936 $\pm$ 0.014	0.045 $\pm$ 0.002
notMIWAE	0.008 $\pm$ 0.000	0.910 $\pm$ 0.013	0.044 $\pm$ 0.002
EGC	0.005 $\pm$ 0.000	1.000 $\pm$ 0.000	0.060 $\pm$ 0.002
MaCoDE(MCAR)	0.008 $\pm$ 0.000	0.950 $\pm$ 0.011	0.053 $\pm$ 0.002
MAR			
Model	Bias ↓	Coverage	Width ↓
MICE	0.010 $\pm$ 0.001	0.845 $\pm$ 0.019	0.040 $\pm$ 0.002
GAIN	0.019 $\pm$ 0.002	0.633 $\pm$ 0.033	0.040 $\pm$ 0.002
missMDA	0.015 $\pm$ 0.001	0.700 $\pm$ 0.022	0.043 $\pm$ 0.002
VAEAC	0.008 $\pm$ 0.001	0.905 $\pm$ 0.016	0.040 $\pm$ 0.002
MIWAE	0.006 $\pm$ 0.000	0.952 $\pm$ 0.012	0.043 $\pm$ 0.002
notMIWAE	0.006 $\pm$ 0.000	0.949 $\pm$ 0.012	0.042 $\pm$ 0.002
EGC	0.006 $\pm$ 0.000	0.996 $\pm$ 0.004	0.058 $\pm$ 0.002
MaCoDE(MAR)	0.006 $\pm$ 0.000	0.963 $\pm$ 0.009	0.051 $\pm$ 0.003
MNARL			
Model	Bias ↓	Coverage	Width ↓
MICE	0.012 $\pm$ 0.001	0.790 $\pm$ 0.026	0.037 $\pm$ 0.002
GAIN	0.026 $\pm$ 0.002	0.538 $\pm$ 0.033	0.040 $\pm$ 0.002
missMDA	0.020 $\pm$ 0.001	0.592 $\pm$ 0.025	0.045 $\pm$ 0.002
VAEAC	0.011 $\pm$ 0.001	0.851 $\pm$ 0.023	0.041 $\pm$ 0.002
MIWAE	0.009 $\pm$ 0.001	0.920 $\pm$ 0.015	0.045 $\pm$ 0.002
notMIWAE	0.009 $\pm$ 0.001	0.908 $\pm$ 0.012	0.044 $\pm$ 0.002
EGC	0.007 $\pm$ 0.000	0.994 $\pm$ 0.004	0.061 $\pm$ 0.002
MaCoDE(MNARL)	0.008 $\pm$ 0.000	0.948 $\pm$ 0.013	0.052 $\pm$ 0.003
MNARQ			
Model	Bias ↓	Coverage	Width ↓
MICE	0.007 $\pm$ 0.001	0.880 $\pm$ 0.018	0.040 $\pm$ 0.002
GAIN	0.011 $\pm$ 0.001	0.819 $\pm$ 0.020	0.040 $\pm$ 0.002
missMDA	0.011 $\pm$ 0.001	0.785 $\pm$ 0.016	0.043 $\pm$ 0.002
VAEAC	0.006 $\pm$ 0.001	0.909 $\pm$ 0.020	0.040 $\pm$ 0.002
MIWAE	0.005 $\pm$ 0.000	0.969 $\pm$ 0.009	0.043 $\pm$ 0.002
notMIWAE	0.005 $\pm$ 0.000	0.958 $\pm$ 0.010	0.043 $\pm$ 0.002
EGC	0.004 $\pm$ 0.000	1.000 $\pm$ 0.000	0.054 $\pm$ 0.002
MaCoDE(MNARQ)	0.005 $\pm$ 0.000	0.944 $\pm$ 0.012	0.051 $\pm$ 0.003

Table 10: **Q3**: Multiple imputation under **MCAR**, **MAR**, **MNARL**, and **MNARQ** at 0.3 missingness. The means and standard errors of the mean across 5 datasets and 10 repeated experiments are reported. ↓ denotes the lower is better.

