



HAL
open science

Under the Hood of Tabular Data Generation Models: the Strong Impact of Hyperparameter Tuning

G. Charbel N. Kindji, Lina Maria Rojas-Barahona, Elisa Fromont, Tanguy
Urvoy

► **To cite this version:**

G. Charbel N. Kindji, Lina Maria Rojas-Barahona, Elisa Fromont, Tanguy Urvoy. Under the Hood of Tabular Data Generation Models: the Strong Impact of Hyperparameter Tuning. 2024. hal-04612244v2

HAL Id: hal-04612244

<https://hal.science/hal-04612244v2>

Preprint submitted on 11 Jul 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Under the Hood of Tabular Data Generation Models: the Strong Impact of Hyperparameter Tuning

G. Charbel N. Kindji^{a,b}, Lina M. Rojas-Barahona^b, Elisa Fromont^a, Tanguy
Urvoy^b,

^a*Univ Rennes, IUF, Inria, CNRS, IRISA, Rennes, 35000, France*

^b*Orange Labs, Lannion, 22300, France*

Abstract

We investigate the impact of dataset-specific hyperparameter, feature encoding, and architecture tuning on five recent model families for tabular data generation through an extensive benchmark on 16 datasets. This study addresses the practical need for a unified evaluation of models that fully considers hyperparameter optimization. Additionally, we propose a reduced search space for each model that allows for quick optimization, achieving nearly equivalent performance at a significantly lower cost.

Our benchmark demonstrates that, for most models, large-scale dataset-specific tuning substantially improves performance compared to the original configurations. Furthermore, we confirm that diffusion-based models generally outperform other models on tabular data. However, this advantage is not significant when the entire tuning and training process is restricted to the same GPU budget for all models.

Keywords: Tabular data generation, Generative Models, Evaluation Metrics, Deep Learning, Hyperparameter Tuning, Neural Architecture Search

1. Introduction

The capability to develop generative models that generate realistic, safe, and useful tabular data is crucial for industries where this type of data is most prevalent. Among the direct applications of tabular data generation we can cite data privacy, imputation, oversampling, explainability or simulation [1, 2]. Another significant application of generative models is their ability to learn data representations that are valuable for pre-training and

fine-tuning for various downstream tasks [3, 4, 5]. However, generating high-quality tabular data presents several specific technical challenges that are not encountered with text or images [6]. First, tabular columns are encoded through heterogeneous data types with distributions that are often non-smooth with mixed (continuous/discrete) behaviours and various modalities. There can also be complex dependencies between columns, and the categorical features are often highly imbalanced. Finally, the wide variety of problems represented in tabular form makes it challenging to establish a universal data encoding and architecture suitable for pre-training across all scenarios.

To handle these challenges, many models have been proposed in the literature, spanning over very diverse approaches that can be probabilistic, GAN-based, diffusion-based, or LLM-based. These models are often evaluated on different datasets with inconsistent metrics, tuning and training budgets. Yet, the performance of the allegedly best tabular generation models published seem very unstable from one dataset to another: these models seem quite sensitive to the feature-encodings and hyperparameter-choices made by the authors.

The main purpose of this work is to study the impact of dataset-specific preprocessing and hyperparameter tuning on tabular data generation models. For each model we want to answer the following questions: (i) is it worth optimizing the hyperparameters/preprocessing specifically for each dataset? (ii) can we propose a reduced search space that fits well for all datasets? (iii) is there a clear trade-off between training/sampling costs, and synthetic data quality? Another goal of this work is to address the practical need for a unified model evaluation that explicitly incorporates hyperparameter optimization

We hence benchmarked 5 model families that are representative of the recent literature on 16 datasets with a strict 3-fold cross-validation procedure. The datasets were chosen based on their size, purpose and diversity. For each fold and dataset, we optimized the model’s hyperparameters, feature encoding, and architecture through hundreds of trials. We conducted two benchmarks: the first involved a large-scale, nearly unconstrained optimization, while the second employed a rapid optimization within a reduced search space, all within an equal compute budget. This work extends surveys like [2, 7] by covering recent diffusion-based models [8, 9] and by providing a larger scale benchmark.

As detailed in section 3.1, we consider multiple facets of the tabular data

generation problem. The key target of our study is the *realism of data* that we mainly assess through the *Classifier two-Sample Test* (C2ST), a metric which quantifies the ability of a classifier to discriminate real from synthetic data [10, 11]. We also consider other metrics such as, the *Machine Learning Efficacy* (ML-efficacy) for utility, and the *Distance to Closest Record Rate* (DCR-rate) for anonymity. We also measure or estimate the cost and carbon footprint at three stages: training (*i.e.* gradient descent), tuning, and sampling.

In Section 4, we analyze the results of these benchmarks and derive some intriguing insights about the models. For example, our experiments indicate that while diffusion-based models like TabSyn and TabDDPM [9] generally outperform other models when left unconstrained, they do not significantly surpass their simpler counterparts when tuning and training budgets are limited. This is because models without Transformers have a smaller memory footprint, allowing for more thorough optimization within the same GPU budget.

In Section 2.1, we quickly review the state-of-the-art with a particular focus on the competitors that were compared in our benchmarks. This review categorizes approaches as non-neural (e.g. SMOTE [12]), non-iterative neural models [11] (e.g. VAE and GANs), and finally, the most recent iterative generative models: diffusion [9] (e.g. TabDDPM and TabSyn) and auto-regressive ones (e.g. Great [13]). Section 3 presents the datasets and the evaluation metrics. Section 4 and Section 5 present respectively the large-scale and limited-budget benchmark results. We conclude in Section 6.

2. Benchmark Challengers

Tabular data generation is a booming research field which gives birth every month to a host of new data synthesis algorithms. In this section we quickly survey the existing families of tabular data generation methods with a specific focus on the models that we selected for our study. This includes models already covered in [2] and [7] as well as more recent diffusion-based models and LLM-based ones [14]. We chose models known for their strong performance, widespread usage, and availability of code that can be easily adapted for both architecture and hyperparameter tuning.

2.1. Tabular Data Generation Model’s Overview

Among the neural approaches for tabular data generation, we make a distinction between the “*push-forward*” models which directly map noise into data, and the iterative models which require a decoding phase. We also consider a few non-neural approaches as baselines.

2.1.1. Non-iterative Neural Approaches

The most popular non-iterative or “*push-forward*” neural networks for tabular data generation are Variational Auto-Encoders (VAE) [15] and Generative Adversarial Networks (GAN) [16]. A few papers also consider self-normalizing flows [17]. One of the simplest VAE architecture for tabular data is TVAE [18]. In its original implementation [19], it consists of a one-hot encoding for categorical variables coupled with a Gaussian Mixture Model normalization scheme (GMM) for continuous features. The encoder/decoder architecture is a simple stack of linear layers. Several variants have been proposed to improve from this baseline. In [20] the GMM normalization is replaced by a two-step training that first fits the marginals then fits the inter-dependencies. In [21, 11], several normalization schemes were tested as a replacement for the GMM normalization. Other variants of VAE use differentiable oblivious trees to ensure privacy [22]. In [23] the VAE is coupled with a Graph Neural Network (GNN). For our benchmark, we followed [11] and used a customized version of TVAE which allows for an optimized choice of the architecture and of the feature encoder.

The most popular method to generate tabular data is certainly adversarial training [24, 25, 26, 27, 28, 18, 29, 21, 30, 7, 31]. It would not be exaggerated to affirm that every exotic variant of GANs has been tested on tabular data generation, but the most successful architectures seem to be the ones based on Wasserstein GANs [32] such as CTGAN [18]. CTGAN is the base architecture that we selected for our benchmark. As mentioned in [21, 11], the feature encoding scheme is critical, especially for numerical features. For this reason we customized the CTGAN code as we did for TVAE, to allow for a choice of architecture and feature encoders.

2.2. Iterative models

It has been shown that iterative generative models (either auto-regressive or by diffusion) almost systematically outperform push forward models when it comes to raw text, sound, or image generation [33]. We confirm here that

tabular data do not escape this rule. However, this performance comes with a cost: the decoding phase is often slow and highly energy consuming.

2.2.1. Auto-regressive Language Models

After the recent breakthrough of large language models (LLM) [34], the usage of token-based language models to generate tabular data seems inevitable. Their main advantage against prior *ad-hoc* models is that they come without specific feature encodings for numerical or categorical columns: these values are directly fed to the model as raw sequences of tokens. Even if there is no clear agreement yet on how tabular instances should be serialized for LLMs, this general encoding ability opens the possibility for a pretraining/finetuning paradigm on heterogeneous tabular datasets [3, 5].

In a recent preprint survey [14], the authors tried to map the exuberant flow of preprint papers on tabular data and LLMs. Several of these preprints present only prompt engineering tricks that generate small tables, nonetheless some of the proposed models seem promising at a larger scale. To name a few of them, GReaT [13] (for "Generation of Realistic Tabular data") proposes to fine-tunes GPT-2 [35] for tabular data generation. RealTabFormer [36] extends GReaT to multiple tables with shared indexes, and TAPTAP [5] experiments a pretraining of GReaT on 450 open tabular datasets.

We made a few experiments with GReaT and its variants and confirmed the remarks of [9, 37], which state that they struggle to capture the joint probability distribution on datasets where the categorical values names do not carry semantic information. Given the prohibitive computational cost of LLMs, the large number of hyperparameters and serialization schemes to consider, as well as the problematic fact that some datasets are already covered (i.e. used in the training data) by the foundation models training sets, we decided to postpone their evaluation for a future work.

2.2.2. Diffusion models

Another recent impressive breakthrough in generative modeling, especially in image generation, was the apparition of diffusion models [38, 39, 40, 41, 42]. The transposition of these models to tabular data gave rise to powerful synthesizers: TableDiffusion [43], Stasy [44], CoDi [45], and TabDDPM [8]. In Tabsyn [9], a more recent proposal, the authors transposed the idea of [46, 47] to make use of a transformer-based VAE in order to embed the diffusion in a latent space. We selected both TabDDPM and Tabsyn for our benchmark.

2.2.3. Non-neural baselines

The most common statistical approaches are based on *copulas* and *Probabilistic Graphical Models* (PGMs). Copulas [48] are functions that join or “couple” multivariate distribution functions to their one-dimensional marginals. They have been widely adopted for tabular data generation because they allow modeling the marginals and the feature inter-dependencies separately [19, 49, 50, 51, 52]. Nevertheless, parametric copulas have been shown to perform poorly on high-dimension data synthesis problems [18, 11].

On the other hand PGMs can model variable dependencies in high-dimension spaces [53, 54, 7, 55, 56, 57]. This approach has been shown to be quite efficient in the data privacy community [58, 59, 60, 37]. However, it often requires a prior knowledge on the dependency graph because graph inference from data is inefficient in high dimension, especially if the sample size is small [61].

The fact that ensembles of trees remain state of the art for predictive tasks on tabular data [6] motivated some interesting attempts to mimic the neural generative approaches with decision trees. It gave rise to adversarial forests [31], Forest-Fow, and Forest-VP (a variance-preserving diffusion algorithm) [62]. We tested Forest-VP but we encountered a severe scalability issue: contrary to neural diffusion models where the noise level is combined with input as an auxiliary variable, the Forest-VP algorithm trains a different ensemble of trees for each level of noise.

Another simple, but quite efficient way to generate new tabular data is to interpolate between existing instances. This geometric “nearest neighbors” approach called *Synthetic Minority Over-sampling Technique* (SMOTE) [12, 63, 64] is frequently used to resample instances before training predictive models on unbalanced datasets. It proceeds by picking a random instance with a fixed target value and finding its k nearest neighbors. New data points are then generated by interpolation between these neighbors. Although very simple, this model is a solid baseline for tabular data generation as shown in [8]. As explained in section 3.2.1, we considered two variants: SMOTE and ucSMOTE (for unconditional SMOTE).

2.3. Selected challengers

We finally selected five families for our benchmark, four of them neural: TVAE, CTGAN, TabDDPM, and TabSyn and one non-neural: SMOTE (with its variant ucSMOTE). As developed in Section 3, for each of these algorithms we had to delve into the code in order to enable a large scale hyperparameters,

features encoding, and architecture optimization. In order to calibrate the metrics, we also reported the trivial baseline which consists in resampling directly the train set. It appears as “*Train Copy*” in result tables.

3. Experimental Settings

In this section, we first present the various evaluation metrics that we report from our experiments. We then present the optimization framework that we developed and we discuss the choices made to wrap the different challengers into this framework. We finally present the different datasets.

3.1. Evaluation Metrics

Our purpose is to assess the quality of tabular data generation through multiple facets. These facets can be summarized with four questions: (i) is synthetic data realistic? does it respect the original distribution’s traits? (ii) is it useful? *i.e.* can it be used to train machine learning models? (iii) do synthesis preserve training data anonymity? does it overfit? (iv) what are the model’s costs and CO2 impacts? Each of these questions is related to specific metrics.

The first and most important question is the realism of the generated data. We rely on *Classifier Two Sample Test*¹ (C2ST) [10] as a primary metric to address it. This metric is also the one that we used for hyperparameter optimization. It evaluates the performance of a classifier at discerning real data from synthetic data². To compute C2ST we use the same protocol as in [11] where the computed value is the mean ROC-AUC of XGBoost [65] over three folds. A C2ST around 1/2 means that XGBoost is unable to discern the test set from the generated set. A high C2ST means on the contrary that XGBoost is able to detect easily the fake data. The C2ST calculation procedure is summarized in Figure 1.

We also consider two other statistical metrics for data realism: *column-wise similarity* and *pair-wise correlation*. To compute these metrics we used the SDMetrics library [19]. The column-wise similarity measures how accurately the synthetic data captures the shape of each column distribution

¹Also mentioned as *Detection test* in SDMetrics <https://docs.sdv.dev/sdmetrics>

²This is similar to what is done in GANs but with a fresh dataset. An overfitted GAN can indeed have a poor C2ST on test set.

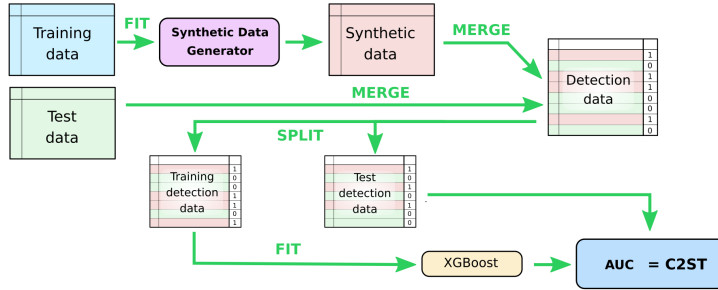


Figure 1: C2ST Metric calculation.

individually. It is reported as "Shape" in the result tables. Pair-wise correlation on the other hand captures how each column varies with each other. Pair-wise correlation is reported as "Pair" in the result tables. A naive generator that would assume independence of the columns might have a high "Shape" score but it would have a low "Pair" score.

The second important question is the utility of generated data. It is commonly measured by ML-Efficacy which evaluates the performance of a predictive model trained on synthetic data. To compute ML-Efficacy we use the same protocol as in [8] where CatBoost [66] is used as a predictor. We report the F1 score for classification tasks and the normalized R2 score for regression tasks. The procedure is summarized in Figure 2. It is important to evaluate the degradation of these scores against the ones obtained when training directly on real data (via *Train Copy*): a large degradation means a low utility.

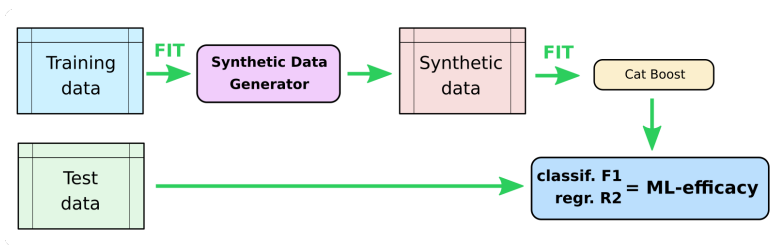


Figure 2: ML-Efficacy Metric calculation.

The third question is anonymity. We want to make sure that our generative model will not leak sensitive information by copying or over-fitting the training instances. To do so we measure respectively the minimum distances

of each generated instance to the train and test sets. The *Distance to Closest Record Rate* (DCR-Rate) counts the proportion of generated instances that are closer to train set than test set [67]. The procedure is summarized in Figure 3.

A synthetic dataset is considered safe if it has a DCR-Rate that is close to $1/2$. On the other hand, a plain copy of the train set, as *Train Copy* baseline does, would have a DCR-Rate of 1. This metric does not guarantee against all privacy breaches, but it provides a reasonable safeguard and ranking criterion for the models. It is also informative about overfitting, as an overfitted model would have samples that are systematically closer to the train set than a model with more generalization capabilities.

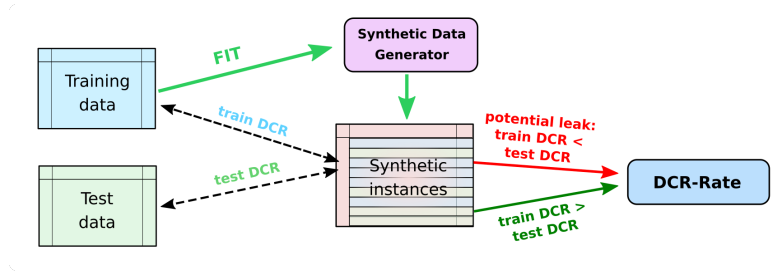


Figure 3: DCR-Rate Metric calculation.

The fourth question is the models’ costs and their carbon impacts. It requires an accurate estimate of three cost values: *time*, *energy consumption*, and *CO₂ impact*. Ideally, these values should be estimated for the three phases of (i) training (gradient descent) (ii) sampling, and (iii) hyperparameters search. For each dataset, each fold and each optimized model architecture, we ran the training and sampling phases on the exact same hardware and software architecture, a single Tesla v100 32 GB, and we measured accurately the cost values with the *CodeCarbon* library³. However, due to the massive nature of the experiments, we could not perform the whole hyperparameter search and training phases on such a uniform hardware and software architecture. We hence estimated the global search costs from the tuning logs by rescaling the training cost measures according to the effective number of steps performed

³<https://codecarbon.io/>

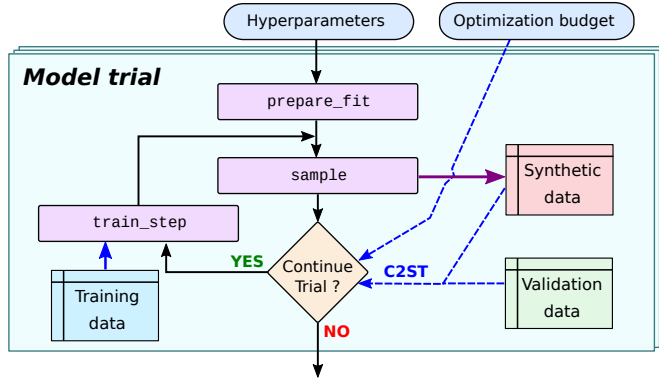


Figure 4: Hyperparameters trial optimization loop.

and the number of GPU used (c.f. equation (1)).

$$\text{total-gpu-cost} \simeq \sum_{t \in \text{trials}} \frac{\text{init-cost}_t + \text{avg-cost-per-step} \times \text{num-steps}_t}{\text{trials-per-gpu}} \quad (1)$$

This is slightly overestimated since: (i) *CodeCarbon* assumes a 100% GPU load while ours was roughly around 95%; and (ii) we measured the step costs on already optimized models which are usually slower because they often count more layers and parameters. Note that the number of trials that can be parallelized on a single GPU depends on the memory footprint of the model⁴.

3.2. Large Scale Optimization Framework and Implementation details

Providing a fair and reliable comparison of the different tabular generative models is a tough technical challenge. For this reason most existing benchmarks like [13, 9] only report the performance of models with their default hyperparameters. A few benchmarks like [8, 37] perform hyperparameters search for all models, but with a simple train/validation/test split, a reduced search space, and a small number of trials (usually from 20 to 50).

We wanted our experiment to be more extensive and more robust, so we decided to deploy it at a large scale on a super-computer equipped with

⁴For instance we could safely run 10 TVAE trials on the same v100 GPU while only 4 TabSyn’s trials could fit because of the VAE transformer’s footprint.

several nodes with 4 GPUs v100 32GB each⁵.

