
Generating Privacy-Preserving Synthetic Tabular Data Using Oblivious Variational Autoencoders

L Vivek Harsha Vardhan¹ Stanley Kok¹

Abstract

A vast amount of the world’s data is stored in tabular form, and a significant portion of it contains personally identifiable attributes (e.g., electronic health records and financial records). The public sharing of such data could expedite AI research for resolving pressing societal, medical and other problems, but attempts to do so are riddled with obstacles because of privacy concerns. To circumvent the privacy issues, a line of statistical and deep learning models seeks to generate synthetic tabular data, which adhere to the distributional statistics of the original data, without disclosing any personally identifiable information. Those preliminary models show promising results, but have not attained their full potential because they have yet to exploit a strong inductive bias that predisposes them to tabular data. We propose a new model called Oblivious Variational Autoencoder (OVAE) that combines variational autoencoders (VAEs) with a differentiable version of oblivious decision trees (ODTs) (Lou & Obukhov, 2017). Boosted ODTs have been highly successful for predictive tasks on tabular data, outperforming even deep learning systems in many Kaggle competitions. However, their use in generative models has largely been overlooked. OVAE incorporates ODTs’ amenability to tabular data as an inductive bias in VAEs, thereby generating synthetic tabular data of high fidelity to the original tables. In an extensive set of experiments, OVAE demonstrates its efficacy and surpasses several state-of-the-art models on a wide range of datasets.

¹Department of Information Systems and Analytics, School of Computing, National University of Singapore, Singapore. Correspondence to: L Vivek Harsha Vardhan <harsha@comp.nus.edu.sg>.

1. Introduction

Businesses, governments, hospitals, and other organizations store a wealth of digital data in the form of tables. If these data could be publicly shared, they would expedite data-intensive AI research to the benefit of the organizations or society at large. However, these tabular data usually contain personal identifiers, which in conjunction with sensitive financial and medical information, pose serious privacy concerns. Thus access to such data are usually hampered by the high walls and deep moats of security processes, such as lengthy reviews by institutional review boards.

To mitigate the privacy risks and to expedite data release, some early approaches sought to suppress, randomize, or perturb potentially identifiable information. However, such techniques were found to be susceptible to re-identification attacks, e.g., via background knowledge (Emam et al., 2011). In recent years, a series of research efforts have addressed the problem from another angle, mitigating the privacy risks through the generation of synthetic data that closely mimic the true underlying tabular data. Because the generated data are completely fake, there is little risk of personal attribute disclosure. Those approaches (Choi et al., 2017; Srivastava et al., 2017; Park et al., 2018; Xu et al., 2019) typically utilize deep generative models (such as variational autoencoders (Kingma & Welling, 2014) and generative adversarial networks (Goodfellow et al., 2014)) in the hope that the models’ success in generating synthetic images (Lin et al., 2019) and text (Guo et al., 2018) would carry over to tabular data. Those approaches have shown promising empirical results, but have yet to achieve their full potential. Many of the deep generative models are transplanted almost wholesale from their initial image or text domains to the tabular one. Consequently, they inherit inductive biases that are more amenable to their original data types than tabular data. We postulate that we can improve the fidelity of the generated data to their true underlying tables by supplementing the deep generative models with an appropriate tabular inductive bias, such as one based on decision trees or their variants.

Decision trees and their variants have proven to be highly successful for the discriminative modeling of tabular data. For example, Catboost (Prokhorenkova et al., 2018), which

uses boosted oblivious decision trees (ODTs) (Lou & Obukhov, 2017), is used in the winning solutions of many Kaggle competitions involving tabular data (frequently beating even deep learning competitors). Though effective for predictive tasks on tables, decision trees and their variants have not been extensively used for the generative modeling of tabular data.

Combining the best of deep generative models and decision tree variants, we propose the Oblivious Variational Autoencoder (OVAE). OVAE embeds ‘softened’ oblivious decision trees (ODTs) in a variational autoencoder (VAE) to both encode tabular data into a latent representation, and to decode (generate) synthetic data from that representation. ODTs are good at representing decision manifolds that approximate the hyperplane boundaries usually present in tabular data. Thus the ODTs imbue a strong inductive bias for tabular data in OVAE, allowing it to preserve the distributional characteristics of the original tabular data in its generated synthetic data.

In sum, our contributions are as follows.

- We propose a new model called OVAE that combines ODTs with a VAE to generate tabular data. To our knowledge, we are the first to adapt a decision tree variant so that it can be used with a VAE for generating privacy-preserving synthetic tables.
- We extensively compare our OVAE model against several state-of-the-art baselines, and show that OVAE compares favorably against the baselines on 12 real-world datasets.

2. Related Work

Early approaches for synthetic table generation utilize statistical models. Those models typically form a multivariate probability distribution over all columns in a table, each of which is regarded as a random variable. To generate synthetic data, the models draw a sample (a row of values) from the distribution, with one value per random variable (column). Examples of such statistical models include Bayesian networks (e.g., CLBN (Chow & Liu, 1968), PrivBayes (Zhang et al., 2017), and (Aviñó et al., 2018)) and copulas (Patki et al., 2016; Sun et al., 2019). The former has been used to generate discrete variables, while the latter has been used to generate (non-linearly correlated) continuous variables. A drawback of those approaches is that they either model discrete data or continuous data but not both, a restriction that makes them unsuitable for a wide range of real-world data that contain both discrete and continuous values. Further, those approaches generally do not exploit the parallel processing afforded by modern graphical processing units (unlike deep learning models), a shortcoming that limits their scalability.