Dataset				abalone
Model	Bias ↓	Coverage	Width ↓	
MICE	0.005 $\pm$ 0.000	1.000 $\pm$ 0.000	0.030 $\pm$ 0.000	
GAIN	0.009 $\pm$ 0.001	0.829 $\pm$ 0.042	0.030 $\pm$ 0.000	
missMDA	0.014 $\pm$ 0.000	0.686 $\pm$ 0.036	0.034 $\pm$ 0.001	
VAEAC	0.005 $\pm$ 0.002	0.943 $\pm$ 0.074	0.030 $\pm$ 0.000	
MIWAE	0.006 $\pm$ 0.001	0.971 $\pm$ 0.019	0.034 $\pm$ 0.001	
not-MIWAE	0.009 $\pm$ 0.001	0.886 $\pm$ 0.019	0.037 $\pm$ 0.001	
EGC	0.004 $\pm$ 0.000	1.000 $\pm$ 0.000	0.040 $\pm$ 0.000	
MaCoDE	0.007 $\pm$ 0.001	0.914 $\pm$ 0.023	0.036 $\pm$ 0.002	
Dataset				banknote
Model	Bias ↓	Coverage	Width ↓	
MICE	0.018 $\pm$ 0.001	0.725 $\pm$ 0.025	0.053 $\pm$ 0.000	
GAIN	0.009 $\pm$ 0.001	0.829 $\pm$ 0.042	0.030 $\pm$ 0.000	
missMDA	0.018 $\pm$ 0.000	0.625 $\pm$ 0.042	0.056 $\pm$ 0.000	
VAEAC	0.009 $\pm$ 0.002	1.000 $\pm$ 0.000	0.054 $\pm$ 0.000	
MIWAE	0.006 $\pm$ 0.000	1.000 $\pm$ 0.000	0.055 $\pm$ 0.000	
not-MIWAE	0.006 $\pm$ 0.000	1.000 $\pm$ 0.000	0.055 $\pm$ 0.000	
EGC	0.004 $\pm$ 0.001	1.000 $\pm$ 0.000	0.066 $\pm$ 0.001	
MaCoDE	0.006 $\pm$ 0.001	1.000 $\pm$ 0.000	0.054 $\pm$ 0.000	
Dataset				breast
Model	Bias ↓	Coverage	Width ↓	
MICE	0.007 $\pm$ 0.000	1.000 $\pm$ 0.000	0.080 $\pm$ 0.000	
GAIN	0.012 $\pm$ 0.000	0.990 $\pm$ 0.005	0.080 $\pm$ 0.000	
missMDA	0.025 $\pm$ 0.000	0.883 $\pm$ 0.013	0.086 $\pm$ 0.000	
VAEAC	0.008 $\pm$ 0.000	1.000 $\pm$ 0.000	0.082 $\pm$ 0.000	
MIWAE	0.005 $\pm$ 0.000	1.000 $\pm$ 0.000	0.084 $\pm$ 0.000	
not-MIWAE	0.006 $\pm$ 0.000	1.000 $\pm$ 0.000	0.082 $\pm$ 0.000	
EGC	0.005 $\pm$ 0.000	1.000 $\pm$ 0.000	0.084 $\pm$ 0.000	
MaCoDE	0.009 $\pm$ 0.000	1.000 $\pm$ 0.000	0.083 $\pm$ 0.000	
Dataset				redwine
Model	Bias ↓	Coverage	Width ↓	
MICE	0.020 $\pm$ 0.000	0.764 $\pm$ 0.020	0.048 $\pm$ 0.000	
GAIN	0.024 $\pm$ 0.001	0.627 $\pm$ 0.050	0.048 $\pm$ 0.000	
missMDA	0.026 $\pm$ 0.000	0.618 $\pm$ 0.012	0.055 $\pm$ 0.000	
VAEAC	0.014 $\pm$ 0.001	0.855 $\pm$ 0.028	0.050 $\pm$ 0.001	
MIWAE	0.009 $\pm$ 0.001	0.945 $\pm$ 0.020	0.054 $\pm$ 0.000	
not-MIWAE	0.010 $\pm$ 0.000	0.909 $\pm$ 0.000	0.053 $\pm$ 0.000	
EGC	0.007 $\pm$ 0.000	1.000 $\pm$ 0.000	0.070 $\pm$ 0.000	
MaCoDE	0.008 $\pm$ 0.001	0.973 $\pm$ 0.014	0.056 $\pm$ 0.000	
Dataset				whitewine
Model	Bias ↓	Coverage	Width ↓	
MICE	0.011 $\pm$ 0.000	0.691 $\pm$ 0.031	0.028 $\pm$ 0.000	
GAIN	0.026 $\pm$ 0.003	0.482 $\pm$ 0.036	0.028 $\pm$ 0.000	
missMDA	0.015 $\pm$ 0.000	0.645 $\pm$ 0.016	0.035 $\pm$ 0.000	
VAEAC	0.010 $\pm$ 0.001	0.700 $\pm$ 0.038	0.028 $\pm$ 0.000	
MIWAE	0.011 $\pm$ 0.001	0.827 $\pm$ 0.029	0.036 $\pm$ 0.001	
not-MIWAE	0.009 $\pm$ 0.000	0.845 $\pm$ 0.033	0.033 $\pm$ 0.001	
EGC	0.006 $\pm$ 0.000	1.000 $\pm$ 0.000	0.062 $\pm$ 0.000	
MaCoDE	0.009 $\pm$ 0.001	0.864 $\pm$ 0.034	0.037 $\pm$ 0.001	

Table 11: **Q3**: Multiple imputation for each dataset under **MCAR** at 0.3 missingness. The means and standard errors of the mean across 5 datasets and 10 repeated experiments are reported. ↓ denotes lower is better.