For this purpose, we used the *Ray tune* library [68] coupled with *Hyperopt* [69] and a median elimination scheduler [70] to optimize hyperparameters and architecture efficiently. Depending on the model’s memory footprints, each GPU could host from four to twelve concurrent tuning trials. For each dataset and each model we performed a strict 3-fold cross validation, which means that for each tuple (dataset, fold, model) we performed an extensive hyperparameters search with 300 trials (except for Tabsyn where we reduced this number to 100 for technical reasons explained in Section 3.2.2). We hence obtained a different optimized architecture for each (dataset, fold, model) tuple.

In [8] the parameters were optimized for ML-efficacy, in [37] the parameters were optimized for an equal combination of realism, utility and privacy. We chose to optimize for realism as in [11] through XGBoost-based C2ST metric (c.f. Section 3.1). To obtain reliable evaluations with variance estimates, for each dataset we evaluated the models by averaging all the metrics on the three folds test sets with five synthetic samples for each fold.

As described in Figure 4, in order to work with our framework each algorithm has to be wrapped into a generic *Synthesizer* class that provides three methods: *prepare_fit* which prepares the dataset and the model according to the hyperparameters, *train_step* which performs a training step roughly equivalent to one or a few epochs, and *sample* which generates synthetic data. After each train step, the model trial was evaluated and it was canceled out by early stopping or by the Ray-Tune scheduler if it performed too poorly or if the time budget was depleted.

3.2.1. The importance of feature encoding

Table 1 presents the encoding schemes that we used for the different benchmark challengers. As pointed out in [6, 21, 11], categorical variables are not the main weakness of neural networks but numerical feature encoding is critical. Most recent neural models use a Quantile-based numerical feature encoding and seem to work well with it. However, the original versions of TVAE and CTGAN rely on a specific Cluster-Based normalization [18]. We hence explored several encoding policies for these two models through hyper-

⁵Complementary experiments were also performed on a workstation with two RTX 4090.

parameters optimization (see Tables A.6 and A.7 in the appendix section).

Model	Num. Encoder	Cat. Encoder	Num. Target
TVAE_base	Cluster-Based ⁶	One hot	-
CTGAN_base	Cluster-Based ⁶	One hot	-
TVAE	<i>Optimized</i>	One hot	-
CTGAN	<i>Optimized</i>	One hot	-
TabDDPM	Quantile	One Hot	Standardize
TabSyn	Quantile	Embedding	-
SMOTE	-	One hot	Median cut
ucSMOTE	-	One hot	Dummy

Table 1: Encoding schemes applied to the benchmark challengers during data preprocessing. We experimented with various encoding schemes on TVAE and CTGAN during hyperparameter search.

We kept the native Cluster-Based⁶ encoder as described in [18], along with the ones proposed in [11], namely: prototype encoding (PTP) [11], piece-wise linear encoder (PLE) [71], continuously distributed residuals (CDF) [72, 11], hybrid (PLE_CDF) [11]. We also added some standard *scikit-learn* transformers: MinMaxScaler and QuantileTransformer.

Contrary to QuantileTransformer which maps values deterministically, the CDF encoding uses randomization to produce continuously distributed residuals even when the original distribution is discrete or partially discrete [72]. The PLE encoder [71] performs a feature binning and normalizes each numerical value depending on the bin it belongs to. The PLE_CDF applies a CDF to the output of a PLE encoding. The PTP encoding, inspired by prototypical networks [73], encodes the input as a weighted average of fixed prototypes.

3.2.2. Model-specific implementation details

Wrapping heterogeneous models within a *Synthesizer* class required some implementation choices from our part, and despite all our efforts to keep a fair comparison, these choices had some impact on the compute time and performance of the models.

⁶<https://docs.sdv.dev/rdt/transformers-glossary/numerical/clusterbasednormalizer>

The first issue is the discrepancy of the training step’s costs. The usual training step unit for most models is the *epoch* which corresponds to a single pass on all instances of the training set. However, depending on the hyperparameters some models like CTGAN perform both one pass through the generator and several passes through the discriminator at each training step. Other models such as TabDDPM are randomized and require several quick passes (almost one for each level of noise) for each instance. We hence had to caliber our wrapper’s `train_step` functions to perform a compute effort that is roughly equivalent to an epoch.

Another issue is the fact that TabSyn [9] combines two models and was not designed for hyperparameters tuning. According to the authors the model does not need hyperparameters tuning⁷. We decided to train a new transformer-based VAE for each trial because it hosts most of the parameters, compute-time and hyperparameters of TabSyn. Training a diffusion model on an unstable latent space would not be meaningful. We thus considered three technical options: (i) optimize first the VAE on a proxy metric (for instance the ability of its decoder to generate realistic data from a standard Gaussian), then optimize the denoiser in the latent space; (ii) wrap the VAE training steps into the *train_step* function and retrain a new denoiser from scratch at each step; (iii) wrap the VAE training phase into the *prepare_fit* function and loose the ability to prune its training steps. The first option is the cheapest and it is probably recommended for most practical applications, but it may be sub-optimal due to the proxy metric. The second option is extremely costly. We hence opted for the third option although it had a non-negligible cost. Indeed, we followed the recommendation of [9] to train the VAE through 4000 epochs which turns out to be huge knowing that most other models only utilized 400 epochs in our benchmark.

We also reduced the number of parallel trials per GPU because of the large memory footprint of the VAE’s transformers. As a consequence, for TabSyn we had both to reduce the number of trials to 100, and to work on a sample of the largest dataset (*Coverttype*) to get the results in a reasonable amount of time. It is worth noting that despite these handicaps, the optimized TabSyn remained better than its non-optimized version. To avoid this experimental

⁷Our experiment show however (c.f. Section 4), that even if the non-optimized TabSyn is quite good on most datasets, the hyperparameter tuning clearly improves the quality of the data it generates.

bias for the second experiment in Section 5.2, we constrained TabSyn’s VAE to use only 10 minutes for training. On *Adult* dataset, for instance, it resulted in roughly 560 epochs.

Finally, the last issue was the different ways the models deal with the target columns. A model conditioned on the target column may improve its ML-efficacy. TabDDPM and SMOTE implementations are natively conditioned on classification targets, while TVAE, CTGAN and TabSyn are not.

We did not modify TabDDPM, but we considered two variants in our experiment for SMOTE: the first, that we call SMOTE, follows the design of [8], it uses the train target distribution for classification datasets and a rough median split for regression targets (see Table 1). The second, that we call ucSMOTE (for unconditional SMOTE), adds a dummy target filled with zeros to the data before calling the SMOTE oversampling library. By doing so, all columns, including the original target, are considered equally.

3.3. Datasets

To evaluate the models, we picked datasets with various characteristics to assess their performances under different scenarios. The experiments were done with a 3-fold cross-validation procedure. Datasets were chosen in order to cover various sizes and dimensions, different types of tasks (regression, binary, and multi-class classification), various types of features (numerical, categorical, or mixed). We also added *Moons* a well known *scikitlearn* synthetic dataset. The complete list of datasets and their characteristics is presented in Table 2.

⁸<https://www.openml.org>

⁹<https://www.kaggle.com/datasets>

Name	Train	Validation	Test	Num	Categ.	Task
Abalone ⁸	2088	696	1393	7	2	Binclass
Adult ⁸	24420	8141	16281	6	9	Binclass
Bank Marketing ⁸	22605	7535	15071	7	10	Regression
Black Friday ⁸	83410	27804	55607	6	4	Regression
Bike Sharing ⁸	8689	2897	5793	9	4	Regression
Coverttype ⁸	290505	96836	193671	10	45	Multiclass
Cardio ⁹	34999	11667	23334	11	1	Binclass
Churn Modelling ⁹	4999	1667	3334	8	4	Binclass
Diamonds ⁸	26970	8990	17980	7	3	Regression
HELOC ⁹	5229	1743	3487	23	1	Binclass
Higgs ⁸	49024	16342	32684	28	1	Binclass
House 16H ⁸	11391	3798	7595	17	0	Regression
Insurance ⁹	669	223	446	4	3	Regression
King ⁹	10806	3602	7205	19	1	Regression
MiniBooNE ⁸	65031	21678	43355	50	1	Binclass
Two Moons	19999	6667	13334	2	1	Binclass

Table 2: List of datasets. Direct links to exact versions of datasets used can be found in Appendix B. The *Coverttype* dataset size was reduced for Tabsyn tuning as follows: 27500 in the training set, 18333 in the validation set, and 9167 in the test set.

4. Extensive Experiment Results

As mentioned in Section 3.1, the evaluation and comparison of tabular generative models is based mainly on four criteria: realism, usefulness, anonymity, and cost. After a global multi-criteria overview of the results, we study and compare the model’s behaviour according to each criterion individually.

4.1. Multi-criteria Overview

Figure 5 shows the average ranking over all datasets and folds of seven model variants over eight metrics (c.f. Section 3.1). To complete these rankings, Table 3 and Table 4 summarize respectively the quality metric and cost distributions among all datasets and folds. Recall that the best scores for the c2ST are around 0.50 (as it means poor AUC for the classifier at telling synthetic data apart from holdout data).

Overall, no model provides the best performance over all considered criteria. We observe, as expected, a strong correlation between energy and



Figure 5: Radar chart of the **extensive experiment** showing optimized model’s average ranking on all datasets according to C2ST, DCR-Rate, ML-Efficacy, *column-wise similarity* (“Shape”), *pair-wise correlation* (“Pair”), training and sampling time and energy costs. For all models except “Tabsyn_base”, which refers to the model using its default hyperparameters, the training costs include both tuning and gradient descent.

GPU time as well as a strong correlation among “quality metrics” (*i.e.* C2ST, *Shape*, and *Pair*). On the one hand, the diffusion models achieve the best performance in terms of quality metrics, especially the tuned version of Tabsyn. As shown in Table 3, this model has a median C2ST value of 0.64. However, Tabsyn is also one of the most expensive models in term of training costs (*Train-Energy* as well as *Train and Sample Times*), as it requires to train both a transformer-based VAE and a denoiser model for each dataset.

On the other hand, the SMOTE baselines obtain the poorest quality and privacy ranking with a high DCR-Rate. The median DCR-Rate score for SMOTE and ucSMOTE is at 0.97 which means that most of the samples from these models are very similar to the training set. However, they achieve strong utility in terms of ML-Efficacy with quartiles very close to *Train Copy*

and the same median value of 0.73. As expected for neighborhood-based algorithms, their training cost is negligible, but their deployment requires a neighborhood search which can be costly on large datasets. Finally the tuned neural push-forward models (CTGAN and TVAE) achieve mitigated results in term of both quality and utility but their deployment is clearly the cheapest. We also note that all neural models achieve reasonable results in term of privacy preservation. Although TabDDPM is clearly the slowest algorithm for deployment, it is fast at training and it obtains homogeneous results over all other metrics. We note that the tuned version of TabDDPM is performing better than the base version of Tabsyn.

For this extensive experiment we only limited mildly the time budget and the number of epochs. However, some models like Tabsyn and CTGAN consumed much more GPU time than others, especially TabDDPM which was very quick at performing an equivalent number of epochs (c.f. Section 3.2.2). It is also important to compare these algorithms with a fair allocation of GPU resource as we do in Section 5.

4.2. Detailed Analysis

In this section we study and compare the model’s behaviour according to each criterion taken individually. For each dataset we computed the performance metrics over 3 folds and 5 synthetic samples per fold to provide a stable central tendency and dispersion estimate. The full dataset-level results are provided in Appendix Table C.12. We summarize these results among all datasets and folds in Table 3 and Table 4. For quality metrics we compare the models through critical difference diagrams.

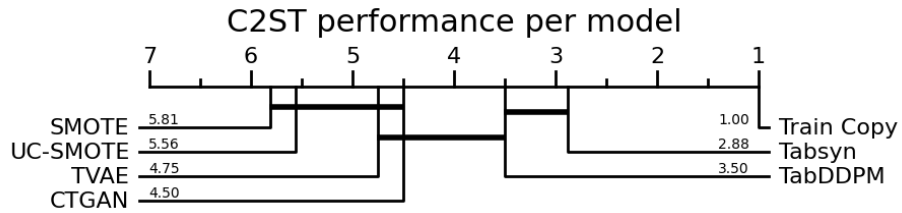
4.2.1. Are synthetic data realistic ?

We show in Figure 6 the *critical difference diagrams* [74] of all tuned models respectively for C2ST, pair-wise correlation and column-wise similarity. These diagrams were obtained by aggregating the ranks of the seven models over all datasets and folds. A thick horizontal line groups the set of models for which the pairwise “*no significant difference*” test hypothesis could not be rejected.

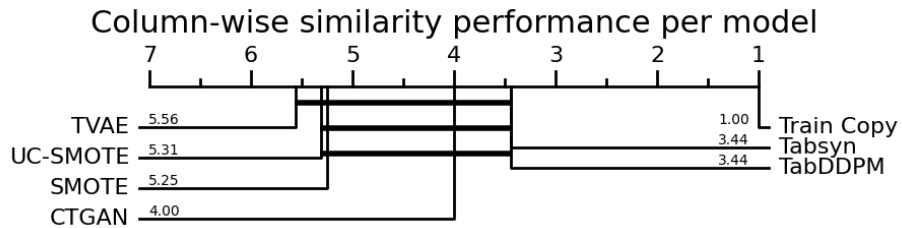
If we except the trivial *Train Copy* policy which is by construction the most realistic generator, we note that Tabsyn is significantly better in terms of C2ST than all other models except TabDDPM. On the other side, the SMOTE baselines are significantly worse than TabDDPM and Tabsyn. The absolute C2ST values in Table 3 corroborate these ranking results with three

Model	Percentiles	C2ST ↓	DCR-R ↓	ML-EF ↑	Shape ↑	Pair ↑
Train Copy	75%	0.50	1.00	0.90	0.99	0.99
	50%	0.50	1.00	0.73	0.99	0.98
	25%	0.50	1.00	0.63	0.99	0.91
TVAE	75%	0.88	0.65	0.79	0.96	0.95
	50%	0.81	0.63	0.71	0.95	0.92
	25%	0.74	0.61	0.50	0.94	0.85
TVAE_base	75%	1.00	0.62	0.70	0.92	0.90
	50%	0.98	0.61	0.65	0.91	0.80
	25%	0.96	0.60	0.44	0.87	0.74
CTGAN	75%	0.90	0.64	0.71	0.98	0.96
	50%	0.83	0.63	0.69	0.97	0.91
	25%	0.71	0.61	0.47	0.96	0.85
CTGAN_base	75%	1.00	0.62	0.63	0.92	0.93
	50%	0.99	0.60	0.47	0.88	0.83
	25%	0.93	0.60	0.30	0.87	0.76
TabDDPM	75%	0.78	0.64	0.80	0.98	0.98
	50%	0.67	0.62	0.69	0.98	0.93
	25%	0.63	0.61	0.58	0.95	0.84
TabDDPM_base	75%	0.86	0.64	0.75	0.99	0.96
	50%	0.77	0.62	0.68	0.98	0.92
	25%	0.65	0.61	0.44	0.94	0.74
TabSyn	75%	0.80	0.63	0.76	0.99	0.97
	50%	0.64	0.62	0.68	0.97	0.93
	25%	0.59	0.61	0.46	0.96	0.74
TabSyn_base	75%	0.86	0.64	0.75	0.98	0.97
	50%	0.71	0.62	0.57	0.97	0.95
	25%	0.63	0.61	0.29	0.95	0.88
SMOTE	75%	0.97	0.98	0.86	0.97	0.99
	50%	0.90	0.97	0.73	0.95	0.95
	25%	0.80	0.86	0.60	0.93	0.85
UC-SMOTE	75%	0.95	0.98	0.87	0.97	0.98
	50%	0.88	0.97	0.73	0.95	0.95
	25%	0.81	0.91	0.59	0.93	0.87

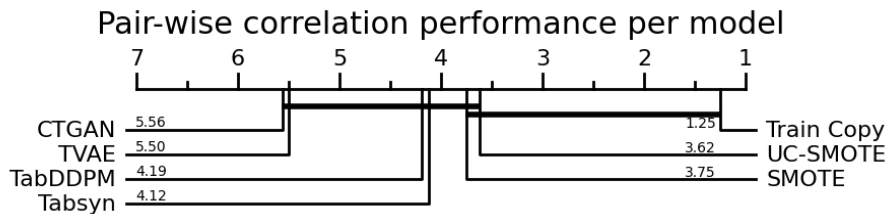
Table 3: Summary results from the *extensive* experiment and the *base* models (using default hyperparameters). We provide performance dispersion at the 25th, 50th, and 75th percentiles over all datasets and folds. To save space, we respectively shortened DCR-rate to DCR-R and ML-efficacy to ML-EF.



(a) C2ST



(b) Column-wise shape similarity



(c) Pair-wise correlation

Figure 6: Models ranking with critical difference diagrams respectively for C2ST, pair-wise correlations, and column-wise shape similarity metrics over all datasets.

sets of models: (i) diffusion-based models TabSyn and TabDDPM with a C2ST lower than 0.80 on most datasets and median around 0.65; (ii) push-forward models CTGAN and TVAE with a median C2ST around 0.82; (iii) SMOTE algorithms with a C2ST higher than 0.80 on most datasets.

The rankings obtained according to the *column-wise similarity* are broadly the same as the one obtained for C2ST. In addition to the rank shown in Figure 6(b), we can see that all the models obtain good scores (usually around 0.96) as compared to the *Train Copy* baseline (0.99). This result suggests that all the models succeed at capturing univariate distributions. For *pair-wise correlation*, we note surprisingly high values for SMOTE and ucSMOTE

baselines while the neural network models obtain broadly the same ranking as for C2ST but the gaps between models are less marked.

A side takeaway from this result is that XGBoost-based C2ST provides a stronger discriminative power than *column-wise similarity* and *pair-wise correlation* metrics.

4.2.2. Can the synthetic data be used to train a machine learning model ?

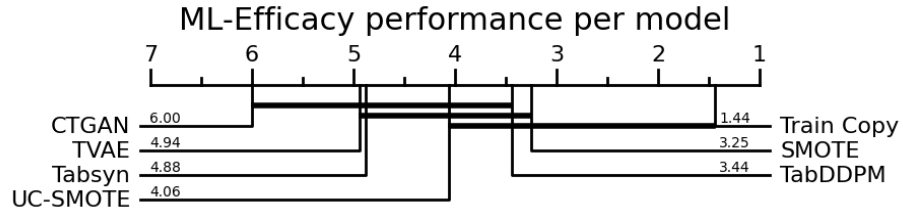


Figure 7: Models ranking with critical difference diagram for Catboost ML-Efficacy over all datasets.

According to the *machine learning efficacy* metric (ML-Efficacy), the most useful generators are the ones that are conditioned on their targets (namely SMOTE and TabDDPM) with median values respectively at 0.73 and 0.69 (against 0.73 for *Train Copy*).

TabDDPM outperforms both base and tuned versions of Tabsyn for this metric. It is also safer than SMOTE and it performs its training iterations faster than the other evaluated models. If ML-Efficacy is of importance, it is advisable to use this model. As expected, all models are far from the performance obtained on real data (*Train Copy*).

4.2.3. Does synthetic data preserve anonymity?

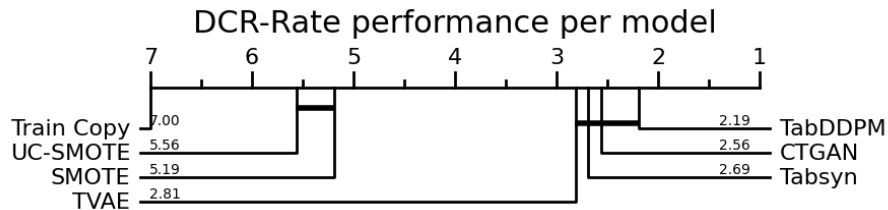


Figure 8: Models ranking with a critical difference diagram for the DCR-Rate metric over all datasets.

A data synthesizer that would only copy its training set would be of little value. If it generates new instances that are too close to its training set, it would obtain a good C2ST score, but it would be prone to over-fitting and it would leak private information from the training set. We assess the ability of a model to generate new data through the DCR-Rate metric (c.f. Section 3.1).

On Figure 8 we observe two significantly distinct groups of models. On the left-hand side a "leaky" group that contains both SMOTE and ucSMOTE, and on the right-hand side, a "safe" group that contains all neural algorithms. The poor performance of SMOTE is mainly due to the way it generates new data points by interpolating between existing ones. Therefore, these models cannot be considered safe concerning data protection. By taking a look at Table 3, the DCR-Rate of the two SMOTE variants is almost always above 0.86. On the other hand, if we exclude the tiny "Insurance" dataset where CTGAN and TVAE overfitted, the DCR-Rate values of the "safe" group are quite uniform around 0.62 and almost always below 0.65: these models can be considered safe.

4.2.4. *What are the models' costs?*

Model training, optimization, and data sampling have a cost and an environmental impact that varies greatly from one model to another. We hence measured and estimated time, energy consumption, and CO₂ impact for each model during three phases: (i) training (measured), (ii) hyperparameters search (estimated), and (iii) sampling (measured). These results are reported fully in Appendix D. We summarize the tuning and training process cost distributions over all datasets and fold in Table 4. As mentioned in Section 3.1, these values are estimated from the tuning logs by taking into account the effective number of training steps, the number of trials, and the average GPU resource usage per training step as reported in Table D.13.

With a median tuning time around 18 minutes as shown in Table 4, TabDDPM is the fastest neural model. As a result, it also consumes less energy and it has less emissions at the training stage. As shown in Figure 5, considering its other performance metrics, it is a suitable choice to achieve good results at a relatively low training cost.

As expected, Figure 5 shows that the push models CTGAN and TVAE are the fastest at sampling stage. We notice that although CTGAN achieves slightly better quality results than TVAE, it is also one of the costliest models at the training stage as shown in Table 4. In order to achieve reasonable performance while reducing training costs, TVAE is hence an option to consider

Model	Percentiles	Duration (HH:MM) ↓	Emission (Kg) ↓	Energy (kWh) ↓
TVAE	75%	01:00	1.05	15.57
	50%	00:24	0.42	6.22
	25%	00:12	0.21	3.08
CTGAN	75%	02:20	2.45	36.36
	50%	01:37	1.67	24.84
	25%	01:06	1.16	17.14
TabDDPM	75%	00:26	0.36	5.30
	50%	00:18	0.29	4.37
	25%	00:12	0.19	2.82
TabSyn	75%	03:25	1.92	22.65
	50%	01:59	1.03	14.20
	25%	01:29	0.72	9.09
SMOTE	75%	00:02	0.00	0.01
	50%	00:00	0.00	0.00
	25%	00:00	0.00	0.00
UC-SMOTE	75%	00:03	0.00	0.01
	50%	00:00	0.00	0.00
	25%	00:00	0.00	0.00

Table 4: Summary of the *costs* of the *extensive* experiment. We provide *costs* dispersion at the 25th, 50th, and 75th percentiles over all datasets and folds. The values are aggregated from the tuning costs estimated per dataset in Table D.14. More information about the estimation process can be found in Section 3.1.

prior to CTGAN.

The training and tuning of Tabsyn is one of the most demanding (with a median time of of 2 hours for tuning). In the end, however, it delivers the best performance in terms of quality metrics. This model also has the advantage of providing a set of default hyperparameters that can have a *reasonable* performance although, if we follow the author’s recommendation of 4000 VAE epochs, its training remains costly by comparison to other models. Indeed, Figure 5 show that, even if we consider the whole tuning+training pipeline TabDDPM and TVAE remain cheaper to train than Tabsyn_base.

In terms of sampling cost, TabDDPM is the worst-performing model. It takes longer than the other models (c.f. Figure 5) and hence consumes more energy with more CO₂ emissions at this step. Tabsyn reduces the number of denoising steps by using VAE embedding and linear noises to reduce its sampling time [9]. It hence achieves better performance at inference than TabDDPM.

The two baselines SMOTE and ucSMOTE being based on neighbourhood interpolation their training and tuning cost is negligible. However, as shown in Appendix Table C.12, their sampling process requiring a nearest neighbor search, it is slower on large datasets than push-forward models like TVAE or CTGAN.

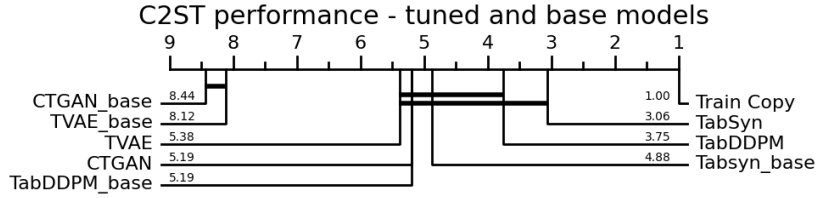
4.3. Is it worth optimizing the hyperparameters for all models ?

As mentioned in the previous section, even with the help of sophisticated search algorithms, hyperparameter tuning is costly and the performance-versus-tuning-budget curve is following a *diminishing returns law*. We were hence interested in comparing heavily tuned models against non-tuned models.

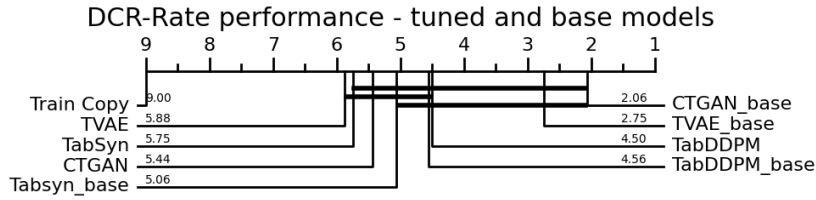
To assess this, we trained the neural models using the hyperparameters provided by the authors in the original papers. These results are detailed in Appendix E.

In Figure 9(a) we can see a, sometimes huge, C2ST performance improvement of all models when tuned. This improvement is statistically significant between optimized TVAE and CTGAN and their base versions. Looking at the absolute C2ST values in Table 3, it confirms that these models should not be used with their default hyperparameters.

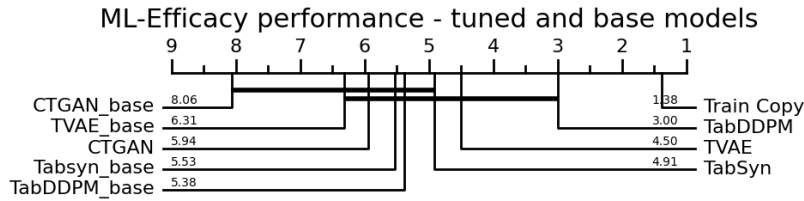
For TabDDPM and Tabsyn, we also notice a 10 points improvement on the median C2ST but this gap is not large enough to provide statistical guarantees.



(a) C2ST



(b) DCR-Rate



(c) ML-Efficacy

Figure 9: Models’ ranking with critical difference diagrams on C2ST, DCR-Rate, and ML-Efficacy metrics over all datasets: base models (using default hyperparameters) versus models obtained after an extensive hyperparameter tuning.

Although it was not the main target for tuning, we also observe an ML-Efficacy gain for all models. This gain is marked at the 25th percentile (*i.e.* for the hardest datasets).

We can therefore conclude that for all neural tabular generative models that we considered (including TabSyn), it is worth optimizing the hyperparameters specifically for each dataset if we want to improve performance. But the trade-off between the optimization cost and the performance gain is highly correlated to the size and design of the hyperparameter’s space (c.f. Appendix A). In the next section we propose a reduced search space and study the impact of a “light” hyperparameters optimization with a limited budget.

5. Limited-Budget Experiment Results

As underlined in Section 4.3, optimizing the hyperparameters for each dataset can significantly improve the quality of the generators. But a large-scale optimization like the one we performed is technically difficult, costly, and it has a non-negligible carbon impact (c.f. Table 4). Researchers and practitioners may hence be interested in reducing this cost without deteriorating too much the model’s performance.

In this section, we leverage the results of our extensive tuning experiment to: (i) suggest reduced search spaces achieving reasonable performance at a much lower cost; (ii) assess and compare the performance of the models when tuned and trained with the same limited budget. By comparing the model’s performance after this ”light” tuning/training with our previous results (respectively with heavy tuning or without tuning at all), we gain new insights into the models.

5.1. Hyperparameters Search Space Reduction

We carried out this experiment for the most hyperparameter-heavy models, namely: TVAE, CTGAN, TabDDPM, and TabSyn. To reduce their search space, we independently considered each hyperparameter/architecture configuration variable and kept only the values that were the most frequently selected during the large-scale tuning phase. For discrete variables we kept the 80% most frequent values, and for continuous variables we kept the value ranging between the 10th and 90th percentiles.

For instance, on CTGAN the large-scale tuning included six encoder options: CDF, PLE_CDF, PTP, Quantile, MinMax, and CBN. But only CDF and PLE_CF were selected on most datasets. We could also drastically reduce TabSyn VAE’s learning rate range from $(10^{-5}, 10^{-2})$ to $(10^{-3}, 7 \cdot 10^{-3})$. We present all these reduced search spaces in Appendix A.

5.2. Putting the Reduced Search Spaces to the Test

To evaluate the reduced search spaces, we ran a new tuning on all datasets with only 50 trials and a strict limit of 20 minutes per trial. Since TabSyn’s training is done in two steps, we allocated 10 minutes for the VAE training and 10 minutes for the denoiser. Each hyperparameter search was performed with the same 3-fold splits and the same methodology as in the extensive experiment.

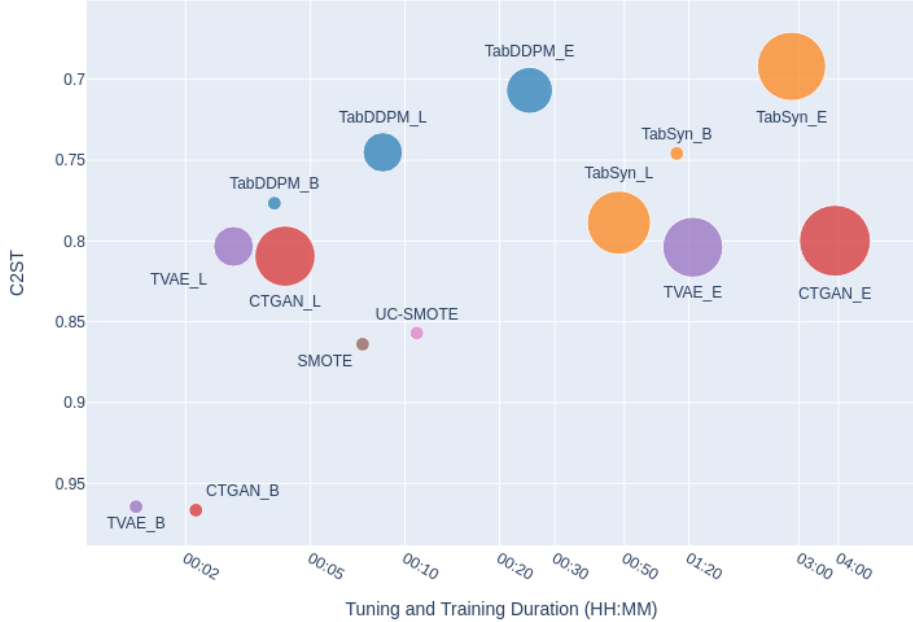


Figure 10: Average C2ST performance of the models under various setups. We show the performance of the base models (appended with ”_B”), the models obtained from the light experiment described in this section (appended with ”_L”), and the ones from the extensive experiment (appended with ”_E”). The C2ST axis is reversed to show the best models at the top. The dot diameter indicates the complexity of the search space. The duration axis and dot diameters are log-scaled for better visualization.

5.2.1. A Cost-Performance Trade-off

In Figure 10 we summarize the performance of the models after this light hyperparameter search against the ones obtained with the base models and the ones obtained after the extensive tuning of Section 3.2. On the x-axis we display the total GPU-time (search+optimization) needed to obtain the corresponding model and on the y-axis we show the C2ST performance. The diameters of the dots represents the number of possible configurations of the hyperparameter search space as described in Appendix A.

If we compare against the non-tuned base models, the C2ST performance

of all models except TabSyn is clearly improved by the light hyperparameter tuning. For CTGAN, and TVAE the gain is huge: after only a few minutes of tuning with the reduced search space we reach the same performance as the one obtained after hours of heavy tuning. For TabDDPM, we also improved the base performance by running the light hyperparameter tuning. However, we did not reach the same performance as we did with an extensive search.

In the light experiment, the TabSyn’s performance degradation against the base model reveals the importance of a well-trained VAE which leads to better samples from the diffusion process. Indeed, on Adult dataset for instance, within its 10 minutes budget, it was trained with less than 600 epochs against 4000 for the base model. This suggests to increase the number of epochs or the training time for TabSyn’s VAE in further experiments.

Finally, as shown in Figure 10, the costs induced by the heavily optimized TabSyn model are high compared to the base one. The performance gain and the gap in cost should be considered depending on the task and constraints as the base model already delivers *decent* performance.

It is worth noting that, due to its quick implementation, TabDDPM performed sometimes more training steps within its 20 minutes time budget on the light experiment than it did on the extensive search where the number of steps was bounded. But on the other hand, its performance was also affected by the reduced number of trials (50 instead of 300).

5.2.2. Multi-criteria Analysis of Light Tuning

In Figure 11 we summarize the model’s relative performances according to the five metrics for the light tuning experiment. To complete this point of view, Table 5 summarizes the quality metric distributions among all datasets and folds. The full results are detailed in Appendix F.

A first remark about Figure 11 is that it is more balanced than Figure 5: with a fair and limited GPU budget allocation, the models tends to perform similarly among all criteria.

There is no more clear leading model although TabDDPM obtains the best C2ST performance and TabSyn slightly dominates in terms of privacy (DCR-Rate). We however have a safe median DCR-Rate score of about 0.63 for this experiment on all neural models (against 0.62 for the extensive experiment). For ML-Efficacy, TabDDPM remains slightly the best but there is no clear leader either.

The reduced search space that we provide is hence a good starting point to perform a quick dataset-specific hyperparameters optimization that will

Model	Percentiles	C2ST ↓	DCR-R ↓	ML-EF ↑	Shape ↑	Pair ↑
Train Copy	75%	0.50	1.00	0.90	0.99	0.99
	50%	0.50	1.00	0.73	0.99	0.98
	25%	0.50	1.00	0.63	0.99	0.91
TVAE	75%	0.92	0.65	0.80	0.96	0.95
	50%	0.80	0.63	0.70	0.95	0.93
	25%	0.70	0.61	0.51	0.94	0.87
CTGAN	75%	0.89	0.64	0.71	0.98	0.96
	50%	0.85	0.63	0.68	0.97	0.92
	25%	0.70	0.62	0.45	0.95	0.85
TabDDPM	75%	0.83	0.67	0.83	0.97	0.95
	50%	0.72	0.63	0.69	0.96	0.91
	25%	0.63	0.62	0.50	0.95	0.81
TabSyn	75%	0.88	0.63	0.71	0.97	0.96
	50%	0.81	0.62	0.65	0.95	0.93
	25%	0.72	0.61	0.34	0.94	0.82

Table 5: Results from the *light* experiment: performance dispersion at the 25th, 50th, and 75th percentiles across all models. These results are obtained by aggregating the scores across all datasets. To save space, we respectively shortened DCR-rate to DCR-R and ML-efficacy to ML-EF.

fit on a medium-size workstation.

The relatively balanced results of this new experiment confirms that the superiority of a tabular data generator is not only due to its model but also to the whole tuning and training compute effort.

Overall, TabDDPM provides a good balance between realism, privacy, utility, and cost. TVAE is a viable alternative if the utility constraint can be relaxed. It can be tuned using the light experiment search space to quickly achieve good performance. If resources are available and cost is not a priority, it is advisable to tune TabSyn, which achieves good realism results but at a higher cost. Finally, for quick results when privacy is not a concern and when the main focus is utility, SMOTE is the recommended approach.

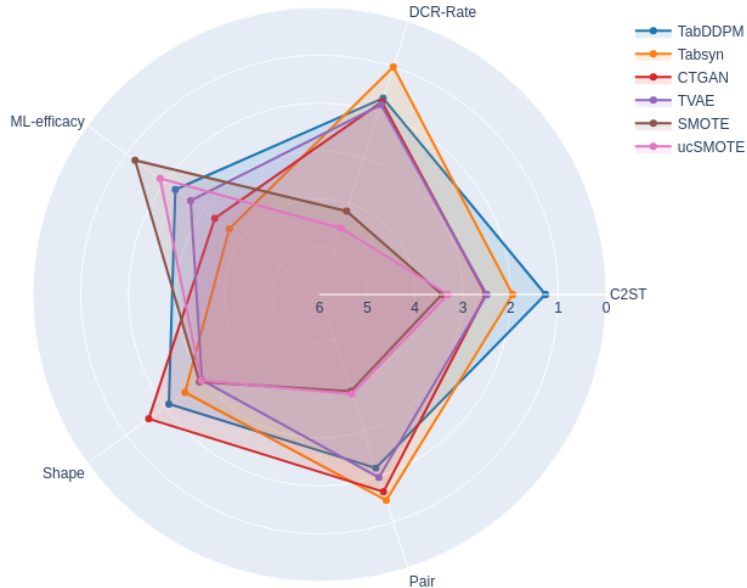


Figure 11: Radar chart of the optimized models under the **light experiment** according to C2ST, DCR-Rate, ML-Efficacy, *column-wise similarity* (named "Shape"), and *pair-wise correlation* (named "Pair"). We take the ranks of the models to obtain normalized metrics. The models were trained with a maximum time budget of 20 minutes and tuned with 50 trials per fold.

6. Conclusion

We benchmarked extensively recent tabular data generation models on 16 datasets with a strict 3-fold cross-validation procedure. We first performed both a large-scale tuning experiment on a super computer from which we derived a reduced search-space. We then performed a quick tuning experiment that fits on a medium-size workstation. Leveraging these experiments we were able to provide several insights on the models while answering to three technical questions: (i) is it worth optimizing the hyper-parameters/preprocessing specifically for each dataset? (ii) can we propose a reduced search space that fits well for all datasets? (iii) is there a clear trade-off between training/sampling costs, and synthetic data quality?

For the two first questions, Figure 10 is certainly the best summary: most models, including Tabsyn, benefit greatly from a dataset-specific tuning. But the whole tuning process is costly (time, money, energy, CO₂ emissions) and

there is clearly a "diminishing return" effect. We conclude that, even for TabSyn, a quick dataset-specific model tuning on a carefully designed search space as we did in Section 5 is enough to get most of the performance at the scale of a medium-size workstation.

Regarding the trade-off question on the multiple considered criteria: realism, privacy, utility and costs, Figure 5 and Figure 11 confirm that if we do not limit the compute power, the two diffusion-based models TabDDPM and TabSyn are the most recommended solutions for tabular data generation. But with a fairly limited compute power, the performance gaps between models become quite narrow.

7. Ackowlegment

This work was granted access to the HPC resources of IDRIS under the allocations 2023-AD011014381, 2023-AD011011407R3, and 2023-AD011012220R2 made by GENCI.

References

- [1] D. Libes, D. Lechevalier, S. Jain, Issues in synthetic data generation for advanced manufacturing, in: 2017 IEEE International Conference on Big Data (Big Data), IEEE, 2017, pp. 1746–1754.
- [2] V. Borisov, T. Leemann, K. Seßler, J. Haug, M. Pawelczyk, G. Kasneci, Deep neural networks and tabular data: A survey, *IEEE Transactions on Neural Networks and Learning Systems* (2022).
- [3] L. Zhang, S. Zhang, K. Balog, Table2vec: Neural word and entity embeddings for table population and retrieval, in: Proceedings of the 42nd international ACM SIGIR conference on research and development in information retrieval, 2019, pp. 1029–1032.
- [4] S. Chitlangia, A. Muralidhar, R. Agarwal, Self supervised pre-training for large scale tabular data, in: *NeurIPS 2022 Workshop on Table Representation Learning*, 2022.
- [5] T. Zhang, S. Wang, S. Yan, L. Jian, Q. Liu, Generative table pre-training empowers models for tabular prediction, in: H. Bouamor,

- J. Pino, K. Bali (Eds.), Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Singapore, 2023, pp. 14836–14854. URL: <https://aclanthology.org/2023.emnlp-main.917>. doi:10.18653/v1/2023.emnlp-main.917.
- [6] L. Grinsztajn, E. Oyallon, G. Varoquaux, Why do tree-based models still outperform deep learning on typical tabular data?, in: Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track, 2022. URL: https://openreview.net/forum?id=Fp7__phQszn.
- [7] J. Fonseca, F. Bacao, Tabular and latent space synthetic data generation: a literature review, *Journal of Big Data* 10 (2023) 115.
- [8] A. Kotelnikov, D. Baranchuk, I. Rubachev, A. Babenko, Tabddpm: Modelling tabular data with diffusion models, in: International Conference on Machine Learning, PMLR, 2023, pp. 17564–17579.
- [9] H. Zhang, J. Zhang, Z. Shen, B. Srinivasan, X. Qin, C. Faloutsos, H. Rangwala, G. Karypis, Mixed-type tabular data synthesis with score-based diffusion in latent space, in: The Twelfth International Conference on Learning Representations, 2024.
- [10] D. Lopez-Paz, M. Oquab, Revisiting classifier two-sample tests, in: International Conference on Learning Representations, 2016.
- [11] E. H. Zein, T. Urvoy, Tabular data generation: Can we fool XGBoost?, in: NeurIPS 2022 First Table Representation Workshop, 2022.
- [12] N. V. Chawla, K. W. Bowyer, L. O. Hall, W. P. Kegelmeyer, Smote: synthetic minority over-sampling technique, *Journal of artificial intelligence research* 16 (2002) 321–357.
- [13] V. Borisov, K. Sessler, T. Leemann, M. Pawelczyk, G. Kasneci, Language models are realistic tabular data generators, in: The Eleventh International Conference on Learning Representations (ICLR), 2023.
- [14] X. Fang, W. Xu, F. A. Tan, J. Zhang, Z. Hu, Y. Qi, S. Nickleach, D. Socolinsky, S. Sengamedu, C. Faloutsos, Large language models on tabular data—a survey, arXiv preprint arXiv:2402.17944 (2024).