Dataset				abalone
Model	Bias ↓	Coverage	Width ↓	
MICE	0.005 $\pm$ 0.001	0.971 $\pm$ 0.019	0.030 $\pm$ 0.000	
GAIN	0.011 $\pm$ 0.002	0.729 $\pm$ 0.072	0.030 $\pm$ 0.000	
missMDA	0.015 $\pm$ 0.001	0.571 $\pm$ 0.048	0.032 $\pm$ 0.001	
VAEAC	0.003 $\pm$ 0.001	0.986 $\pm$ 0.014	0.031 $\pm$ 0.000	
MIWAE	0.005 $\pm$ 0.001	0.971 $\pm$ 0.019	0.033 $\pm$ 0.001	
not-MIWAE	0.006 $\pm$ 0.001	0.943 $\pm$ 0.023	0.032 $\pm$ 0.001	
EGC	0.004 $\pm$ 0.001	0.986 $\pm$ 0.014	0.035 $\pm$ 0.001	
MaCoDE	0.005 $\pm$ 0.001	0.957 $\pm$ 0.022	0.033 $\pm$ 0.001	
Dataset				banknote
Model	Bias ↓	Coverage	Width ↓	
MICE	0.012 $\pm$ 0.002	0.800 $\pm$ 0.033	0.052 $\pm$ 0.000	
GAIN	0.021 $\pm$ 0.004	0.675 $\pm$ 0.084	0.052 $\pm$ 0.000	
missMDA	0.014 $\pm$ 0.001	0.800 $\pm$ 0.033	0.055 $\pm$ 0.000	
VAEAC	0.009 $\pm$ 0.001	0.925 $\pm$ 0.038	0.053 $\pm$ 0.000	
MIWAE	0.004 $\pm$ 0.001	1.000 $\pm$ 0.000	0.054 $\pm$ 0.000	
not-MIWAE	0.004 $\pm$ 0.001	1.000 $\pm$ 0.000	0.054 $\pm$ 0.000	
EGC	0.005 $\pm$ 0.001	1.000 $\pm$ 0.000	0.071 $\pm$ 0.003	
MaCoDE	0.004 $\pm$ 0.001	1.000 $\pm$ 0.000	0.054 $\pm$ 0.000	
Dataset				breast
Model	Bias ↓	Coverage	Width ↓	
MICE	0.005 $\pm$ 0.000	1.000 $\pm$ 0.000	0.080 $\pm$ 0.000	
GAIN	0.011 $\pm$ 0.001	0.960 $\pm$ 0.016	0.080 $\pm$ 0.000	
missMDA	0.019 $\pm$ 0.001	0.887 $\pm$ 0.015	0.084 $\pm$ 0.000	
VAEAC	0.006 $\pm$ 0.000	0.997 $\pm$ 0.002	0.081 $\pm$ 0.000	
MIWAE	0.004 $\pm$ 0.000	1.000 $\pm$ 0.000	0.082 $\pm$ 0.000	
not-MIWAE	0.005 $\pm$ 0.001	0.993 $\pm$ 0.007	0.081 $\pm$ 0.000	
EGC	0.006 $\pm$ 0.000	1.000 $\pm$ 0.000	0.083 $\pm$ 0.000	
MaCoDE	0.007 $\pm$ 0.000	0.997 $\pm$ 0.003	0.083 $\pm$ 0.000	
Dataset				redwine
Model	Bias ↓	Coverage	Width ↓	
MICE	0.015 $\pm$ 0.001	0.836 $\pm$ 0.035	0.048 $\pm$ 0.000	
GAIN	0.020 $\pm$ 0.001	0.627 $\pm$ 0.021	0.048 $\pm$ 0.000	
missMDA	0.020 $\pm$ 0.002	0.709 $\pm$ 0.035	0.053 $\pm$ 0.000	
VAEAC	0.011 $\pm$ 0.001	0.882 $\pm$ 0.019	0.049 $\pm$ 0.000	
MIWAE	0.008 $\pm$ 0.001	0.945 $\pm$ 0.024	0.053 $\pm$ 0.001	