- [15] D. P. Kingma, M. Welling, Auto-encoding variational bayes, *stat* 1050 (2014) 1.
- [16] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative adversarial networks, *Communications of the ACM* 63 (2020) 139–144.
- [17] E. G. Tabak, C. V. Turner, A family of nonparametric density estimation algorithms, *Communications on Pure and Applied Mathematics* 66 (2013) 145–164.
- [18] L. Xu, M. Skoularidou, A. Cuesta-Infante, K. Veeramachaneni, Modeling tabular data using conditional gan, in: H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, R. Garnett (Eds.), *Advances in Neural Information Processing Systems (NeurIPS)*, volume 32, Curran Associates, Inc., 2019.
- [19] N. Patki, R. Wedge, K. Veeramachaneni, The synthetic data vault, in: *IEEE International Conference on Data Science and Advanced Analytics (DSAA)*, 2016, pp. 399–410. doi:10.1109/DSAA.2016.49.
- [20] C. Ma, S. Tschitschek, R. Turner, J. M. Hernández-Lobato, C. Zhang, Vaem: a deep generative model for heterogeneous mixed type data, *Advances in Neural Information Processing Systems* 33 (2020) 11237–11247.
- [21] Z. Zhao, A. Kurnar, R. Birke, L. Y. Chen, Ctab-gan: Effective table data synthesizing, in: *Asian Conference on Machine Learning*, PMLR, 2021, pp. 97–112.
- [22] L. V. H. Vardhan, S. Kok, Generating privacy-preserving synthetic tabular data using oblivious variational autoencoders, in: *Proceedings of the Workshop on Economics of Privacy and Data Labor at the 37 th International Conference on Machine Learning*, 2020.
- [23] T. Liu, Z. Qian, J. Berrevoets, M. van der Schaar, GOGGLE: Generative modelling for tabular data by learning relational structure, in: *The Eleventh International Conference on Learning Representations*, 2023. URL: <https://openreview.net/forum?id=fPVRcJqspu>.

- [24] J. Jordon, J. Yoon, M. van der Schaar, Pate-gan: Generating synthetic data with differential privacy guarantees, in: International Conference on Learning Representations, 2018. URL: <https://api.semanticscholar.org/CorpusID:53342261>.
- [25] N. Park, M. Mohammadi, K. Gorde, S. Jajodia, H. Park, Y. Kim, Data synthesis based on generative adversarial networks, Proceedings of the VLDB Endowment 11 (2018).
- [26] R. D. Camino, C. Hammerschmidt, et al., Generating multi-categorical samples with generative adversarial networks, in: ICML 2018 workshop on Theoretical Foundations and Applications of Deep Generative Models, 2018.
- [27] M. K. Baowaly, C.-C. Lin, C.-L. Liu, K.-T. Chen, Synthesizing electronic health records using improved generative adversarial networks, Journal of the American Medical Informatics Association: JAMIA 26 (2019) 228–241.
- [28] H. Chen, S. Jajodia, J. Liu, N. Park, V. Sokolov, V. Subrahmanian, Faketables: Using gans to generate functional dependency preserving tables with bounded real data, in: 28th International Joint Conference on Artificial Intelligence, IJCAI 2019, International Joint Conferences on Artificial Intelligence, 2019, pp. 2074–2080.
- [29] A. Koivu, M. Sairanen, A. Airola, T. Pahikkala, Synthetic minority oversampling of vital statistics data with generative adversarial networks, Journal of the American Medical Informatics Association 27 (2020) 1667–1674.
- [30] R. Sauber-Cole, T. M. Khoshgoftaar, The use of generative adversarial networks to alleviate class imbalance in tabular data: a survey, Journal of Big Data 9 (2022) 98.
- [31] D. S. Watson, K. Blesch, J. Kapar, M. N. Wright, Adversarial random forests for density estimation and generative modeling, in: International Conference on Artificial Intelligence and Statistics, PMLR, 2023, pp. 5357–5375.

- [32] M. Arjovsky, S. Chintala, L. Bottou, Wasserstein generative adversarial networks, in: International conference on machine learning, PMLR, 2017, pp. 214–223.
- [33] S. Bond-Taylor, A. Leach, Y. Long, C. G. Willcocks, Deep generative modelling: A comparative review of vaes, gans, normalizing flows, energy-based and autoregressive models, IEEE Transactions on Pattern Analysis and Machine Intelligence 44 (2022) 7327–7347. doi:10.1109/TPAMI.2021.3116668.
- [34] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Gray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell, P. Welinder, P. Christiano, J. Leike, R. Lowe, Training language models to follow instructions with human feedback, in: A. H. Oh, A. Agarwal, D. Belgrave, K. Cho (Eds.), Advances in Neural Information Processing Systems, 2022. URL: <https://openreview.net/forum?id=TG8KACxEON>.
- [35] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, J. Dean, S. Ghemawat, Language models are unsupervised multitask learners, in: OSDI’04: Sixth Symposium on Operating System Design and Implementation, 2004, pp. 137–150.
- [36] A. V. Solatorio, O. Dupriez, Realtabformer: Generating realistic relational and tabular data using transformers, arXiv preprint arXiv:2302.02041 (2023).
- [37] Y. Du, N. Li, Towards principled assessment of tabular data synthesis algorithms, arXiv preprint arXiv:2402.06806 (2024).
- [38] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, S. Ganguli, Deep unsupervised learning using nonequilibrium thermodynamics, in: International Conference on Machine Learning, PMLR, 2015, pp. 2256–2265.
- [39] J. Ho, A. Jain, P. Abbeel, Denoising diffusion probabilistic models, Advances in neural information processing systems 33 (2020) 6840–6851.
- [40] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, B. Poole, Score-based generative modeling through stochastic differential equations, in: International Conference on Learning Representations, 2021. URL: <https://openreview.net/forum?id=PXTIG12RRHS>.

- [41] P. Dhariwal, A. Nichol, Diffusion models beat gans on image synthesis, *Advances in neural information processing systems* 34 (2021) 8780–8794.
- [42] T. Karras, M. Aittala, T. Aila, S. Laine, Elucidating the design space of diffusion-based generative models, *Advances in Neural Information Processing Systems* 35 (2022) 26565–26577.
- [43] G. Truda, Generating tabular datasets under differential privacy, *arXiv preprint arXiv:2308.14784* (2023). TableDiffusion.
- [44] J. Kim, C. Lee, N. Park, STaSy: Score-based tabular data synthesis (2022).
- [45] C. Lee, J. Kim, N. 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.
- [46] A. Vahdat, K. Kreis, J. Kautz, Score-based generative modeling in latent space, *Advances in neural information processing systems* 34 (2021) 11287–11302.
- [47] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, B. Ommer, High-resolution image synthesis with latent diffusion models, in: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2022*, pp. 10684–10695.
- [48] R. B. Nelsen, *An Introduction to Copulas*, second ed., Springer, New York, NY, USA, 2006.
- [49] Z. Li, Y. Zhao, J. Fu, Sync: A copula based framework for generating synthetic data from aggregated sources, in: *2020 International Conference on Data Mining Workshops (ICDMW)*, IEEE, 2020, pp. 571–578.
- [50] S. Kamthe, S. Assefa, M. Deisenroth, Copula flows for synthetic data generation, *arXiv preprint arXiv:2101.00598* (2021).
- [51] D. Meyer, T. Nagler, Synthia: multidimensional synthetic data generation in python, *Journal of Open Source Software* 6 (2021) 2863. URL: <https://doi.org/10.21105/joss.02863>. doi:10.21105/joss.02863.

- [52] D. Meyer, T. Nagler, R. J. Hogan, Copula-based synthetic data augmentation for machine-learning emulators, *Geoscientific Model Development* 14 (2021) 5205–5215. URL: <https://doi.org/10.5194/gmd-14-5205-2021>. doi:10.5194/gmd-14-5205-2021.
- [53] R. Mckenna, D. Sheldon, G. Miklau, Graphical-model based estimation and inference for differential privacy, in: K. Chaudhuri, R. Salakhutdinov (Eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, PMLR, 2019, pp. 4435–4444. URL: <https://proceedings.mlr.press/v97/mckenna19a.html>.
- [54] Z. Zhang, T. Wang, N. Li, J. Honorio, M. Backes, S. He, J. Chen, Y. Zhang, {PrivSyn}: Differentially private data synthesis, in: *30th USENIX Security Symposium (USENIX Security 21)*, 2021, pp. 929–946.
- [55] J. Dahmen, D. J. Cook, Synsys: A synthetic data generation system for healthcare applications, *Sensors (Basel, Switzerland)* 19 (2019). URL: <https://api.semanticscholar.org/CorpusID:75136844>.
- [56] B. Tang, H. He, Kerneladasyn: Kernel based adaptive synthetic data generation for imbalanced learning, in: *2015 IEEE Congress on Evolutionary Computation (CEC)*, 2015, pp. 664–671. doi:10.1109/CEC.2015.7256954.
- [57] M. Baak, S. Brugman, I. Fridman Rojas, L. Dalmeida, R. E.Q. Urlus, J.-B. Oger, Synthsonic: Fast, probabilistic modeling and synthesis of tabular data, in: G. Camps-Valls, F. J. R. Ruiz, I. Valera (Eds.), *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, PMLR, 2022, pp. 4747–4763. URL: <https://proceedings.mlr.press/v151/baak22a.html>.
- [58] R. McKenna, G. Miklau, D. Sheldon, Winning the nist contest: A scalable and general approach to differentially private synthetic data, *Journal of Privacy and Confidentiality* 11 (2021). URL: <https://journalprivacyconfidentiality.org/index.php/jpc/article/view/778>. doi:10.29012/jpc.778.

- [59] Y. Tao, R. McKenna, M. Hay, A. Machanavajjhala, G. Miklau, Benchmarking differentially private synthetic data generation algorithms (2022).
- [60] Y. Hu, F. Wu, Q. Li, Y. Long, G. Garrido, C. Ge, B. Ding, D. Forsyth, B. Li, D. Song, Sok: Privacy-preserving data synthesis, in: 2024 IEEE Symposium on Security and Privacy (SP), IEEE Computer Society, 2023, pp. 2–2.
- [61] J. Jewson, L. Li, L. Battaglia, S. Hansen, D. Rossell, P. Zwiernik, Graphical model inference with external network data, cemmap working paper CWP20/22, London, 2022. URL: <https://hdl.handle.net/10419/272832>. doi:10.47004/wp.cem.2022.2022.
- [62] A. Jolicoeur-Martineau, K. Fatras, T. Kachman, Generating and imputing tabular data via diffusion and flow-based gradient-boosted trees, in: S. Dasgupta, S. Mandt, Y. Li (Eds.), Proceedings of The 27th International Conference on Artificial Intelligence and Statistics, volume 238 of *Proceedings of Machine Learning Research*, PMLR, 2024, pp. 1288–1296. URL: <https://proceedings.mlr.press/v238/jolicoeur-martineau24a.html>.
- [63] L. Torgo, R. Ribeiro, B. Pfahringer, P. Branco, Smote for regression, volume 8154, 2013, pp. 378–389. doi:10.1007/978-3-642-40669-0_33.
- [64] L. Camacho, G. Douzas, F. Bacao, Geometric smote for regression, *Expert Systems with Applications* 193 (2022) 1–8. doi:10.1016/j.eswa.2021.116387, camacho, L., Douzas, G., & Bacao, F. (2022). Geometric SMOTE for regression. *Expert Systems with Applications*, 193(May), 1–8. [116387]. <https://doi.org/10.1016/j.eswa.2021.116387>.
- [65] T. Chen, C. Guestrin, Xgboost: A scalable tree boosting system, in: Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining, 2016, pp. 785–794.
- [66] L. Prokhorenkova, G. Gusev, A. Vorobev, A. V. Dorogush, A. Gulin, Catboost: unbiased boosting with categorical features, *Advances in neural information processing systems* 31 (2018).
- [67] M. Platzner, T. Reutterer, Holdout-based empirical assessment of mixed-type synthetic data, *Frontiers in big Data* 4 (2021) 679939.

- [68] R. Liaw, E. Liang, R. Nishihara, P. Moritz, J. E. Gonzalez, I. Stoica, Tune: A research platform for distributed model selection and training, arXiv preprint arXiv:1807.05118 (2018).
- [69] J. Bergstra, D. Yamins, D. Cox, Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures, in: International conference on machine learning, PMLR, 2013, pp. 115–123.
- [70] L. Li, K. Jamieson, A. Rostamizadeh, E. Gonina, J. Ben-tzur, M. Hardt, B. Recht, A. Talwalkar, A system for massively parallel hyperparameter tuning, in: I. Dhillon, D. Papailiopoulos, V. Sze (Eds.), Proceedings of Machine Learning and Systems, volume 2, 2020, pp. 230–246. URL: https://proceedings.mlsys.org/paper_files/paper/2020/file/a06f20b349c6cf09a6b171c71b88bbfc-Paper.pdf.
- [71] Y. Gorishniy, I. Rubachev, A. Babenko, On embeddings for numerical features in tabular deep learning, Advances in Neural Information Processing Systems 35 (2022) 24991–25004.
- [72] P. K. Dunn, G. K. Smyth, Randomized quantile residuals, Journal of Computational and graphical statistics 5 (1996) 236–244.
- [73] P. Mettes, E. van der Pol, C. G. M. Snoek, Hyperspherical Prototype Networks, Curran Associates Inc., Red Hook, NY, USA, 2019.
- [74] J. Demšar, Statistical comparisons of classifiers over multiple data sets, The Journal of Machine learning research 7 (2006) 1–30.

Appendix A. Benchmark challengers Hyperparameters

Hyperparameters search space of TVAE is in Table A.6, for CTGAN in A.7, for Tabsyn’s VAE and MLP in A.8, for TabDDPM in A.9, for SMOTE and ucSMOTE in A.10. These tables also present the reduced search spaces suggested and applied for the experiment of Section 5.

¹⁰ClusterBasedNormalizer

Parameter	Possible values	
	Extensive	Reduced
Learning rate	qLogUniform(1e-4, 1e-2, 1e-4)	qLogUniform(1e-4, 7.3e-03, 1e-4)
Batch size	[100, 500, 2000]	[100]
Embedding dim.	[16, 32, 64, 128, 256, 512]	[16, 32, 64]
Encoder dim.	[64, 128, 256, 512]	[256, 512]
Encoder depth	[2, 3, 4]	[2]
Decoder dim.	[64, 128, 256, 512]	[256, 512]
Decoder depth	[2, 3, 4]	[2, 4]
Loss factor	[3, 2, 1, 0.5]	[3, 2]
L2 scale	qLogUniform(1e-5, 1e-4, 1e-5)	qLogUniform(1e-5, 6.3e-5, 1e-6)
Numerical encoder	[CDF, PLE.CDF, PTP, QuantileTransformer, MinMaxScaler, CBN ¹⁰]	[CDF]
Categorical encoder	[one-hot-encoder]	[one-hot-encoder]
Epochs	[400]	∞
Number of trials per fold	300	50

Table A.6: Hyperparameter search spaces of TVAE: extensive and reduced. "Extensive" refers to the search space used during the extensive tuning done in this paper and "Reduced" refers to the reduced experiment performed in Section 5. for the extensive and reduced tuning experiments. For the reduced experiment the number of epochs was bounded by the time budget. A description of numerical encoders can be found in Section 3.2.1 and [11].

Parameter	Possible values	
	Extensive	Reduced
Discriminator learning rate	qLogUniform(5e-5, 1e-2, 5e-5)	qLogUniform(4e-4, 2.1e-03, 5e-5)
Generator learning rate	qLogUniform(5e-5, 1e-2, 5e-5)	qLogUniform(5e-5, 1.3e-3, 5e-5)
Batch size	[50, 100, 250, 500, 1000]	[100, 500, 1000]
Embedding dim.	[32, 64, 128, 256]	[32, 128]
Generator dim.	[128, 256]	[128]
Generator depth	[2, 3, 4]	[3, 4]
Discriminator dim.	[128, 256]	[256]
Discriminator depth	[2, 3]	[2, 3]
Generator decay	qLogUniform(1e-6, 1e-5, 1e-6)	qLogUniform(1e-6, 6.4e-6, 1e-7)
Discriminator decay	qLogUniform(1e-6, 1e-5, 1e-6)	qLogUniform(1e-6, 8e-6, 1e-6)
Log frequency	[False, True]	[False, True]
Numerical encoder	[CDF, PLE.CDF, PTP, QuantileTransformer, MinMaxScaler, CBN ¹⁰]	[CDF, PLE.CDF]
Categorical encoder	[one-hot-encoder]	[one-hot-encoder]
Epochs	[400]	∞
Number of trials per fold	300	50

Table A.7: Hyperparameter search spaces of CTGAN for the extensive and reduced tuning experiments. For the reduced experiment the number of epochs was bounded by the time budget.

Model	Parameter	Possible values	
		Extensive	Reduced
VAE	Learning rate	qLogUniform(5e-5, 1e-2, 5e-5)	qLogUniform(1.1e-3, 7.2e-3, 5e-5)
	Batch size	[1024, 2048, 4096]	[1024, 2048]
	Weight decay	qLogUniform(1e-6, 1e-5, 1e-6)	qLogUniform(1e-6, 5e-6, 1e-6)
	Token dim.	[2, 4]	[4]
	Number of head	[1, 2]	[1, 2]
	Factor	[8, 16, 32, 64]	[8, 16, 64]
	Number of layers	[1, 2, 3, 4]	[1, 2, 4]
	Max. beta	[1e-2]	[1e-2]
	Min. beta	[1e-5]	[1e-5]
	Lambda	[0.7, 0.8, 0.85, 0.9, 0.95]	[0.8, 0.85, 0.9, 0.95]
	Epochs	[4000]	∞
MLP	Learning Rate	qLogUniform(5e-5, 1e-2, 5e-5)	qLogUniform(7.7e-4, 2.5e-3, 1e-5)
	Weight Decay	qLogUniform(1e-6, 1e-5, 1e-6)	qLogUniform(1e-6, 3.3e-6, 1e-7)
	Batch size	[1024, 2048, 4096]	[1024, 4096]
	MLP’s hidden dimension	[512, 1024]	[1024]
	Epochs	[2000]	∞
Number of trials per fold		100	50

Table A.8: Hyperparameter search spaces of Tabsyn’s VAE and MLP for the extensive and reduced tuning experiments. For the reduced experiment the number of epochs was bounded by a time budget (10 minutes for the VAE, 10 minutes for the denoiser).

Parameter	Possible values	
	Extensive	Reduced
Batch size	[256, 4096]	[4096]
Dropout	[0.0]	[0.0]
Number of timesteps	[1000]	[1000]
Learning rate	qLogUniform(1e-5, 1e-3, 1e-5)	qLogUniform(3.5e-4, 9.2e-4, 1e-5)
Number of layers	[2, 4, 6, 8]	[2, 4, 6]
First layer’s dim.	[128, 256, 512, 1024]	[256, 512, 1024]
Middle layer’s dim.	[128, 256, 512, 1024]	[512, 1024]
Last layer’s dim.	[128, 256, 512, 1024]	[256, 512, 1024]
Training iterations	[20000]	∞
Number of trials per fold	300	50

Table A.9: Hyperparameter search spaces of TabDDPM for the extensive and reduced tuning experiments.

Parameter	Possible values
K-Neighbors	Grid search in range [2, 20]
Number of trials per fold	38

Table A.10: Hyperparameter search space of SMOTE and ucSMOTE.

Appendix B. Datasets Links

We provide the list of links toward the datasets that were used in this paper in Table B.11 below.

Dataset	URL
Abalone	www.openml.org/d/183
Adult	www.openml.org/d/1590
Bank	www.openml.org/d/1461
Black Friday	www.openml.org/d/41540
Bike Sharing	www.openml.org/d/42712
Covertypes	www.openml.org/d/150
Cardio	www.kaggle.com/sulianova/datasets
Churn	www.kaggle.com/datasets/shrutimechlearn
Diamonds	www.openml.org/d/42225
HELOC	www.kaggle.com/averkiyoliabev/datasets
Higgs	www.openml.org/d/4532
House 16H	www.openml.org/d/574
Insurance	www.kaggle.com/datasets/mirichoi0218/insurance
King	www.kaggle.com/datasets/harlfoxem/housesalesprediction
MiniBooNE	www.openml.org/d/41150
Moons	scikit-learn.org/stable/modules/classes.html

Table B.11: Links to the datasets

Appendix C. Dataset-Level Results for Large-Scale Experiment

Table C.12 present the per-dataset performance according to the metrics described in Section 3.1 averaged on 3-folds with 5 samples per-fold.

Dataset	Model	Metrics						
		C2ST ↓	DCR-Rate ↓	ML-Efficacy ↑	Shape ↑	Pair ↑	Train time ↓	Sample time ↓
Abalone	Train Copy	0.51 ± 0.00	1.00 ± 0.00	0.23 ± 0.01	0.96 ± 0.01	0.88 ± 0.01	-	-
	CTGAN	0.73 ± 0.03	0.63 ± 0.02	0.17 ± 0.01	0.93 ± 0.01	0.86 ± 0.03	685 ± 398.05	00 ± 0.01
	TVAE	0.75 ± 0.07	0.64 ± 0.05	0.23 ± 0.02	0.91 ± 0.02	0.86 ± 0.04	113 ± 3.57	00 ± 0.00
	TabDDPM	0.78 ± 0.01	0.67 ± 0.01	0.23 ± 0.01	0.94 ± 0.00	0.85 ± 0.03	169 ± 28.82	03 ± 0.62
	TabSyn	0.78 ± 0.01	0.63 ± 0.01	0.22 ± 0.01	0.94 ± 0.02	0.87 ± 0.00	1056 ± 183.20	00 ± 0.08
	SMOTE	1.00 ± 0.00	0.85 ± 0.02	0.50 ± 0.02	0.85 ± 0.01	0.72 ± 0.01	-	00 ± 0.01
	UC-SMOTE	0.89 ± 0.01	0.92 ± 0.02	0.51 ± 0.03	0.95 ± 0.01	0.88 ± 0.03	-	00 ± 0.00

Results per datasets and models under diverse metrics for the extensive search (continued).

Dataset	Model	Metrics						
		C2ST ↓	DCR-Rate ↓	ML-Efficacy ↑	Shape ↑	Pair ↑	Train time ↓	Sample time ↓
Adult	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.71 ± 0.01	0.99 ± 0.00	0.98 ± 0.00	-	-
	CTGAN	0.74 ± 0.02	0.72 ± 0.00	0.67 ± 0.01	0.97 ± 0.00	0.88 ± 0.01	2755 ± 920.97	00 ± 0.01
	TVAE	0.77 ± 0.01	0.72 ± 0.00	0.63 ± 0.03	0.96 ± 0.00	0.91 ± 0.02	1404 ± 2.70	00 ± 0.00
	TabDDPM	0.65 ± 0.00	0.62 ± 0.01	0.67 ± 0.01	0.98 ± 0.00	0.95 ± 0.01	443 ± 44.97	10 ± 2.22
	TabSyn	0.64 ± 0.01	0.62 ± 0.00	0.66 ± 0.01	0.98 ± 0.00	0.96 ± 0.01	5042 ± 2119.47	01 ± 1.06
	SMOTE	0.93 ± 0.00	0.85 ± 0.00	0.69 ± 0.01	0.95 ± 0.00	0.90 ± 0.01	-	06 ± 0.11
UC-SMOTE	0.93 ± 0.00	0.86 ± 0.01	0.67 ± 0.01	0.95 ± 0.00	0.90 ± 0.01	-	10 ± 0.03	
Bank marketing	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.54 ± 0.01	0.99 ± 0.00	0.98 ± 0.01	-	-
	CTGAN	0.72 ± 0.01	0.63 ± 0.00	0.47 ± 0.02	0.98 ± 0.00	0.95 ± 0.01	3159 ± 65.25	00 ± 0.01
	TVAE	0.81 ± 0.04	0.63 ± 0.00	0.45 ± 0.08	0.96 ± 0.01	0.92 ± 0.02	1089 ± 628.11	00 ± 0.02
	TabDDPM	0.65 ± 0.01	0.61 ± 0.00	0.52 ± 0.01	0.99 ± 0.01	0.96 ± 0.01	461 ± 42.63	10 ± 1.98
	TabSyn	0.61 ± 0.02	0.62 ± 0.01	0.49 ± 0.02	0.99 ± 0.01	0.97 ± 0.01	2783 ± 471.56	02 ± 0.03
	SMOTE	0.79 ± 0.01	0.98 ± 0.00	0.53 ± 0.02	0.97 ± 0.00	0.95 ± 0.00	-	09 ± 0.16
UC-SMOTE	0.79 ± 0.00	0.98 ± 0.00	0.46 ± 0.02	0.97 ± 0.00	0.95 ± 0.00	-	09 ± 0.03	
Bike sharing	Train Copy	0.50 ± 0.01	1.00 ± 0.00	0.94 ± 0.00	0.99 ± 0.00	0.53 ± 0.00	-	-
	CTGAN	0.90 ± 0.01	0.61 ± 0.00	0.68 ± 0.02	0.97 ± 0.01	0.53 ± 0.00	3772 ± 2301.33	00 ± 0.05
	TVAE	0.81 ± 0.01	0.61 ± 0.00	0.78 ± 0.03	0.96 ± 0.01	0.53 ± 0.00	540 ± 15.46	00 ± 0.00
	TabDDPM	0.81 ± 0.01	0.62 ± 0.00	0.79 ± 0.01	0.97 ± 0.01	0.53 ± 0.00	332 ± 13.26	07 ± 0.90
	TabSyn	0.86 ± 0.02	0.62 ± 0.01	0.51 ± 0.07	0.96 ± 0.00	0.53 ± 0.00	1843 ± 556.82	00 ± 0.01
	SMOTE	0.98 ± 0.00	0.99 ± 0.01	0.85 ± 0.00	0.95 ± 0.00	0.67 ± 0.08	-	01 ± 0.06
UC-SMOTE	0.98 ± 0.00	0.98 ± 0.00	0.86 ± 0.01	0.95 ± 0.01	0.67 ± 0.08	-	00 ± 0.09	
Black friday	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.53 ± 0.01	1.00 ± 0.00	0.99 ± 0.00	-	-
	CTGAN	0.82 ± 0.02	0.66 ± 0.01	0.46 ± 0.02	0.97 ± 0.01	0.87 ± 0.00	10111 ± 8081.87	00 ± 0.12
	TVAE	0.87 ± 0.01	0.68 ± 0.01	0.42 ± 0.02	0.95 ± 0.01	0.91 ± 0.02	4978 ± 1.32	00 ± 0.02
	TabDDPM	0.87 ± 0.01	0.64 ± 0.00	0.47 ± 0.02	0.99 ± 0.01	0.98 ± 0.01	366 ± 13.15	39 ± 14.32
	TabSyn	0.87 ± 0.01	0.68 ± 0.00	0.17 ± 0.01	0.99 ± 0.00	0.46 ± 0.00	7533 ± 563.20	06 ± 3.54
	SMOTE	0.80 ± 0.00	0.97 ± 0.00	0.50 ± 0.01	0.95 ± 0.00	0.94 ± 0.00	-	31 ± 0.16
UC-SMOTE	0.81 ± 0.01	0.97 ± 0.00	0.49 ± 0.01	0.94 ± 0.00	0.93 ± 0.00	-	106 ± 0.07	
Cardio	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.72 ± 0.01	1.00 ± 0.00	0.98 ± 0.01	-	-
	CTGAN	0.62 ± 0.01	0.64 ± 0.00	0.70 ± 0.02	0.99 ± 0.01	0.96 ± 0.01	4678 ± 59.79	00 ± 0.01
	TVAE	0.72 ± 0.02	0.64 ± 0.01	0.72 ± 0.01	0.97 ± 0.01	0.95 ± 0.02	1708 ± 984.53	00 ± 0.02
	TabDDPM	0.55 ± 0.01	0.62 ± 0.00	0.72 ± 0.01	0.99 ± 0.00	0.98 ± 0.01	185 ± 73.66	13 ± 8.85
	TabSyn	0.56 ± 0.00	0.64 ± 0.00	0.72 ± 0.00	0.99 ± 0.00	0.98 ± 0.01	2960 ± 578.82	03 ± 0.01
	SMOTE	0.95 ± 0.00	0.98 ± 0.00	0.73 ± 0.00	0.93 ± 0.00	0.97 ± 0.01	-	00 ± 0.02
UC-SMOTE	0.94 ± 0.00	0.98 ± 0.00	0.72 ± 0.01	0.93 ± 0.01	0.97 ± 0.01	-	01 ± 0.05	
Churn	Train Copy	0.50 ± 0.01	1.00 ± 0.00	0.59 ± 0.02	0.95 ± 0.00	0.87 ± 0.01	-	-
	CTGAN	0.63 ± 0.02	0.64 ± 0.01	0.36 ± 0.03	0.93 ± 0.01	0.84 ± 0.02	2505 ± 953.49	00 ± 0.03
	TVAE	0.64 ± 0.01	0.63 ± 0.00	0.51 ± 0.00	0.92 ± 0.01	0.84 ± 0.01	294 ± 6.20	00 ± 0.01
	TabDDPM	0.66 ± 0.07	0.64 ± 0.05	0.50 ± 0.05	0.89 ± 0.06	0.82 ± 0.06	618 ± 35.47	25 ± 2.36
	TabSyn	0.58 ± 0.02	0.61 ± 0.01	0.56 ± 0.03	0.93 ± 0.01	0.48 ± 0.00	1627 ± 340.95	00 ± 0.22
	SMOTE	0.76 ± 0.01	0.86 ± 0.03	0.50 ± 0.02	0.87 ± 0.01	0.81 ± 0.02	-	01 ± 0.06
UC-SMOTE	0.78 ± 0.02	0.93 ± 0.02	0.12 ± 0.09	0.87 ± 0.01	0.80 ± 0.02	-	00 ± 0.05	

Results per datasets and models under diverse metrics for the extensive search (continued).

Dataset	Model	Metrics						
		C2ST ↓	DCR-Rate ↓	ML-Efficacy ↑	Shape ↑	Pair ↑	Train time ↓	Sample time ↓
Coverttype	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.90 ± 0.00	1.00 ± 0.00	1.00 ± 0.01	-	-
	CTGAN	0.97 ± 0.01	0.60 ± 0.01	0.70 ± 0.01	0.98 ± 0.00	0.96 ± 0.01	51644 ± 17257.41	04 ± 0.60
	TVAE	0.90 ± 0.01	0.60 ± 0.00	0.77 ± 0.01	0.98 ± 0.00	0.96 ± 0.01	23219 ± 18238.02	01 ± 0.44
	TabDDPM	0.94 ± 0.01	0.59 ± 0.04	0.66 ± 0.04	0.95 ± 0.01	0.91 ± 0.01	868 ± 30.61	249 ± 25.67
	TabSyn	0.63 ± 0.02	0.62 ± 0.01	0.75 ± 0.02	0.99 ± 0.00	0.69 ± 0.00	4937 ± 1556.02	03 ± 0.05
	SMOTE	0.97 ± 0.00	0.97 ± 0.00	0.90 ± 0.00	1.00 ± 0.00	1.00 ± 0.01	-	125 ± 0.69
Diamonds	UC-SMOTE	0.97 ± 0.00	0.97 ± 0.00	0.90 ± 0.00	1.00 ± 0.00	0.99 ± 0.01	-	138 ± 0.71
	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.98 ± 0.00	0.99 ± 0.00	0.77 ± 0.01	-	-
	CTGAN	0.86 ± 0.01	0.65 ± 0.01	0.94 ± 0.01	0.97 ± 0.01	0.80 ± 0.04	3389 ± 2.57	00 ± 0.01
	TVAE	0.79 ± 0.03	0.64 ± 0.00	0.96 ± 0.01	0.94 ± 0.01	0.74 ± 0.04	1104 ± 637.85	00 ± 0.02
	TabDDPM	0.71 ± 0.01	0.61 ± 0.00	0.97 ± 0.00	0.98 ± 0.01	0.70 ± 0.02	374 ± 13.27	16 ± 6.28
	TabSyn	0.87 ± 0.01	0.65 ± 0.00	0.79 ± 0.10	0.98 ± 0.01	0.76 ± 0.04	3255 ± 360.89	02 ± 1.12
Heloc	SMOTE	0.97 ± 0.00	0.93 ± 0.01	0.92 ± 0.01	0.97 ± 0.00	0.77 ± 0.02	-	03 ± 0.02
	UC-SMOTE	0.97 ± 0.00	0.94 ± 0.01	0.93 ± 0.01	0.97 ± 0.00	0.80 ± 0.04	-	01 ± 0.04
	Train Copy	0.50 ± 0.01	1.00 ± 0.00	0.70 ± 0.01	0.99 ± 0.00	0.97 ± 0.01	-	-
	CTGAN	0.92 ± 0.01	0.64 ± 0.01	0.66 ± 0.05	0.97 ± 0.01	0.94 ± 0.02	4330 ± 3393.17	00 ± 0.10
	TVAE	0.87 ± 0.00	0.66 ± 0.01	0.70 ± 0.01	0.95 ± 0.00	0.94 ± 0.01	421 ± 249.85	00 ± 0.00
	TabDDPM	0.72 ± 0.01	0.64 ± 0.01	0.70 ± 0.01	0.97 ± 0.00	0.95 ± 0.01	80 ± 9.69	01 ± 0.53
Higgs	TabSyn	0.70 ± 0.02	0.66 ± 0.01	0.70 ± 0.02	0.97 ± 0.00	0.95 ± 0.01	1870 ± 302.36	00 ± 0.00
	SMOTE	0.92 ± 0.01	1.00 ± 0.00	0.69 ± 0.01	0.94 ± 0.00	0.95 ± 0.02	-	00 ± 0.02
	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.74 ± 0.00	0.99 ± 0.00	0.99 ± 0.00	-	-
	CTGAN	0.85 ± 0.02	0.60 ± 0.00	0.70 ± 0.02	0.99 ± 0.01	0.98 ± 0.00	12696 ± 4197.49	00 ± 0.13
	TVAE	0.92 ± 0.01	0.60 ± 0.00	0.70 ± 0.02	0.94 ± 0.01	0.98 ± 0.01	5752 ± 116.50	00 ± 0.08
	TabDDPM	0.57 ± 0.00	0.63 ± 0.00	0.73 ± 0.00	0.99 ± 0.00	0.99 ± 0.00	223 ± 10.78	22 ± 2.25
House 16h	TabSyn	0.57 ± 0.01	0.61 ± 0.00	0.73 ± 0.00	0.94 ± 0.00	0.99 ± 0.00	7503 ± 2581.02	04 ± 2.17
	SMOTE	0.84 ± 0.00	1.00 ± 0.00	0.73 ± 0.00	0.96 ± 0.00	0.98 ± 0.01	-	00 ± 0.05
	UC-SMOTE	0.83 ± 0.00	1.00 ± 0.00	0.73 ± 0.00	0.96 ± 0.00	0.98 ± 0.01	-	01 ± 0.04
	Train Copy	0.50 ± 0.01	1.00 ± 0.00	0.64 ± 0.01	0.99 ± 0.01	0.99 ± 0.00	-	-
	CTGAN	0.84 ± 0.01	0.62 ± 0.00	0.47 ± 0.01	0.98 ± 0.01	0.97 ± 0.01	983 ± 559.75	00 ± 0.01
	TVAE	0.82 ± 0.03	0.61 ± 0.01	0.48 ± 0.06	0.95 ± 0.01	0.98 ± 0.01	922 ± 6.68	00 ± 0.01
Insurance	TabDDPM	0.60 ± 0.01	0.61 ± 0.00	0.61 ± 0.02	0.98 ± 0.00	0.99 ± 0.01	95 ± 21.99	02 ± 0.88
	TabSyn	0.74 ± 0.01	0.62 ± 0.01	0.40 ± 0.02	0.97 ± 0.00	0.99 ± 0.00	2370 ± 810.03	01 ± 0.01
	SMOTE	0.87 ± 0.00	0.86 ± 0.05	0.62 ± 0.01	0.92 ± 0.01	0.99 ± 0.00	-	00 ± 0.03
	UC-SMOTE	0.87 ± 0.01	0.83 ± 0.05	0.62 ± 0.02	0.92 ± 0.01	0.99 ± 0.00	-	00 ± 0.02
	Train Copy	0.48 ± 0.01	1.00 ± 0.00	0.85 ± 0.03	0.96 ± 0.01	0.91 ± 0.01	-	-
	CTGAN	0.67 ± 0.03	0.92 ± 0.01	0.71 ± 0.01	0.94 ± 0.00	0.87 ± 0.02	265 ± 23.11	00 ± 0.00
Insurance	TVAE	0.67 ± 0.02	0.92 ± 0.01	0.77 ± 0.01	0.91 ± 0.01	0.86 ± 0.02	32 ± 4.50	00 ± 0.00
	TabDDPM	0.68 ± 0.02	0.62 ± 0.01	0.83 ± 0.03	0.93 ± 0.02	0.89 ± 0.02	181 ± 18.27	03 ± 1.01
	TabSyn	0.60 ± 0.02	0.61 ± 0.02	0.82 ± 0.00	0.94 ± 0.01	0.88 ± 0.00	1156 ± 74.54	00 ± 0.03
	SMOTE	0.68 ± 0.01	0.80 ± 0.03	0.81 ± 0.01	0.94 ± 0.01	0.87 ± 0.01	-	00 ± 0.01
	UC-SMOTE	0.67 ± 0.01	0.85 ± 0.02	0.80 ± 0.02	0.93 ± 0.02	0.87 ± 0.02	-	00 ± 0.00

Results per datasets and models under diverse metrics for the extensive search (continued).

Dataset	Model	Metrics						
		C2ST ↓	DCR-Rate ↓	ML-Efficacy ↑	Shape ↑	Pair ↑	Train time ↓	Sample time ↓
King	Train Copy	0.50 ± 0.01	1.00 ± 0.00	0.87 ± 0.00	0.98 ± 0.01	0.97 ± 0.01	-	-
	CTGAN	0.93 ± 0.01	0.62 ± 0.00	0.72 ± 0.04	0.97 ± 0.01	0.94 ± 0.01	1805 ± 8.66	00 ± 0.01
	TVAE	0.94 ± 0.01	0.61 ± 0.01	0.81 ± 0.02	0.94 ± 0.01	0.94 ± 0.01	990 ± 6.97	00 ± 0.00
	TabDDPM	0.97 ± 0.01	0.73 ± 0.06	0.64 ± 0.20	0.71 ± 0.06	0.83 ± 0.01	244 ± 14.03	06 ± 0.61
	TabSyn	0.93 ± 0.01	0.62 ± 0.00	0.12 ± 0.03	0.96 ± 0.01	0.94 ± 0.00	2717 ± 648.93	01 ± 0.01
	SMOTE	0.98 ± 0.00	0.97 ± 0.00	0.83 ± 0.01	0.92 ± 0.01	0.95 ± 0.01	-	02 ± 0.23
	UC-SMOTE	0.98 ± 0.00	0.96 ± 0.01	0.82 ± 0.03	0.92 ± 0.01	0.96 ± 0.01	-	00 ± 0.03
Miniboo ne	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.89 ± 0.01	0.99 ± 0.00	0.99 ± 0.01	-	-
	CTGAN	0.91 ± 0.02	0.60 ± 0.01	0.83 ± 0.03	0.96 ± 0.02	0.58 ± 0.01	6879 ± 2388.04	00 ± 0.06
	TVAE	0.95 ± 0.03	0.60 ± 0.00	0.86 ± 0.00	0.93 ± 0.01	0.59 ± 0.01	4876 ± 4730.59	00 ± 0.05
	TabDDPM	0.64 ± 0.02	0.61 ± 0.01	0.89 ± 0.00	0.99 ± 0.00	0.91 ± 0.03	193 ± 46.01	23 ± 8.73
	TabSyn	0.61 ± 0.01	0.61 ± 0.00	0.89 ± 0.00	0.99 ± 0.01	0.91 ± 0.03	7861 ± 1700.48	07 ± 0.08
	SMOTE	0.83 ± 0.00	1.00 ± 0.00	0.89 ± 0.00	0.97 ± 0.01	0.99 ± 0.01	-	01 ± 0.10
	UC-SMOTE	0.83 ± 0.00	1.00 ± 0.00	0.89 ± 0.00	0.97 ± 0.01	0.98 ± 0.01	-	03 ± 0.04
Moons	Train Copy	0.50 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	0.99 ± 0.00	0.99 ± 0.00	-	-
	CTGAN	0.68 ± 0.01	0.61 ± 0.01	1.00 ± 0.00	0.96 ± 0.01	0.97 ± 0.01	2571 ± 967.05	00 ± 0.02
	TVAE	0.63 ± 0.01	0.61 ± 0.02	1.00 ± 0.00	0.98 ± 0.00	0.97 ± 0.00	752 ± 51.10	00 ± 0.01
	TabDDPM	0.52 ± 0.01	0.61 ± 0.01	1.00 ± 0.00	0.99 ± 0.00	0.99 ± 0.00	183 ± 92.44	09 ± 6.74
	TabSyn	0.54 ± 0.01	0.61 ± 0.01	1.00 ± 0.00	0.99 ± 0.00	0.99 ± 0.01	2647 ± 656.79	01 ± 0.81
	SMOTE	0.54 ± 0.00	0.84 ± 0.05	1.00 ± 0.00	0.99 ± 0.00	0.99 ± 0.00	-	00 ± 0.02
	UC-SMOTE	0.54 ± 0.01	0.80 ± 0.03	1.00 ± 0.00	0.99 ± 0.00	0.99 ± 0.01	-	00 ± 0.03

Table C.12: Results per dataset and model under diverse metrics. The training times (gradient descent) and sampling times are given in seconds. The sampling times are given for 5 samples. The best values per metric are formatted in **bold green** and the worse values are in **red**.