not-MIWAE	0.009 $\pm$ 0.001	0.909 $\pm$ 0.030	0.051 $\pm$ 0.000	
EGC	0.007 $\pm$ 0.001	1.000 $\pm$ 0.000	0.067 $\pm$ 0.001	
MaCoDE	0.008 $\pm$ 0.000	0.980 $\pm$ 0.013	0.054 $\pm$ 0.001	
Dataset				whitewine
Model	Bias ↓	Coverage	Width ↓	
MICE	0.008 $\pm$ 0.001	0.773 $\pm$ 0.034	0.028 $\pm$ 0.000	
GAIN	0.025 $\pm$ 0.003	0.500 $\pm$ 0.053	0.028 $\pm$ 0.000	
missMDA	0.011 $\pm$ 0.000	0.718 $\pm$ 0.032	0.033 $\pm$ 0.001	
VAEAC	0.008 $\pm$ 0.001	0.827 $\pm$ 0.029	0.029 $\pm$ 0.000	
MIWAE	0.007 $\pm$ 0.001	0.891 $\pm$ 0.030	0.033 $\pm$ 0.001	
not-MIWAE	0.006 $\pm$ 0.001	0.945 $\pm$ 0.024	0.031 $\pm$ 0.000	
EGC	0.006 $\pm$ 0.000	1.000 $\pm$ 0.000	0.058 $\pm$ 0.001	
MaCoDE	0.007 $\pm$ 0.003	0.882 $\pm$ 0.075	0.033 $\pm$ 0.003	

Table 12: **Q3**: Multiple imputation for each dataset under **MAR** at 0.3 missingness. The means and standard errors of the mean across 5 datasets and 10 repeated experiments are reported. ↓ denotes lower is better. ↓ denotes lower is better.

Dataset				abalone
Model	Bias ↓	Coverage	Width ↓	
MICE	0.007 $\pm$ 0.001	0.929 $\pm$ 0.024	0.030 $\pm$ 0.000	
GAIN	0.015 $\pm$ 0.002	0.629 $\pm$ 0.061	0.030 $\pm$ 0.000	
missMDA	0.018 $\pm$ 0.001	0.471 $\pm$ 0.057	0.034 $\pm$ 0.001	
VAEAC	0.005 $\pm$ 0.001	0.957 $\pm$ 0.022	0.030 $\pm$ 0.000	
MIWAE	0.008 $\pm$ 0.001	0.900 $\pm$ 0.043	0.034 $\pm$ 0.001	
not-MIWAE	0.010 $\pm$ 0.002	0.886 $\pm$ 0.060	0.035 $\pm$ 0.003	
EGC	0.005 $\pm$ 0.000	0.986 $\pm$ 0.014	0.039 $\pm$ 0.001	
MaCoDE	0.007 $\pm$ 0.001	0.929 $\pm$ 0.024	0.035 $\pm$ 0.001	
Dataset				banknote
Model	Bias ↓	Coverage	Width ↓	
MICE	0.017 $\pm$ 0.001	0.775 $\pm$ 0.025	0.053 $\pm$ 0.000	
GAIN	0.032 $\pm$ 0.005	0.625 $\pm$ 0.077	0.053 $\pm$ 0.000	
missMDA	0.019 $\pm$ 0.001	0.725 $\pm$ 0.045	0.056 $\pm$ 0.000	
VAEAC	0.010 $\pm$ 0.001	0.975 $\pm$ 0.025	0.053 $\pm$ 0.000	
MIWAE	0.005 $\pm$ 0.001	1.000 $\pm$ 0.000	0.055 $\pm$ 0.000	
not-MIWAE	0.005 $\pm$ 0.003	1.000 $\pm$ 0.000	0.055 $\pm$ 0.000	