Appendix D. Training and Sampling Impact

We performed large-scale experiments including a quite costly hyperparameter tuning step. Note that the term ”costs” in this section refers to the duration, energy consumption, and CO₂ emissions. We evaluate the training and sampling costs of the models with their optimized hyperparameters, as well as the full hyperparameter tuning costs.

Appendix D.1. Training and Sampling Costs

Due to the massive nature of the experiments, the hyperparameter search could not all be run on the same hardware. We hence estimated the training (gradient descent) cost by running all models on a single Tesla V100-SXM2 32 GB. Table D.13 provides the raw training and sampling energy consumption

and emissions for reference. All of the Tabsyn costs shown are a combination of the VAE cost and the denoiser cost, which are estimated separately and then added together. Also, all Tabsyn costs on the *Covertime* dataset are estimated considering the complete dataset size.

Dataset	Model	Emissions Training ↓	Energy Training ↓	Emissions Sampling ↓	Energy Sampling ↓
Abalone	CTGAN	$3.05 \pm 1.77 \cdot 10^{-3}$	$4.53 \pm 2.62 \cdot 10^{-2}$	$9.74 \pm 2.37 \cdot 10^{-8}$	$1.45 \pm 0.35 \cdot 10^{-6}$
	TabDDPM	$9.16 \pm 1.42 \cdot 10^{-6}$	$1.36 \pm 0.21 \cdot 10^{-4}$	$2.90 \pm 0.50 \cdot 10^{-5}$	$4.31 \pm 0.74 \cdot 10^{-4}$
	TabSyn	$3.18 \pm 0.32 \cdot 10^{-4}$	$4.27 \pm 0.32 \cdot 10^{-3}$	$1.49 \pm 0.78 \cdot 10^{-6}$	$2.21 \pm 1.16 \cdot 10^{-5}$
	TVAE	$5.06 \pm 0.13 \cdot 10^{-4}$	$7.51 \pm 0.19 \cdot 10^{-3}$	$6.23 \pm 2.14 \cdot 10^{-8}$	$9.24 \pm 3.17 \cdot 10^{-7}$
	SMOTE	-	-	$4.04 \pm 0.35 \cdot 10^{-7}$	$5.99 \pm 0.51 \cdot 10^{-6}$
	UC-SMOTE	-	-	$1.37 \pm 0.07 \cdot 10^{-7}$	$2.03 \pm 0.11 \cdot 10^{-6}$
Adult	CTGAN	$1.24 \pm 0.41 \cdot 10^{-2}$	$1.84 \pm 0.61 \cdot 10^{-1}$	$9.67 \pm 0.41 \cdot 10^{-7}$	$1.43 \pm 0.06 \cdot 10^{-5}$
	TabDDPM	$2.68 \pm 0.35 \cdot 10^{-4}$	$3.98 \pm 0.51 \cdot 10^{-3}$	$8.41 \pm 2.01 \cdot 10^{-5}$	$1.25 \pm 0.30 \cdot 10^{-3}$
	TabSyn	$1.68 \pm 0.72 \cdot 10^{-3}$	$1.42 \pm 0.36 \cdot 10^{-2}$	$1.56 \pm 0.85 \cdot 10^{-5}$	$2.31 \pm 1.26 \cdot 10^{-4}$
	TVAE	$6.30 \pm 0.01 \cdot 10^{-3}$	$9.34 \pm 0.01 \cdot 10^{-2}$	$3.54 \pm 0.07 \cdot 10^{-7}$	$5.25 \pm 0.11 \cdot 10^{-6}$
	SMOTE	-	-	$2.46 \pm 0.04 \cdot 10^{-5}$	$3.66 \pm 0.06 \cdot 10^{-4}$
	UC-SMOTE	-	-	$3.82 \pm 0.01 \cdot 10^{-5}$	$5.67 \pm 0.01 \cdot 10^{-4}$
Bank marketing	CTGAN	$1.41 \pm 0.03 \cdot 10^{-2}$	$2.09 \pm 0.04 \cdot 10^{-1}$	$8.88 \pm 0.32 \cdot 10^{-7}$	$1.32 \pm 0.05 \cdot 10^{-5}$
	TabDDPM	$2.50 \pm 0.22 \cdot 10^{-4}$	$3.71 \pm 0.33 \cdot 10^{-3}$	$7.81 \pm 1.91 \cdot 10^{-5}$	$1.16 \pm 0.28 \cdot 10^{-3}$
	TabSyn	$1.06 \pm 0.08 \cdot 10^{-3}$	$1.46 \pm 0.04 \cdot 10^{-2}$	$1.80 \pm 0.02 \cdot 10^{-5}$	$2.67 \pm 0.30 \cdot 10^{-4}$
	TVAE	$4.90 \pm 2.83 \cdot 10^{-3}$	$7.27 \pm 4.19 \cdot 10^{-2}$	$2.84 \pm 0.47 \cdot 10^{-7}$	$4.22 \pm 0.69 \cdot 10^{-6}$
	SMOTE	-	-	$3.46 \pm 0.05 \cdot 10^{-5}$	$5.13 \pm 0.08 \cdot 10^{-4}$
	UC-SMOTE	-	-	$3.32 \pm 0.01 \cdot 10^{-5}$	$4.92 \pm 0.02 \cdot 10^{-4}$
Bike sharing	CTGAN	$1.68 \pm 1.03 \cdot 10^{-2}$	$2.50 \pm 1.52 \cdot 10^{-1}$	$5.79 \pm 2.48 \cdot 10^{-7}$	$8.59 \pm 3.68 \cdot 10^{-6}$
	TabDDPM	$8.73 \pm 0.70 \cdot 10^{-5}$	$1.30 \pm 0.10 \cdot 10^{-3}$	$6.04 \pm 0.72 \cdot 10^{-5}$	$8.96 \pm 1.07 \cdot 10^{-4}$
	TabSyn	$6.29 \pm 1.72 \cdot 10^{-4}$	$8.11 \pm 0.27 \cdot 10^{-3}$	$7.02 \pm 0.03 \cdot 10^{-6}$	$1.04 \pm 0.00 \cdot 10^{-4}$
	TVAE	$2.42 \pm 0.07 \cdot 10^{-3}$	$3.59 \pm 0.11 \cdot 10^{-2}$	$1.37 \pm 0.27 \cdot 10^{-7}$	$2.03 \pm 0.40 \cdot 10^{-6}$
	SMOTE	-	-	$4.53 \pm 0.22 \cdot 10^{-6}$	$6.72 \pm 0.33 \cdot 10^{-5}$
	UC-SMOTE	-	-	$1.48 \pm 0.32 \cdot 10^{-6}$	$2.20 \pm 0.48 \cdot 10^{-5}$
Black friday	CTGAN	$4.55 \pm 3.60 \cdot 10^{-2}$	$6.75 \pm 5.34 \cdot 10^{-1}$	$2.70 \pm 0.56 \cdot 10^{-6}$	$4.01 \pm 0.83 \cdot 10^{-5}$
	TabDDPM	$8.97 \pm 1.75 \cdot 10^{-4}$	$1.33 \pm 0.26 \cdot 10^{-2}$	$2.98 \pm 1.09 \cdot 10^{-4}$	$4.42 \pm 1.62 \cdot 10^{-3}$
	TabSyn	$2.78 \pm 0.26 \cdot 10^{-3}$	$2.58 \pm 0.90 \cdot 10^{-2}$	$5.18 \pm 2.81 \cdot 10^{-5}$	$7.69 \pm 4.17 \cdot 10^{-4}$
	TVAE	$2.23 \pm 0.01 \cdot 10^{-2}$	$3.31 \pm 0.01 \cdot 10^{-1}$	$8.59 \pm 0.65 \cdot 10^{-7}$	$1.27 \pm 0.10 \cdot 10^{-5}$
	SMOTE	-	-	$1.14 \pm 0.01 \cdot 10^{-4}$	$1.68 \pm 0.01 \cdot 10^{-3}$
	UC-SMOTE	-	-	$3.84 \pm 0.02 \cdot 10^{-4}$	$5.70 \pm 0.03 \cdot 10^{-3}$
Cardio	CTGAN	$2.11 \pm 0.03 \cdot 10^{-2}$	$3.13 \pm 0.04 \cdot 10^{-1}$	$1.30 \pm 0.03 \cdot 10^{-6}$	$1.93 \pm 0.04 \cdot 10^{-5}$
	TabDDPM	$2.22 \pm 1.13 \cdot 10^{-4}$	$3.30 \pm 1.68 \cdot 10^{-3}$	$1.05 \pm 0.70 \cdot 10^{-4}$	$1.56 \pm 1.04 \cdot 10^{-3}$
	TabSyn	$1.15 \pm 0.16 \cdot 10^{-3}$	$1.67 \pm 0.03 \cdot 10^{-2}$	$2.89 \pm 0.01 \cdot 10^{-5}$	$4.29 \pm 0.02 \cdot 10^{-4}$
	TVAE	$7.64 \pm 4.39 \cdot 10^{-3}$	$1.13 \pm 0.65 \cdot 10^{-1}$	$2.70 \pm 0.94 \cdot 10^{-7}$	$4.00 \pm 1.40 \cdot 10^{-6}$
	SMOTE	-	-	$2.27 \pm 0.05 \cdot 10^{-6}$	$3.37 \pm 0.07 \cdot 10^{-5}$
	UC-SMOTE	-	-	$4.27 \pm 0.16 \cdot 10^{-6}$	$6.34 \pm 0.24 \cdot 10^{-5}$
Churn	CTGAN	$1.15 \pm 0.43 \cdot 10^{-2}$	$1.71 \pm 0.63 \cdot 10^{-1}$	$6.61 \pm 1.15 \cdot 10^{-7}$	$9.81 \pm 1.71 \cdot 10^{-6}$
	TabDDPM	$1.10 \pm 0.07 \cdot 10^{-4}$	$1.63 \pm 0.10 \cdot 10^{-3}$	$1.96 \pm 0.17 \cdot 10^{-4}$	$2.90 \pm 0.25 \cdot 10^{-3}$
	TabSyn	$4.97 \pm 0.79 \cdot 10^{-4}$	$6.25 \pm 0.62 \cdot 10^{-3}$	$3.43 \pm 1.81 \cdot 10^{-6}$	$5.09 \pm 2.68 \cdot 10^{-5}$
	TVAE	$1.34 \pm 0.03 \cdot 10^{-3}$	$1.98 \pm 0.05 \cdot 10^{-2}$	$3.08 \pm 0.29 \cdot 10^{-7}$	$4.57 \pm 0.42 \cdot 10^{-6}$
	SMOTE	-	-	$4.13 \pm 0.20 \cdot 10^{-6}$	$6.12 \pm 0.29 \cdot 10^{-5}$
	UC-SMOTE	-	-	$3.53 \pm 0.19 \cdot 10^{-6}$	$5.24 \pm 0.28 \cdot 10^{-5}$

Tuning costs estimation for the best models (continued).

Dataset	Model	Emissions Training ↓	Energy Training ↓	Emissions Sampling ↓	Energy Sampling ↓
Coverttype	CTGAN	$2.33 \pm 0.78 \cdot 10^{-1}$	3.46 ± 1.15	$1.92 \pm 0.27 \cdot 10^{-5}$	$2.85 \pm 0.41 \cdot 10^{-4}$
	TabDDPM	$6.22 \pm 0.35 \cdot 10^{-3}$	$9.23 \pm 0.53 \cdot 10^{-2}$	$1.80 \pm 0.26 \cdot 10^{-3}$	$2.67 \pm 0.38 \cdot 10^{-2}$
	TabSyn	$1.82 \pm 0.68 \cdot 10^{-3}$	$1.26 \pm 0.72 \cdot 10^{-2}$	$2.37 \pm 0.04 \cdot 10^{-5}$	$3.52 \pm 0.06 \cdot 10^{-4}$
	TVAE	$1.04 \pm 0.81 \cdot 10^{-1}$	1.54 ± 1.21	$8.00 \pm 1.93 \cdot 10^{-6}$	$1.19 \pm 0.29 \cdot 10^{-4}$
	SMOTE	-	-	$4.52 \pm 0.02 \cdot 10^{-4}$	$6.71 \pm 0.04 \cdot 10^{-3}$
	UC-SMOTE	-	-	$5.01 \pm 0.03 \cdot 10^{-4}$	$7.43 \pm 0.04 \cdot 10^{-3}$
Diamonds	CTGAN	$1.52 \pm 0.00 \cdot 10^{-2}$	$2.26 \pm 0.01 \cdot 10^{-1}$	$8.44 \pm 0.13 \cdot 10^{-7}$	$1.25 \pm 0.02 \cdot 10^{-5}$
	TabDDPM	$3.33 \pm 0.71 \cdot 10^{-4}$	$4.93 \pm 1.06 \cdot 10^{-3}$	$1.30 \pm 0.47 \cdot 10^{-4}$	$1.92 \pm 0.70 \cdot 10^{-3}$
	TabSyn	$1.21 \pm 0.18 \cdot 10^{-3}$	$1.36 \pm 0.49 \cdot 10^{-2}$	$1.67 \pm 0.89 \cdot 10^{-5}$	$2.48 \pm 1.32 \cdot 10^{-4}$
	TVAE	$4.93 \pm 2.84 \cdot 10^{-3}$	$7.31 \pm 4.21 \cdot 10^{-2}$	$2.42 \pm 0.76 \cdot 10^{-7}$	$3.59 \pm 1.12 \cdot 10^{-6}$
	SMOTE	-	-	$1.32 \pm 0.01 \cdot 10^{-5}$	$1.96 \pm 0.01 \cdot 10^{-4}$
	UC-SMOTE	-	-	$3.67 \pm 0.14 \cdot 10^{-6}$	$5.44 \pm 0.21 \cdot 10^{-5}$
Heloc	CTGAN	$1.93 \pm 1.51 \cdot 10^{-2}$	$2.87 \pm 2.24 \cdot 10^{-1}$	$6.03 \pm 4.19 \cdot 10^{-7}$	$8.95 \pm 6.22 \cdot 10^{-6}$
	TabDDPM	$1.21 \pm 0.29 \cdot 10^{-5}$	$1.80 \pm 0.43 \cdot 10^{-4}$	$1.50 \pm 0.42 \cdot 10^{-5}$	$2.23 \pm 0.63 \cdot 10^{-4}$
	TabSyn	$5.78 \pm 0.66 \cdot 10^{-4}$	$6.91 \pm 0.18 \cdot 10^{-3}$	$4.60 \pm 0.11 \cdot 10^{-6}$	$6.83 \pm 0.16 \cdot 10^{-5}$
	TVAE	$1.88 \pm 1.11 \cdot 10^{-3}$	$2.80 \pm 1.65 \cdot 10^{-2}$	$1.15 \pm 0.43 \cdot 10^{-7}$	$1.71 \pm 0.63 \cdot 10^{-6}$
	SMOTE	-	-	$9.32 \pm 1.28 \cdot 10^{-7}$	$1.38 \pm 0.19 \cdot 10^{-5}$
	UC-SMOTE	-	-	$2.02 \pm 0.07 \cdot 10^{-6}$	$3.00 \pm 0.10 \cdot 10^{-5}$
Higgs	CTGAN	$5.71 \pm 1.88 \cdot 10^{-2}$	$8.47 \pm 2.79 \cdot 10^{-1}$	$2.92 \pm 0.56 \cdot 10^{-6}$	$4.33 \pm 0.84 \cdot 10^{-5}$
	TabDDPM	$4.01 \pm 0.25 \cdot 10^{-4}$	$5.96 \pm 0.37 \cdot 10^{-3}$	$1.77 \pm 0.18 \cdot 10^{-4}$	$2.63 \pm 0.27 \cdot 10^{-3}$
	TabSyn	$2.88 \pm 1.12 \cdot 10^{-3}$	$2.07 \pm 0.74 \cdot 10^{-2}$	$3.20 \pm 1.72 \cdot 10^{-5}$	$4.75 \pm 2.55 \cdot 10^{-4}$
	TVAE	$2.57 \pm 0.05 \cdot 10^{-2}$	$3.81 \pm 0.08 \cdot 10^{-1}$	$1.06 \pm 0.34 \cdot 10^{-6}$	$1.58 \pm 0.51 \cdot 10^{-5}$
	SMOTE	-	-	$2.01 \pm 0.19 \cdot 10^{-6}$	$2.98 \pm 0.28 \cdot 10^{-5}$
	UC-SMOTE	-	-	$6.19 \pm 0.13 \cdot 10^{-6}$	$9.18 \pm 0.20 \cdot 10^{-5}$
House 16h	CTGAN	$4.42 \pm 2.48 \cdot 10^{-3}$	$6.55 \pm 3.68 \cdot 10^{-2}$	$1.24 \pm 0.16 \cdot 10^{-7}$	$1.85 \pm 0.24 \cdot 10^{-6}$
	TabDDPM	$3.81 \pm 1.14 \cdot 10^{-5}$	$5.66 \pm 1.69 \cdot 10^{-4}$	$2.03 \pm 0.70 \cdot 10^{-5}$	$3.01 \pm 1.04 \cdot 10^{-4}$
	TabSyn	$8.56 \pm 2.73 \cdot 10^{-4}$	$9.40 \pm 1.12 \cdot 10^{-3}$	$9.68 \pm 0.08 \cdot 10^{-6}$	$1.44 \pm 0.01 \cdot 10^{-4}$
	TVAE	$4.12 \pm 0.02 \cdot 10^{-3}$	$6.11 \pm 0.03 \cdot 10^{-2}$	$1.40 \pm 0.17 \cdot 10^{-7}$	$2.08 \pm 0.26 \cdot 10^{-6}$
	SMOTE	-	-	$2.92 \pm 1.19 \cdot 10^{-7}$	$4.34 \pm 1.77 \cdot 10^{-6}$
	UC-SMOTE	-	-	$4.04 \pm 0.48 \cdot 10^{-7}$	$5.99 \pm 0.72 \cdot 10^{-6}$
Insurance	CTGAN	$1.18 \pm 0.10 \cdot 10^{-3}$	$1.75 \pm 0.15 \cdot 10^{-2}$	$7.40 \pm 1.62 \cdot 10^{-8}$	$1.10 \pm 0.24 \cdot 10^{-6}$
	TabDDPM	$2.83 \pm 0.29 \cdot 10^{-6}$	$4.20 \pm 0.43 \cdot 10^{-5}$	$2.65 \pm 0.94 \cdot 10^{-5}$	$3.93 \pm 1.40 \cdot 10^{-4}$
	TabSyn	$3.45 \pm 0.04 \cdot 10^{-4}$	$4.64 \pm 0.36 \cdot 10^{-3}$	$5.87 \pm 2.03 \cdot 10^{-7}$	$8.71 \pm 3.01 \cdot 10^{-6}$
	TVAE	$1.44 \pm 0.20 \cdot 10^{-4}$	$2.14 \pm 0.29 \cdot 10^{-3}$	$4.70 \pm 1.02 \cdot 10^{-8}$	$6.98 \pm 1.51 \cdot 10^{-7}$
	SMOTE	-	-	$1.83 \pm 0.26 \cdot 10^{-7}$	$2.71 \pm 0.39 \cdot 10^{-6}$
	UC-SMOTE	-	-	$1.02 \pm 0.04 \cdot 10^{-7}$	$1.51 \pm 0.06 \cdot 10^{-6}$
King	CTGAN	$8.15 \pm 0.04 \cdot 10^{-3}$	$1.21 \pm 0.01 \cdot 10^{-1}$	$5.72 \pm 0.34 \cdot 10^{-7}$	$8.49 \pm 0.50 \cdot 10^{-6}$
	TabDDPM	$9.27 \pm 0.56 \cdot 10^{-5}$	$1.38 \pm 0.08 \cdot 10^{-3}$	$5.08 \pm 0.49 \cdot 10^{-5}$	$7.55 \pm 0.73 \cdot 10^{-4}$
	TabSyn	$9.40 \pm 1.26 \cdot 10^{-4}$	$9.57 \pm 0.35 \cdot 10^{-3}$	$9.15 \pm 0.02 \cdot 10^{-6}$	$1.36 \pm 0.00 \cdot 10^{-4}$
	TVAE	$4.43 \pm 0.04 \cdot 10^{-3}$	$6.57 \pm 0.07 \cdot 10^{-2}$	$2.59 \pm 0.18 \cdot 10^{-7}$	$3.85 \pm 0.26 \cdot 10^{-6}$
	SMOTE	-	-	$8.77 \pm 0.84 \cdot 10^{-6}$	$1.30 \pm 0.12 \cdot 10^{-4}$
	UC-SMOTE	-	-	$2.25 \pm 0.10 \cdot 10^{-6}$	$3.34 \pm 0.14 \cdot 10^{-5}$

Tuning costs estimation for the best models (continued).