EGC	0.007 $\pm$ 0.001	1.000 $\pm$ 0.000	0.072 $\pm$ 0.002	
MaCoDE	0.006 $\pm$ 0.001	1.000 $\pm$ 0.000	0.054 $\pm$ 0.000	
Dataset				breast
Model	Bias ↓	Coverage	Width ↓	
MICE	0.008 $\pm$ 0.000	1.000 $\pm$ 0.000	0.080 $\pm$ 0.000	
GAIN	0.014 $\pm$ 0.001	0.960 $\pm$ 0.011	0.080 $\pm$ 0.000	
missMDA	0.026 $\pm$ 0.001	0.853 $\pm$ 0.015	0.086 $\pm$ 0.000	
VAEAC	0.009 $\pm$ 0.000	1.000 $\pm$ 0.000	0.082 $\pm$ 0.000	
MIWAE	0.006 $\pm$ 0.000	1.000 $\pm$ 0.000	0.084 $\pm$ 0.000	
not-MIWAE	0.007 $\pm$ 0.002	0.990 $\pm$ 0.016	0.082 $\pm$ 0.000	
EGC	0.007 $\pm$ 0.000	1.000 $\pm$ 0.000	0.084 $\pm$ 0.000	
MaCoDE	0.009 $\pm$ 0.000	1.000 $\pm$ 0.000	0.084 $\pm$ 0.000	
Dataset				redwine
Model	Bias ↓	Coverage	Width ↓	
MICE	0.021 $\pm$ 0.001	0.673 $\pm$ 0.034	0.048 $\pm$ 0.000	
GAIN	0.026 $\pm$ 0.002	0.473 $\pm$ 0.057	0.048 $\pm$ 0.000	
missMDA	0.028 $\pm$ 0.000	0.600 $\pm$ 0.028	0.055 $\pm$ 0.000	
VAEAC	0.016 $\pm$ 0.001	0.782 $\pm$ 0.036	0.050 $\pm$ 0.000	
MIWAE	0.012 $\pm$ 0.001	0.909 $\pm$ 0.019	0.055 $\pm$ 0.001	
not-MIWAE	0.013 $\pm$ 0.001	0.855 $\pm$ 0.047	0.053 $\pm$ 0.001	
EGC	0.009 $\pm$ 0.000	1.000 $\pm$ 0.000	0.072 $\pm$ 0.001	
MaCoDE	0.011 $\pm$ 0.001	0.990 $\pm$ 0.010	0.055 $\pm$ 0.001	
Dataset				whitewine
Model	Bias ↓	Coverage	Width ↓	
MICE	0.012 $\pm$ 0.001	0.673 $\pm$ 0.039	0.028 $\pm$ 0.000	
GAIN	0.029 $\pm$ 0.004	0.427 $\pm$ 0.045	0.028 $\pm$ 0.000	
missMDA	0.015 $\pm$ 0.000	0.582 $\pm$ 0.034	0.035 $\pm$ 0.001	
VAEAC	0.012 $\pm$ 0.001	0.691 $\pm$ 0.024	0.029 $\pm$ 0.000	
MIWAE	0.010 $\pm$ 0.001	0.873 $\pm$ 0.024	0.036 $\pm$ 0.001	
not-MIWAE	0.008 $\pm$ 0.001	0.891 $\pm$ 0.072	0.033 $\pm$ 0.002	
EGC	0.008 $\pm$ 0.000	0.991 $\pm$ 0.009	0.063 $\pm$ 0.001	
MaCoDE	0.009 $\pm$ 0.001	0.827 $\pm$ 0.037	0.034 $\pm$ 0.001	

Table 13: **Q3**: Multiple imputation for each dataset under **MNARL** at 0.3 missingness. The means and standard errors of the mean across 5 datasets and 10 repeated experiments are reported. ↓ denotes lower is better.