Dataset	Model	Emissions Training ↓	Energy Training ↓	Emissions Sampling ↓	Energy Sampling ↓
Miniboo ne	CTGAN	3.16±1.09 · 10⁻²	4.69±1.62 · 10⁻¹	3.50±0.24 · 10 ⁻⁶	5.20±0.35 · 10 ⁻⁵
	TabDDPM	4.50±1.33 · 10⁻⁴	6.68±1.97 · 10⁻³	1.84±0.69 · 10⁻⁴	2.74±1.03 · 10⁻³
	TabSyn	3.41±0.61 · 10 ⁻³	3.18±0.79 · 10 ⁻²	5.57±0.07 · 10 ⁻⁵	8.26±0.11 · 10 ⁻⁴
	TVAE	2.18±2.12 · 10 ⁻²	3.24±3.14 · 10 ⁻¹	1.94±0.22 · 10⁻⁶	2.88±0.32 · 10⁻⁵
	SMOTE	-	-	4.39±0.34 · 10 ⁻⁶	6.51±0.51 · 10 ⁻⁵
	UC-SMOTE	-	-	1.32±0.01 · 10 ⁻⁵	1.95±0.02 · 10 ⁻⁴
Moons	CTGAN	1.15±0.43 · 10⁻²	1.71±0.64 · 10⁻¹	5.02±1.07 · 10 ⁻⁷	7.45±1.59 · 10 ⁻⁶
	TabDDPM	1.31±0.77 · 10⁻⁴	1.94±1.14 · 10⁻³	7.39±5.33 · 10⁻⁵	1.10±0.79 · 10⁻³
	TabSyn	9.26±2.41 · 10 ⁻⁴	1.14±0.36 · 10 ⁻²	1.22±0.65 · 10 ⁻⁵	1.81±0.96 · 10 ⁻⁴
	TVAE	3.37±0.24 · 10 ⁻³	5.00±0.36 · 10 ⁻²	1.78±0.43 · 10⁻⁷	2.65±0.64 · 10⁻⁶
	SMOTE	-	-	2.44±0.52 · 10 ⁻⁷	3.63±0.78 · 10 ⁻⁶
	UC-SMOTE	-	-	1.03±0.12 · 10 ⁻⁶	1.53±0.18 · 10 ⁻⁵

Table D.13: CO₂ emissions (in Kg) and Energy Consumption (in kWh) of the benchmark challengers. The energy consumption is obtained by summing the CPU, GPU and RAM energy. Training costs are given for all models on the same basis of 400 epochs. Sampling costs are given for 5 samples. The best values are formatted in **bold green** and the worse are in **red**.

Appendix D.2. Whole Tuning Cost Estimation

As mentioned in Section 3.1 we could not perform the hyperparameter search and training phases on a uniform hardware and software architecture and we estimated the tuning cost with Equation (1).

Each trial can be stopped based on three conditions: an early stopping decided by the model, a poor C2ST performance or a time limit. Therefore, to get an accurate estimate of the init-cost and the cost-per-step from the single-GPU mentioned in Appendix D.1, we needed to extract from our logs the effective number of training steps performed per trial.

In addition, we also considered the parallelization scheme applied during the tuning procedure. One issue with TabSyn was that we observed a general slowdown when the model was parallelized too heavily. This issue was even more marked on datasets with a large number of columns. We hence reduced the number trials per GPU and the global number of trials to 100 to fall back to a reasonable time for this model.

Finally, we measured cost on a typical configuration used during our tuning experiment: 8 Tesla V100-SXM2 32 GB. Considering those 8 GPUS, we used the following parallel allocation of trials: 64 for TVAE and CTGAN, 40 for TabDDPM, and 16 for TabSyn. A model that can be easily parallelized during the hyperparameter tuning phase offers a cost advantage. It is hence

important to consider this aspect during the evaluation process. The results are presented in Table D.14.

Dataset	Model	Metrics		
		Energy (kWh) ↓	Emissions (Kg) ↓	Duration (HH:MM) ↓
Abalone	TVAE	0.00	0.01	00:02
	CTGAN	0.01	0.08	00:19
	TabDDPM	0.00	0.05	00:10
	TabSyn	0.03	0.40	01:31
	SMOTE	0.00	0.00	00:00
	UC-SMOTE	0.00	0.00	00:00
Adult	TVAE	0.01	0.17	00:42
	CTGAN	0.02	0.30	01:15
	TabDDPM	0.01	0.12	00:25
	TabSyn	0.08	0.91	03:27
	SMOTE	0.00	0.01	00:04
	UC-SMOTE	0.00	0.02	00:06
Bank marketing	TVAE	0.01	0.10	00:25
	CTGAN	0.03	0.38	01:35
	TabDDPM	0.01	0.13	00:29
	TabSyn	0.08	0.96	03:27
	SMOTE	0.00	0.02	00:06
	UC-SMOTE	0.00	0.02	00:05
Bike sharing	TVAE	0.00	0.05	00:12
	CTGAN	0.03	0.40	01:41
	TabDDPM	0.01	0.12	00:22
	TabSyn	0.05	0.61	02:18
	SMOTE	0.00	0.00	00:00
	UC-SMOTE	0.00	0.00	00:00
Black friday	TVAE	0.03	0.46	01:55
	CTGAN	0.08	1.18	04:54
	TabDDPM	0.01	0.13	00:23
	TabSyn	0.21	2.40	09:19
	SMOTE	0.00	0.06	00:19
	UC-SMOTE	0.01	0.22	01:07
Cardio	TVAE	0.01	0.16	00:41
	CTGAN	0.03	0.51	02:07
	TabDDPM	0.01	0.08	00:12
	TabSyn	0.10	1.38	04:13
	SMOTE	0.00	0.00	00:00
	UC-SMOTE	0.00	0.00	00:00
Churn	TVAE	0.00	0.03	00:07
	CTGAN	0.02	0.29	01:11
	TabDDPM	0.02	0.31	00:48
	TabSyn	0.04	0.53	02:02
	SMOTE	0.00	0.00	00:00
	UC-SMOTE	0.00	0.00	00:00

Tuning costs estimation (continued).

Dataset	Model	Metrics		
		Energy (kWh) ↓	Emissions (Kg) ↓	Duration (HH:MM) ↓
Covertypes	TVAE	0.19	2.77	11:34
	CTGAN	0.54	8.04	33:17
	TabDDPM	0.03	0.49	01:48
	TabSyn	0.19	1.53	07:54
	SMOTE	0.02	0.26	01:19
	UC-SMOTE	0.02	0.28	01:27
Diamonds	TVAE	0.01	0.10	00:25
	CTGAN	0.03	0.40	01:39
	TabDDPM	0.01	0.15	00:24
	TabSyn	0.08	1.01	03:22
	SMOTE	0.00	0.01	00:02
	UC-SMOTE	0.00	0.00	00:00
Heloc	TVAE	0.00	0.04	00:10
	CTGAN	0.03	0.49	02:02
	TabDDPM	0.00	0.03	00:05
	TabSyn	0.04	0.54	02:08
	SMOTE	0.00	0.00	00:00
	UC-SMOTE	0.00	0.00	00:00
Higgs	TVAE	0.04	0.53	02:13
	CTGAN	0.11	1.61	06:42
	TabDDPM	0.01	0.10	00:15
	TabSyn	0.18	1.58	07:50
	SMOTE	0.00	0.00	00:00
	UC-SMOTE	0.00	0.00	00:01
House 16h	TVAE	0.01	0.09	00:23
	CTGAN	0.01	0.10	00:25
	TabDDPM	0.00	0.04	00:06
	TabSyn	0.05	0.58	02:16
	SMOTE	0.00	0.00	00:00
	UC-SMOTE	0.00	0.00	00:00
Insurance	TVAE	0.00	0.00	00:00
	CTGAN	0.00	0.03	00:06
	TabDDPM	0.00	0.05	00:12
	TabSyn	0.03	0.35	01:25
	SMOTE	0.00	0.00	00:00
	UC-SMOTE	0.00	0.00	00:00
King	TVAE	0.01	0.09	00:23
	CTGAN	0.01	0.21	00:51
	TabDDPM	0.01	0.17	00:27
	TabSyn	0.05	0.66	02:26
	SMOTE	0.00	0.00	00:01
	UC-SMOTE	0.00	0.00	00:00

Tuning costs estimation (continued).

Dataset	Model	Metrics		
		Energy (kWh) ↓	Emissions (Kg) ↓	Duration (HH:MM) ↓
Miniboo ne	TVAE	0.04	0.56	02:19
	CTGAN	0.05	0.73	02:59
	TabDDPM	0.01	0.10	00:14
	TabSyn	0.31	2.02	11:44
	SMOTE	0.00	0.00	00:00
	UC-SMOTE	0.00	0.01	00:02
Moons	TVAE	0.01	0.08	00:20
	CTGAN	0.02	0.29	01:12
	TabDDPM	0.01	0.08	00:12
	TabSyn	0.05	0.87	02:35
	SMOTE	0.00	0.00	00:00
	UC-SMOTE	0.00	0.00	00:00

Table D.14: Estimated hyperparameter search cost based on the estimated training cost of the best models presented in Table D.13. All costs associated with TabSyn include those incurred by the VAE and the denoiser. The energy and emissions values are rounded to two decimals and take into account the number of trials we could run in parallel per model. The best values per metric are formatted in **bold green** and the worse are in **red**.

Appendix E. Base Models and their Tuned Versions

We also trained the models using their native codes and hyperparameters to provide an additional reference for comparison with their tuned versions. Table E.15 presents the results as evaluated under the same procedure as the tuned models for C2ST, DCR-Rate, ML-Efficacy, *column-wise similarity* (named "Shape"), and *pair-wise correlation* (named "Pair"). For TabDDPM, since there is no default hyperparameters provided [8], we fixed one based on the base configuration provided in the authors' GitHub repository.

Dataset	Model	Metrics				
		C2ST ↓	DCR-Rate ↓	ML-Efficacy ↑	Shape ↑	Pair ↑
Abalone	Train Copy	0.51 ± 0.00	1.00 ± 0.00	0.23 ± 0.01	0.96 ± 0.01	0.88 ± 0.01
	CTGAN	0.99 ± 0.01	0.60 ± 0.01	0.12 ± 0.03	0.87 ± 0.02	0.76 ± 0.01
	TVAE	0.96 ± 0.00	0.61 ± 0.02	0.22 ± 0.02	0.91 ± 0.02	0.83 ± 0.03
	TabDDPM	1.00 ± 0.00	0.64 ± 0.01	0.00 ± 0.00	0.85 ± 0.01	0.70 ± 0.01
	TabSyn	0.78 ± 0.02	0.64 ± 0.01	0.22 ± 0.01	0.95 ± 0.01	0.88 ± 0.01
Adult	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.71 ± 0.01	0.99 ± 0.00	0.98 ± 0.00
	CTGAN	0.96 ± 0.01	0.61 ± 0.01	0.61 ± 0.05	0.88 ± 0.02	0.82 ± 0.02
	TVAE	0.94 ± 0.01	0.61 ± 0.00	0.63 ± 0.03	0.92 ± 0.00	0.85 ± 0.01
	TabDDPM	0.66 ± 0.01	0.62 ± 0.00	0.67 ± 0.00	0.98 ± 0.00	0.95 ± 0.00

Results per datasets and models under diverse metrics for the *base models* (continued).

Dataset	Model	Metrics				
		C2ST ↓	DCR-Rate ↓	ML-Efficacy ↑	Shape ↑	Pair ↑
	TabSyn	0.71 ± 0.06	0.62 ± 0.00	0.66 ± 0.01	0.98 ± 0.01	0.95 ± 0.02
Bank marketing	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.54 ± 0.01	0.99 ± 0.00	0.98 ± 0.01
	CTGAN	0.89 ± 0.01	0.62 ± 0.01	0.33 ± 0.05	0.92 ± 0.00	0.86 ± 0.01
	TVAE	0.95 ± 0.00	0.62 ± 0.00	0.52 ± 0.02	0.91 ± 0.01	0.84 ± 0.02
	TabDDPM	0.68 ± 0.01	0.63 ± 0.01	0.48 ± 0.02	0.99 ± 0.01	0.96 ± 0.00
	TabSyn	0.65 ± 0.04	0.63 ± 0.01	0.47 ± 0.01	0.98 ± 0.01	0.96 ± 0.00
Bike sharing	Train Copy	0.50 ± 0.01	1.00 ± 0.00	0.94 ± 0.00	0.99 ± 0.00	0.53 ± 0.00
	CTGAN	1.00 ± 0.00	0.61 ± 0.00	0.37 ± 0.03	0.93 ± 0.01	0.51 ± 0.01
	TVAE	1.00 ± 0.00	0.61 ± 0.01	0.42 ± 0.08	0.91 ± 0.01	0.51 ± 0.01
	TabDDPM	0.85 ± 0.01	0.61 ± 0.00	0.67 ± 0.03	0.98 ± 0.00	0.53 ± 0.00
	TabSyn	0.94 ± 0.03	0.61 ± 0.00	0.31 ± 0.04	0.95 ± 0.01	0.52 ± 0.01
Black friday	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.53 ± 0.01	1.00 ± 0.00	0.99 ± 0.00
	CTGAN	0.92 ± 0.01	0.66 ± 0.01	0.40 ± 0.01	0.94 ± 0.01	0.84 ± 0.01
	TVAE	0.98 ± 0.00	0.67 ± 0.01	0.29 ± 0.08	0.83 ± 0.01	0.73 ± 0.01
	TabDDPM	0.92 ± 0.00	0.67 ± 0.01	0.31 ± 0.03	0.99 ± 0.00	0.98 ± 0.00
	TabSyn	0.91 ± 0.01	0.68 ± 0.00	0.13 ± 0.03	0.98 ± 0.00	0.96 ± 0.01
Cardio	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.72 ± 0.01	1.00 ± 0.00	0.98 ± 0.01
	CTGAN	1.00 ± 0.01	0.62 ± 0.01	0.67 ± 0.02	0.93 ± 0.01	0.93 ± 0.01
	TVAE	1.00 ± 0.00	0.63 ± 0.00	0.67 ± 0.04	0.88 ± 0.01	0.90 ± 0.02
	TabDDPM	0.59 ± 0.00	0.64 ± 0.00	0.72 ± 0.01	0.99 ± 0.00	0.96 ± 0.02
	TabSyn	0.58 ± 0.01	0.64 ± 0.01	0.72 ± 0.01	0.99 ± 0.01	0.97 ± 0.01
Churn	Train Copy	0.50 ± 0.01	1.00 ± 0.00	0.59 ± 0.02	0.95 ± 0.00	0.87 ± 0.01
	CTGAN	0.87 ± 0.02	0.60 ± 0.00	0.19 ± 0.10	0.85 ± 0.02	0.79 ± 0.00
	TVAE	0.98 ± 0.01	0.60 ± 0.01	0.45 ± 0.10	0.75 ± 0.03	0.63 ± 0.03
	TabDDPM	1.00 ± 0.00	0.65 ± 0.40	0.00 ± 0.00	0.59 ± 0.02	0.56 ± 0.02
	TabSyn	0.64 ± 0.03	0.61 ± 0.01	0.43 ± 0.13	0.92 ± 0.01	0.85 ± 0.00
Covertypes	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.90 ± 0.00	1.00 ± 0.00	1.00 ± 0.01
	CTGAN	1.00 ± 0.00	0.60 ± 0.00	0.66 ± 0.01	0.97 ± 0.00	0.94 ± 0.01
	TVAE	1.00 ± 0.00	0.60 ± 0.00	0.72 ± 0.01	0.98 ± 0.00	0.95 ± 0.01
	TabDDPM	0.82 ± 0.01	0.60 ± 0.00	0.73 ± 0.01	1.00 ± 0.00	0.99 ± 0.00
	TabSyn	0.72 ± 0.04	0.60 ± 0.00	0.79 ± 0.03	0.99 ± 0.00	0.99 ± 0.00
Diamonds	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.98 ± 0.00	0.99 ± 0.00	0.77 ± 0.01
	CTGAN	0.98 ± 0.00	0.60 ± 0.00	0.87 ± 0.01	0.91 ± 0.01	0.75 ± 0.04
	TVAE	0.97 ± 0.01	0.60 ± 0.01	0.91 ± 0.01	0.92 ± 0.02	0.76 ± 0.03
	TabDDPM	0.76 ± 0.01	0.61 ± 0.00	0.97 ± 0.00	0.99 ± 0.01	0.71 ± 0.01
	TabSyn	0.95 ± 0.02	0.64 ± 0.00	0.19 ± 0.06	0.96 ± 0.00	0.72 ± 0.01
Heloc	Train Copy	0.50 ± 0.01	1.00 ± 0.00	0.70 ± 0.01	0.99 ± 0.00	0.97 ± 0.01
	CTGAN	1.00 ± 0.00	0.63 ± 0.00	0.40 ± 0.14	0.87 ± 0.05	0.87 ± 0.07
	TVAE	0.98 ± 0.01	0.63 ± 0.02	0.69 ± 0.01	0.90 ± 0.00	0.75 ± 0.01
	TabDDPM	0.71 ± 0.01	0.64 ± 0.02	0.70 ± 0.01	0.97 ± 0.00	0.94 ± 0.01
	TabSyn	0.75 ± 0.01	0.65 ± 0.01	0.69 ± 0.01	0.97 ± 0.00	0.96 ± 0.01
Higgs	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.74 ± 0.00	0.99 ± 0.00	0.99 ± 0.00
	CTGAN	1.00 ± 0.00	0.60 ± 0.00	0.61 ± 0.08	0.87 ± 0.02	0.95 ± 0.00
	TVAE	1.00 ± 0.00	0.61 ± 0.02	0.68 ± 0.02	0.78 ± 0.00	0.94 ± 0.01
	TabDDPM	0.79 ± 0.02	0.60 ± 0.01	0.72 ± 0.01	0.98 ± 0.01	0.93 ± 0.02
	TabSyn	0.59 ± 0.01	0.61 ± 0.00	0.73 ± 0.00	0.94 ± 0.00	0.99 ± 0.00

Results per datasets and models under diverse metrics for the *base models* (continued).