Dataset				abalone
Model	Bias ↓	Coverage	Width ↓	
MICE	0.003 $\pm$ 0.000	1.000 $\pm$ 0.000	0.030 $\pm$ 0.000	
GAIN	0.008 $\pm$ 0.001	0.886 $\pm$ 0.036	0.031 $\pm$ 0.001	
missMDA	0.009 $\pm$ 0.000	0.800 $\pm$ 0.032	0.033 $\pm$ 0.000	
VAEAC	0.002 $\pm$ 0.000	1.000 $\pm$ 0.000	0.030 $\pm$ 0.000	
MIWAE	0.005 $\pm$ 0.001	0.957 $\pm$ 0.022	0.032 $\pm$ 0.001	
not-MIWAE	0.007 $\pm$ 0.001	0.943 $\pm$ 0.023	0.035 $\pm$ 0.001	
EGC	0.003 $\pm$ 0.000	1.000 $\pm$ 0.000	0.038 $\pm$ 0.001	
MaCoDE	0.007 $\pm$ 0.001	0.886 $\pm$ 0.019	0.033 $\pm$ 0.001	
Dataset				banknote
Model	Bias ↓	Coverage	Width ↓	
MICE	0.008 $\pm$ 0.001	0.875 $\pm$ 0.042	0.052 $\pm$ 0.000	
GAIN	0.013 $\pm$ 0.002	0.825 $\pm$ 0.053	0.052 $\pm$ 0.000	
missMDA	0.009 $\pm$ 0.001	0.775 $\pm$ 0.045	0.054 $\pm$ 0.000	
VAEAC	0.005 $\pm$ 0.000	1.000 $\pm$ 0.000	0.053 $\pm$ 0.000	
MIWAE	0.003 $\pm$ 0.000	1.000 $\pm$ 0.000	0.054 $\pm$ 0.000	
not-MIWAE	0.003 $\pm$ 0.000	1.000 $\pm$ 0.000	0.054 $\pm$ 0.000	
EGC	0.004 $\pm$ 0.001	1.000 $\pm$ 0.000	0.066 $\pm$ 0.001	
MaCoDE	0.003 $\pm$ 0.000	1.000 $\pm$ 0.000	0.053 $\pm$ 0.000	
Dataset				breast
Model	Bias ↓	Coverage	Width ↓	
MICE	0.004 $\pm$ 0.000	0.997 $\pm$ 0.003	0.080 $\pm$ 0.000	
GAIN	0.007 $\pm$ 0.000	0.997 $\pm$ 0.003	0.080 $\pm$ 0.000	
missMDA	0.017 $\pm$ 0.001	0.927 $\pm$ 0.012	0.084 $\pm$ 0.000	
VAEAC	0.006 $\pm$ 0.000	1.000 $\pm$ 0.000	0.082 $\pm$ 0.000	
MIWAE	0.004 $\pm$ 0.000	1.000 $\pm$ 0.000	0.082 $\pm$ 0.000	
not-MIWAE	0.004 $\pm$ 0.000	0.997 $\pm$ 0.003	0.081 $\pm$ 0.000	
EGC	0.004 $\pm$ 0.000	1.000 $\pm$ 0.000	0.082 $\pm$ 0.000	
MaCoDE	0.005 $\pm$ 0.000	1.000 $\pm$ 0.000	0.083 $\pm$ 0.000	
Dataset				redwine
Model	Bias ↓	Coverage	Width ↓	
MICE	0.010 $\pm$ 0.001	0.873 $\pm$ 0.028	0.048 $\pm$ 0.000	
GAIN	0.011 $\pm$ 0.001	0.827 $\pm$ 0.032	0.048 $\pm$ 0.000	
missMDA	0.016 $\pm$ 0.002	0.755 $\pm$ 0.024	0.052 $\pm$ 0.000	
VAEAC	0.009 $\pm$ 0.001	0.900 $\pm$ 0.025	0.050 $\pm$ 0.000	
MIWAE	0.006 $\pm$ 0.001	0.955 $\pm$ 0.020	0.052 $\pm$ 0.000	
not-MIWAE	0.006 $\pm$ 0.001	0.973 $\pm$ 0.014	0.050 $\pm$ 0.000	
EGC	0.005 $\pm$ 0.001	1.000 $\pm$ 0.000	0.063 $\pm$ 0.001	
MaCoDE	0.006 $\pm$ 0.001	0.982 $\pm$ 0.018	0.053 $\pm$ 0.000	
Dataset				whitewine
Model	Bias ↓	Coverage	Width ↓	
MICE	0.008 $\pm$ 0.000	0.773 $\pm$ 0.020	0.028 $\pm$ 0.000	
GAIN	0.012 $\pm$ 0.002	0.736 $\pm$ 0.029	0.028 $\pm$ 0.000	
missMDA	0.010 $\pm$ 0.000	0.809 $\pm$ 0.021	0.033 $\pm$ 0.001	
VAEAC	0.009 $\pm$ 0.001	0.736 $\pm$ 0.029	0.028 $\pm$ 0.000	
MIWAE	0.006 $\pm$ 0.001	0.964 $\pm$ 0.020	0.033 $\pm$ 0.001	
not-MIWAE	0.006 $\pm$ 0.000	0.918 $\pm$ 0.025	0.031 $\pm$ 0.001	
EGC	0.004 $\pm$ 0.000	1.000 $\pm$ 0.000	0.055 $\pm$ 0.001	
MaCoDE	0.006 $\pm$ 0.001	0.855 $\pm$ 0.031	0.032 $\pm$ 0.001	

Table 14: **Q3**: Multiple imputation for each dataset under **MNARQ** at 0.3 missingness. The means and standard errors of the mean across 5 datasets and 10 repeated experiments are reported. ↓ denotes lower is better.