Dataset	Model	Metrics				
		C2ST ↓	DCR-Rate ↓	ML-Efficacy ↑	Shape ↑	Pair ↑
House 16h	Train Copy	0.50 ± 0.01	1.00 ± 0.00	0.64 ± 0.01	0.99 ± 0.01	0.99 ± 0.00
	CTGAN	1.00 ± 0.00	0.61 ± 0.02	0.20 ± 0.05	0.86 ± 0.01	0.95 ± 0.01
	TVAE	0.98 ± 0.00	0.61 ± 0.01	0.36 ± 0.02	0.89 ± 0.00	0.96 ± 0.01
	TabDDPM	0.63 ± 0.01	0.61 ± 0.01	0.57 ± 0.01	0.98 ± 0.00	0.97 ± 0.01
	TabSyn	0.84 ± 0.03	0.62 ± 0.00	0.33 ± 0.06	0.96 ± 0.01	0.98 ± 0.01
Insurance	Train Copy	0.48 ± 0.01	1.00 ± 0.00	0.85 ± 0.03	0.96 ± 0.01	0.91 ± 0.01
	CTGAN	0.92 ± 0.01	0.67 ± 0.06	-0.18 ± 0.05	0.84 ± 0.00	0.81 ± 0.01
	TVAE	0.87 ± 0.02	0.62 ± 0.02	0.61 ± 0.03	0.85 ± 0.03	0.77 ± 0.03
	TabDDPM	0.58 ± 0.01	0.60 ± 0.02	0.83 ± 0.02	0.95 ± 0.01	0.90 ± 0.01
	TabSyn	0.60 ± 0.03	0.61 ± 0.02	0.83 ± 0.01	0.94 ± 0.01	0.89 ± 0.01
King	Train Copy	0.50 ± 0.01	1.00 ± 0.00	0.87 ± 0.00	0.98 ± 0.01	0.97 ± 0.01
	CTGAN	1.00 ± 0.00	0.60 ± 0.01	0.54 ± 0.04	0.89 ± 0.01	0.93 ± 0.00
	TVAE	0.99 ± 0.01	0.60 ± 0.00	0.70 ± 0.02	0.92 ± 0.01	0.91 ± 0.01
	TabDDPM	1.00 ± 0.00	0.91 ± 0.14	-192.96 ± 153.37	0.33 ± 0.03	0.74 ± 0.02
	TabSyn	0.97 ± 0.01	0.62 ± 0.01	0.05 ± 0.04	0.96 ± 0.01	0.94 ± 0.00
Miniboo ne	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.89 ± 0.01	0.99 ± 0.00	0.99 ± 0.01
	CTGAN	1.00 ± 0.00	0.60 ± 0.00	0.54 ± 0.05	0.80 ± 0.03	0.56 ± 0.01
	TVAE	1.00 ± 0.00	0.60 ± 0.00	0.82 ± 0.01	0.91 ± 0.02	0.56 ± 0.00
	TabDDPM	0.82 ± 0.01	0.63 ± 0.01	0.88 ± 0.00	0.95 ± 0.02	0.83 ± 0.03
	TabSyn	0.71 ± 0.09	0.60 ± 0.00	0.89 ± 0.01	0.98 ± 0.01	0.89 ± 0.05
Moons	Train Copy	0.50 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	0.99 ± 0.00	0.99 ± 0.00
	CTGAN	0.93 ± 0.01	0.60 ± 0.05	1.00 ± 0.00	0.94 ± 0.01	0.70 ± 0.03
	TVAE	0.82 ± 0.02	0.61 ± 0.03	1.00 ± 0.00	0.97 ± 0.01	0.76 ± 0.02
	TabDDPM	0.61 ± 0.16	0.53 ± 0.15	0.97 ± 0.06	0.93 ± 0.10	0.90 ± 0.16
	TabSyn	0.58 ± 0.03	0.61 ± 0.01	1.00 ± 0.00	0.98 ± 0.01	0.98 ± 0.01

Table E.15: Results for *base models*. Models are trained using their default hyperparameters as provided by the authors in their papers. The best values per metric are formatted in **bold green** and the worse values are in **red**.

Appendix F. Quick Search on Reduced Hyperparameter Space

With the reduced hyperparameters search spaces presented in Appendix A, we ran a *light* hyperparameter tuning described in Section 5. The experiment was done on all dataset (Table 2). We performed 50 trials per fold with 3 fold, meaning we ran a total of 150 trials for this light experiment. The results per datasets are shown in Table F.16.

Dataset	Model	Metrics					
		C2ST ↓	DCR-Rate ↓	ML-Efficacy ↑	Shape ↑	Pair ↑	Sampling time ↓
Abalone	Train Copy	0.51 ± 0.00	1.00 ± 0.00	0.23 ± 0.01	0.96 ± 0.01	0.88 ± 0.01	-
	CTGAN	0.71 ± 0.03	0.63 ± 0.00	0.17 ± 0.01	0.93 ± 0.02	0.86 ± 0.04	00 ± 0.01
	TVAE	0.64 ± 0.02	0.67 ± 0.00	0.23 ± 0.01	0.93 ± 0.02	0.89 ± 0.03	00 ± 0.00
	TabDDPM	0.78 ± 0.01	0.69 ± 0.03	0.23 ± 0.01	0.95 ± 0.00	0.88 ± 0.02	03 ± 0.88
	TabSyn	0.80 ± 0.00	0.62 ± 0.01	0.19 ± 0.03	0.93 ± 0.01	0.85 ± 0.01	00 ± 0.01
Adult	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.71 ± 0.01	0.99 ± 0.00	0.98 ± 0.00	-
	CTGAN	0.77 ± 0.01	0.72 ± 0.00	0.65 ± 0.00	0.96 ± 0.01	0.89 ± 0.02	00 ± 0.02
	TVAE	0.77 ± 0.01	0.72 ± 0.01	0.65 ± 0.02	0.96 ± 0.01	0.93 ± 0.00	00 ± 0.00
	TabDDPM	0.67 ± 0.01	0.62 ± 0.00	0.67 ± 0.01	0.97 ± 0.01	0.94 ± 0.01	11 ± 0.63
	TabSyn	0.73 ± 0.01	0.62 ± 0.00	0.66 ± 0.01	0.97 ± 0.01	0.94 ± 0.01	02 ± 0.01
Bank marketing	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.54 ± 0.01	0.99 ± 0.00	0.98 ± 0.01	-
	CTGAN	0.72 ± 0.01	0.63 ± 0.00	0.45 ± 0.04	0.98 ± 0.00	0.94 ± 0.01	00 ± 0.01
	TVAE	0.78 ± 0.02	0.63 ± 0.00	0.47 ± 0.03	0.97 ± 0.01	0.94 ± 0.00	00 ± 0.00
	TabDDPM	0.68 ± 0.00	0.63 ± 0.01	0.52 ± 0.03	0.97 ± 0.00	0.95 ± 0.01	16 ± 2.18
	TabSyn	0.75 ± 0.00	0.62 ± 0.00	0.42 ± 0.04	0.97 ± 0.01	0.95 ± 0.00	02 ± 0.01
Bike sharing	Train Copy	0.50 ± 0.01	1.00 ± 0.00	0.94 ± 0.00	0.99 ± 0.00	0.53 ± 0.00	-
	CTGAN	0.89 ± 0.01	0.61 ± 0.01	0.71 ± 0.01	0.98 ± 0.01	0.53 ± 0.00	00 ± 0.00
	TVAE	0.81 ± 0.01	0.61 ± 0.01	0.69 ± 0.03	0.96 ± 0.00	0.53 ± 0.00	00 ± 0.00
	TabDDPM	0.80 ± 0.02	0.63 ± 0.01	0.83 ± 0.01	0.96 ± 0.02	0.53 ± 0.01	06 ± 1.34
	TabSyn	0.91 ± 0.01	0.61 ± 0.01	0.32 ± 0.06	0.95 ± 0.01	0.52 ± 0.00	00 ± 0.01
Black friday	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.53 ± 0.01	1.00 ± 0.00	0.99 ± 0.00	-
	CTGAN	0.87 ± 0.01	0.68 ± 0.00	0.36 ± 0.01	0.97 ± 0.01	0.91 ± 0.02	01 ± 0.17
	TVAE	0.95 ± 0.01	0.68 ± 0.00	0.27 ± 0.02	0.97 ± 0.00	0.95 ± 0.00	00 ± 0.02
	TabDDPM	0.89 ± 0.01	0.67 ± 0.00	0.44 ± 0.02	0.97 ± 0.01	0.96 ± 0.01	50 ± 17.19
	TabSyn	0.89 ± 0.01	0.68 ± 0.00	0.16 ± 0.01	0.97 ± 0.01	0.95 ± 0.01	07 ± 0.01
Cardio	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.72 ± 0.01	1.00 ± 0.00	0.98 ± 0.01	-
	CTGAN	0.64 ± 0.01	0.64 ± 0.01	0.72 ± 0.01	0.99 ± 0.00	0.96 ± 0.00	00 ± 0.06
	TVAE	0.71 ± 0.01	0.64 ± 0.01	0.73 ± 0.01	0.96 ± 0.01	0.96 ± 0.01	00 ± 0.01
	TabDDPM	0.57 ± 0.01	0.64 ± 0.00	0.72 ± 0.01	0.99 ± 0.00	0.98 ± 0.01	12 ± 4.58
	TabSyn	0.59 ± 0.01	0.64 ± 0.00	0.72 ± 0.01	0.98 ± 0.01	0.97 ± 0.01	03 ± 0.01
Churn	Train Copy	0.50 ± 0.01	1.00 ± 0.00	0.59 ± 0.02	0.95 ± 0.00	0.87 ± 0.01	-
	CTGAN	0.63 ± 0.01	0.63 ± 0.01	0.39 ± 0.01	0.93 ± 0.00	0.84 ± 0.01	00 ± 0.01
	TVAE	0.64 ± 0.00	0.63 ± 0.01	0.53 ± 0.01	0.92 ± 0.00	0.84 ± 0.00	00 ± 0.00
	TabDDPM	0.98 ± 0.03	0.67 ± 0.18	0.02 ± 0.03	0.59 ± 0.05	0.55 ± 0.02	18 ± 1.06
	TabSyn	0.68 ± 0.06	0.61 ± 0.00	0.46 ± 0.03	0.91 ± 0.01	0.84 ± 0.01	00 ± 0.01

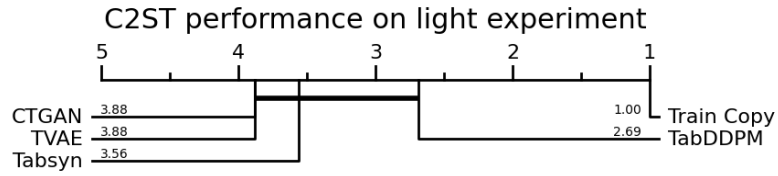
Results per datasets and models under diverse metrics for the light search (continued).

Dataset	Model	Metrics					
		C2ST ↓	DCR-Rate ↓	ML-Efficacy ↑	Shape ↑	Pair ↑	Sampling time ↓
Coverttype	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.90 ± 0.00	1.00 ± 0.00	1.00 ± 0.01	-
	CTGAN	0.98 ± 0.00	0.60 ± 0.00	0.69 ± 0.00	0.98 ± 0.00	0.95 ± 0.01	06 ± 1.58
	TVAE	0.90 ± 0.01	0.60 ± 0.00	0.77 ± 0.01	0.98 ± 0.00	0.96 ± 0.00	02 ± 0.07
	TabDDPM	0.97 ± 0.01	0.64 ± 0.02	0.69 ± 0.00	0.94 ± 0.02	0.89 ± 0.03	355 ± 10.90
	TabSyn	0.87 ± 0.00	0.61 ± 0.01	0.65 ± 0.03	0.98 ± 0.00	0.68 ± 0.00	02 ± 0.02
Diamonds	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.98 ± 0.00	0.99 ± 0.00	0.77 ± 0.01	-
	CTGAN	0.89 ± 0.01	0.64 ± 0.01	0.94 ± 0.00	0.95 ± 0.00	0.74 ± 0.02	00 ± 0.08
	TVAE	0.76 ± 0.02	0.64 ± 0.00	0.96 ± 0.01	0.95 ± 0.01	0.73 ± 0.01	00 ± 0.01
	TabDDPM	0.75 ± 0.01	0.61 ± 0.01	0.97 ± 0.00	0.97 ± 0.01	0.70 ± 0.02	10 ± 1.26
	TabSyn	0.92 ± 0.01	0.64 ± 0.01	0.65 ± 0.09	0.95 ± 0.01	0.72 ± 0.02	02 ± 0.01
Heloc	Train Copy	0.50 ± 0.01	1.00 ± 0.00	0.70 ± 0.01	0.99 ± 0.00	0.97 ± 0.01	-
	CTGAN	0.93 ± 0.02	0.64 ± 0.01	0.70 ± 0.01	0.97 ± 0.00	0.93 ± 0.01	00 ± 0.02
	TVAE	0.92 ± 0.01	0.63 ± 0.01	0.69 ± 0.02	0.94 ± 0.00	0.88 ± 0.01	00 ± 0.00
	TabDDPM	0.71 ± 0.02	0.67 ± 0.01	0.69 ± 0.02	0.95 ± 0.01	0.92 ± 0.03	03 ± 1.99
	TabSyn	0.82 ± 0.01	0.65 ± 0.01	0.69 ± 0.01	0.95 ± 0.01	0.95 ± 0.01	00 ± 0.01
Higgs	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.74 ± 0.00	0.99 ± 0.00	0.99 ± 0.00	-
	CTGAN	0.87 ± 0.00	0.60 ± 0.01	0.70 ± 0.01	0.98 ± 0.00	0.98 ± 0.00	00 ± 0.09
	TVAE	0.92 ± 0.01	0.60 ± 0.01	0.71 ± 0.02	0.94 ± 0.01	0.97 ± 0.00	00 ± 0.03
	TabDDPM	0.61 ± 0.02	0.62 ± 0.01	0.73 ± 0.00	0.98 ± 0.01	0.95 ± 0.01	08 ± 1.00
	TabSyn	0.76 ± 0.05	0.60 ± 0.01	0.71 ± 0.02	0.92 ± 0.01	0.98 ± 0.00	04 ± 0.01
House 16h	Train Copy	0.50 ± 0.01	1.00 ± 0.00	0.64 ± 0.01	0.99 ± 0.01	0.99 ± 0.00	-
	CTGAN	0.84 ± 0.02	0.62 ± 0.00	0.45 ± 0.02	0.97 ± 0.01	0.97 ± 0.01	00 ± 0.01
	TVAE	0.84 ± 0.00	0.62 ± 0.01	0.46 ± 0.03	0.94 ± 0.00	0.98 ± 0.00	00 ± 0.00
	TabDDPM	0.64 ± 0.01	0.61 ± 0.01	0.61 ± 0.01	0.96 ± 0.01	0.94 ± 0.03	02 ± 0.75
	TabSyn	0.83 ± 0.03	0.62 ± 0.01	0.35 ± 0.01	0.95 ± 0.01	0.99 ± 0.01	01 ± 0.01
Insurance	Train Copy	0.48 ± 0.01	1.00 ± 0.00	0.85 ± 0.03	0.96 ± 0.01	0.91 ± 0.01	-
	CTGAN	0.67 ± 0.02	0.92 ± 0.01	0.66 ± 0.04	0.94 ± 0.01	0.85 ± 0.02	00 ± 0.01
	TVAE	0.65 ± 0.02	0.91 ± 0.02	0.80 ± 0.03	0.93 ± 0.01	0.88 ± 0.01	00 ± 0.00
	TabDDPM	0.61 ± 0.02	0.62 ± 0.02	0.83 ± 0.03	0.94 ± 0.01	0.89 ± 0.02	06 ± 2.76
	TabSyn	0.63 ± 0.02	0.60 ± 0.03	0.80 ± 0.01	0.93 ± 0.01	0.88 ± 0.02	00 ± 0.01
King	Train Copy	0.50 ± 0.01	1.00 ± 0.00	0.87 ± 0.00	0.98 ± 0.01	0.97 ± 0.01	-
	CTGAN	0.96 ± 0.01	0.62 ± 0.00	0.66 ± 0.03	0.97 ± 0.00	0.94 ± 0.01	00 ± 0.13
	TVAE	0.95 ± 0.01	0.61 ± 0.01	0.80 ± 0.01	0.95 ± 0.01	0.94 ± 0.00	00 ± 0.01
	TabDDPM	1.00 ± 0.00	0.70 ± 0.46	0.03 ± 0.68	0.28 ± 0.07	0.76 ± 0.00	05 ± 0.43
	TabSyn	0.98 ± 0.01	0.62 ± 0.01	-0.01 ± 0.01	0.94 ± 0.01	0.93 ± 0.00	01 ± 0.01
Miniboo ne	Train Copy	0.50 ± 0.00	1.00 ± 0.00	0.89 ± 0.01	0.99 ± 0.00	0.99 ± 0.01	-
	CTGAN	0.90 ± 0.01	0.60 ± 0.01	0.87 ± 0.01	0.93 ± 0.01	0.58 ± 0.01	01 ± 0.03
	TVAE	0.94 ± 0.01	0.60 ± 0.00	0.87 ± 0.00	0.94 ± 0.00	0.59 ± 0.01	00 ± 0.03
	TabDDPM	0.73 ± 0.01	0.62 ± 0.00	0.89 ± 0.01	0.96 ± 0.01	0.83 ± 0.05	12 ± 1.35
	TabSyn	0.87 ± 0.01	0.60 ± 0.01	0.87 ± 0.00	0.96 ± 0.00	0.76 ± 0.14	06 ± 0.08

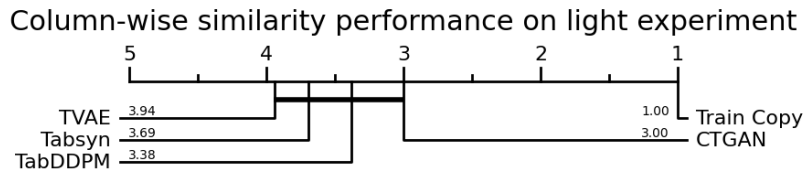
Results per datasets and models under diverse metrics for the light search (continued).

Dataset	Model	Metrics					
		C2ST ↓	DCR-Rate ↓	ML-Efficacy ↑	Shape ↑	Pair ↑	Sampling time ↓
Moons	Train Copy	0.50 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	0.99 ± 0.00	0.99 ± 0.00	-
	CTGAN	0.70 ± 0.02	0.62 ± 0.02	1.00 ± 0.00	0.96 ± 0.01	0.96 ± 0.01	00 ± 0.02
	TVAE	0.69 ± 0.01	0.62 ± 0.01	1.00 ± 0.00	0.96 ± 0.01	0.93 ± 0.02	00 ± 0.00
	TabDDPM	0.53 ± 0.01	0.61 ± 0.01	1.00 ± 0.00	0.99 ± 0.00	0.98 ± 0.01	08 ± 2.13
	TabSyn	0.59 ± 0.01	0.61 ± 0.01	1.00 ± 0.00	0.98 ± 0.01	0.97 ± 0.01	01 ± 0.00

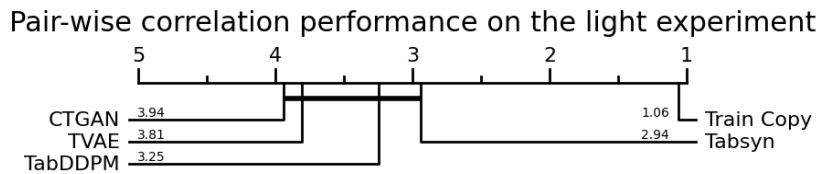
Table F.16: Light experiment results under various metrics. Results are averaged over 3 folds with 5 synthetic samples per fold as done in the extensive hyperparameter tuning. The best values per metric are formatted in **bold green** and the worse values are in **red**.



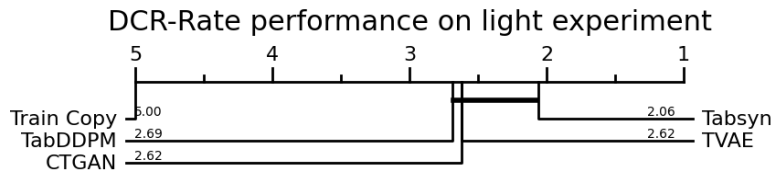
(a) C2ST



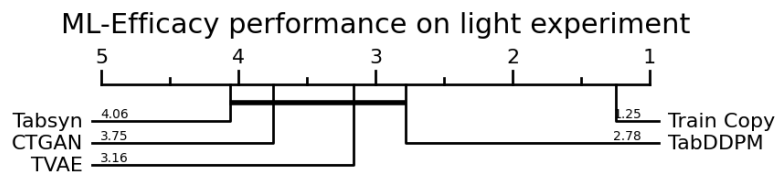
(b) Column-wise similarity



(c) Pair-wise correlation



(d) DCR-Rate



(e) ML-Efficacy

Figure F.12: Models' ranking under the light experiment setup with critical difference diagrams on C2ST, *column-wise similarity*, *pair-wise correlation*, DCR-Rate, and ML-Efficacy metrics over all datasets.