## References

- [1] Seunghwan An and Jong-June Jeon. Distributional learning of variational autoencoder: Application to synthetic data generation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [2] Jock Blackard. Coverttype. UCI Machine Learning Repository, 1998. DOI: <https://doi.org/10.24432/C50K5N>.
- [3] Yuri Burda, Roger Baker Grosse, and Ruslan Salakhutdinov. Importance weighted autoencoders. *ICLR*, abs/1509.00519, 2015.
- [4] Shoja’eddin Chenouri, Majid Mojirsheibani, and Zahra Montazeri. Empirical measures for incomplete data with applications. *Electronic Journal of Statistics*, 3:1021–1038, 2009.
- [5] Edward Choi, Siddharth Biswal, Bradley Malin, Jon Duke, Walter F Stewart, and Jimeng Sun. Generating multi-label discrete patient records using generative adversarial networks. In *Machine learning for healthcare conference*, pages 286–305. PMLR, 2017.
- [6] Paulo Cortez, A. Cerdeira, F. Almeida, T. Matos, and J. Reis. Wine Quality. UCI Machine Learning Repository, 2009. DOI: <https://doi.org/10.24432/C56S3T>.
- [7] Tianyu Du, Luca Melis, and Ting Wang. Remasker: Imputing tabular data with masked autoencoding. In *The Twelfth International Conference on Learning Representations*, 2024.
- [8] Manbir S Gulati and Paul F Roysdon. Tabmt: Generating tabular data with masked transformers. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [9] Niels Bruun Ipsen, Pierre-Alexandre Mattei, and Jes Frellsen. not- $\{miwae\}$ : Deep generative modelling with missing not at random data. In *International Conference on Learning Representations*, 2021.
- [10] Oleg Ivanov, Michael Figurnov, and Dmitry Vetrov. Variational autoencoder with arbitrary conditioning. In *International Conference on Learning Representations*, 2019.
- [11] Daniel Jarrett, Bogdan C Cebere, Tennison Liu, Alicia Curth, and Mihaela van der Schaar. HyperImpute: Generalized iterative imputation with automatic model selection. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 9916–9937. PMLR, 17–23 Jul 2022.
- [12] Julie Josse and François Husson. missMDA: A package for handling missing values in multivariate data analysis. *Journal of Statistical Software*, 70(1):1–31, 2016.
- [13] Jayoung Kim, Chaejeong Lee, and Noseong Park. STasy: Score-based tabular data synthesis. In *The Eleventh International Conference on Learning Representations*, 2023.
- [14] Akim Kotelnikov, Dmitry Baranchuk, Ivan Rubachev, and Artem Babenko. Tabddpm: Modelling tabular data with diffusion models. In *International Conference on Machine Learning*, pages 17564–17579. PMLR, 2023.
- [15] Chaejeong Lee, Jayoung Kim, and Noseong Park. Codi: co-evolving contrastive diffusion models for mixed-type tabular synthesis. In *Proceedings of the 40th International Conference on Machine Learning*, ICML’23. JMLR.org, 2023.
- [16] Jin Hyuk Lee and J. Charles Huber. Evaluation of multiple imputation with large proportions of missing data: How much is too much? *Iranian Journal of Public Health*, 50:1372 – 1380, 2021.
- [17] Steven Cheng-Xian Li, Bo Jiang, and Benjamin Marlin. Learning from incomplete data with generative adversarial networks. In *International Conference on Learning Representations*, 2019.
- [18] Volker Lohweg. Banknote Authentication. UCI Machine Learning Repository, 2013. DOI: <https://doi.org/10.24432/C55P57>.
- [19] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [20] Pierre-Alexandre Mattei and Jes Frellsen. Miwae: Deep generative modelling and imputation of incomplete data sets. In *International Conference on Machine Learning*, 2019.
- [21] Stan Matwin, Jordi Nin, Morvarid Sehatkar, and Tomasz Szapiro. A review of attribute disclosure control. In *Advanced Research in Data Privacy*, 2015.
- [22] Boris Muzellec, Julie Josse, Claire Boyer, and Marco Cuturi. Missing data imputation using optimal transport. In *International Conference on Machine Learning*, 2020.
- [23] Warwick Nash, Tracy Sellers, Simon Talbot, Andrew Cawthorn, and Wes Ford. Abalone. UCI Machine Learning Repository, 1995. DOI: <https://doi.org/10.24432/C55C7W>.



- [24] Alfredo Nazábal, Pablo M. Olmos, Zoubin Ghahramani, and Isabel Valera. Handling incomplete heterogeneous data using vaes. *Pattern Recognition*, 107:107501, 2020.
- [25] Noseong Park, Mahmoud Mohammadi, Kshitij Gorde, Sushil Jajodia, Hongkyu Park, and Youngmin Kim. Data synthesis based on generative adversarial networks. *Proc. VLDB Endow.*, 11:1071–1083, 2018.
- [26] Zhaozhi Qian, Rob Davis, and Mihaela van der Schaar. Synthcity: a benchmark framework for diverse use cases of tabular synthetic data. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023.
- [27] Donald B. Rubin and Nathaniel Schenker. Multiple imputation for interval estimation from simple random samples with ignorable nonresponse. *Journal of the American Statistical Association*, 81:366–374, 1986.
- [28] David Slate. Letter Recognition. UCI Machine Learning Repository, 1991. DOI: <https://doi.org/10.24432/C5ZP40>.
- [29] Latanya Sweeney. k-anonymity: A model for protecting privacy. *Int. J. Uncertain. Fuzziness Knowl. Based Syst.*, 10:557–570, 2002.
- [30] Stef van Buuren. Flexible imputation of missing data. 2012.
- [31] Stef van Buuren and Karin G. M. Groothuis-Oudshoorn. Mice: Multivariate imputation by chained equations in r. *Journal of Statistical Software*, 45:1–67, 2011.
- [32] William Wolberg, Olvi Mangasarian, Nick Street, , and W. Street. Breast Cancer Wisconsin (Diagnostic). UCI Machine Learning Repository, 1995. DOI: <https://doi.org/10.24432/C5DW2B>.
- [33] Lei Xu, Maria Skoularidou, Alfredo Cuesta-Infante, and Kalyan Veeramachaneni. Modeling tabular data using conditional gan. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [34] I-Cheng Yeh. Concrete Compressive Strength. UCI Machine Learning Repository, 2007. DOI: <https://doi.org/10.24432/C5PK67>.
- [35] Jinsung Yoon, James Jordon, and Mihaela van der Schaar. GAIN: Missing data imputation using generative adversarial nets. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 5689–5698. PMLR, 10–15 Jul 2018.
- [36] Hengrui Zhang, Jiani Zhang, Zhengyuan Shen, Balasubramaniam Srinivasan, Xiao Qin, Christos Faloutsos, Huzefa Rangwala, and George Karypis. Mixed-type tabular data synthesis with score-based diffusion in latent space. In *The Twelfth International Conference on Learning Representations*, 2024.
- [37] Nanhua Zhang, Chunyan Liu, Steven J Steiner, Richard B Colletti, Robert N. Baldassano, Shiran Chen, Stanley Cohen, Michael D. Kappelman, Shehzad Ahmed Saeed, Laurie S. Conklin, Richard Strauss, Sheri Volger, Eileen C. King, and Kim Hung Lo. Using multiple imputation of real-world data to estimate clinical remission in pediatric inflammatory bowel disease. *Journal of Comparative Effectiveness Research*, 12, 2023.
- [38] He Zhao, Ke Sun, Amir Dezfouli, and Edwin V. Bonilla. Transformed distribution matching for missing value imputation. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 42159–42186. PMLR, 23–29 Jul 2023.
- [39] Yuxuan Zhao, Alex Townsend, and Madeleine Udell. Probabilistic missing value imputation for mixed categorical and ordered data. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [40] Zilong Zhao, Aditya Kunar, Robert Birke, and Lydia Y. Chen. Ctab-gan: Effective table data synthesizing. In Vineeth N. Balasubramanian and Ivor Tsang, editors, *Proceedings of The 13th Asian Conference on Machine Learning*, volume 157 of *Proceedings of Machine Learning Research*, pages 97–112. PMLR, 17–19 Nov 2021.
- [41] Zilong Zhao, Aditya Kunar, Robert Birke, and Lydia Yiyu Chen. Ctab-gan+: Enhancing tabular data synthesis. *Frontiers in Big Data*, abs/2204.00401, 2023.