Another line of models for tabular data generation is based on deep learning, and is predominantly built upon generative adversarial networks (GANs) (Goodfellow et al., 2014). GANs are proposed to generate multi-categorical columns of discrete data in (Camino et al., 2018), and continuous laboratory time series data in (Yahi et al., 2017). ehrGAN (Che et al., 2017) generates data to augment scarce medical records in a semi-supervised manner. medGAN (Choi et al., 2017) combines a GAN and an autoencoder to model both continuous and binary table columns. VEEGAN (Srivastava et al., 2017) ameliorates the mode collapse problem in GANs through variational learning. tableGAN (Park et al., 2018) uses a convolutional neural network as the discriminator in its GAN to maximize the quality of a table’s label column. PATE-GAN (Yoon et al., 2019) modifies the Private Aggregation of Teacher Ensembles (PATE) framework (Papernot et al., 2017), and uses it to enforce differential privacy in the generator component of a GAN. Most recently, CTGAN (Xu et al., 2019) uses *mode-specific normalization* (MSN) for modeling continuous data with multiple modes, and introduces techniques for correctly modeling the minor categories in categorical data with skewed distributions.

Another deep learning model that has been used for synthetic tabular data generation is the variational autoencoder (VAE). TVAE (Xu et al., 2019) is a vanilla VAE with fully connected neural networks in both its encoder and decoder. It also uses MSN to address the problems of mode collapse. In an extensive empirical comparison, TVAE outperforms CTGAN and other state-of-the-art GAN models, and is a frontrunner that we compare against.

Unlike all of the aforementioned models that do not impose strong tabular constraints, our proposed oblivious variational autoencoder (OVAE) incorporates a strong inductive bias for tabular data in the form of “softened” oblivious decision trees (ODTs) into a VAE. OVAE is similar to TVAE in that both use VAE as a building block, but OVAE differs by using ODTs in place of TVAE’s standard feed-forward neural networks.

There exists a body of work integrating inductive biases in the form of grammars, templates, or other constraints into deep generative models, but those biases are primarily applicable to domains other than tabular data. GVAE (Kusner et al., 2017) and TES-AE (Paassen et al., 2020) use context-free grammars to constrain both their encoders and decoders so that their generated outputs are syntactically plausible. SD-VAE (Dai et al., 2018) borrows the idea of syntax-directed definition from compiler theory, and integrates it into SD-VAE’s decoder, thereby generating outputs that are semantically meaningful. GVAE, TES-AE, and SD-VAE are primarily used to generate sequential data such as molecular strings, arithmetic expressions, and programming languages. (Hu et al., 2018) adapts principles of

reinforcement learning to deep generative models so as to tune pre-specified high-level constraints, and has been used to generate human images constrained by pose templates, and to generate missing words in a sentence constrained by text templates. Similar to that body of research that uses inductive biases for text, image and molecular data, our proposed OVAE model also utilizes a strong inductive bias; but unlike the other research, OVAE incorporates an inductive bias that squares particularly well with tabular data.

3. Background

Our OVAE model is created upon the building blocks of variational autoencoders and differentiable oblivious decision trees. We describe each in turn.

3.1. Variational Autoencoder

A variational autoencoder (VAE) (Kingma & Welling, 2014) is a hierarchical Bayesian model (Gelman & Su, 2007) that postulates the existence of latent variables \mathbf{z} to help model observed variables \mathbf{x} . VAE learns a generative model $p_\theta(\mathbf{x}, \mathbf{z})$ by approximating the (intractable) posterior $p_\theta(\mathbf{z}|\mathbf{x})$ with a proposal distribution $q_\phi(\mathbf{z}|\mathbf{x})$, and by maximizing the data log-likelihood $\log p_\theta(\mathbf{x})$. The log-likelihood decomposes into two terms: the evidence lower bound (ELBO) and the KL-divergence between $q_\phi(\mathbf{z}|\mathbf{x})$ and $p_\theta(\mathbf{z}|\mathbf{x})$.

$$\log p_\theta(\mathbf{x}) = \underbrace{E_{q_\phi(\mathbf{z}|\mathbf{x})} \log \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})}}_{\text{ELBO}} + D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x}))$$

VAE uses a vanilla fully connected feed-forward neural network to parameterize the variational distribution $q_\phi(\mathbf{z}|\mathbf{x})$. It jointly finds the model’s parameters θ and the variational distribution $q_\phi(\mathbf{z}|\mathbf{x})$ that maximize the ELBO via stochastic gradient descent, using the reparameterization trick (Kingma & Welling, 2014) to calculate the gradients with respect to the variational parameters ϕ . Assuming a prior distribution $p(\mathbf{z})$ on \mathbf{z} , the generative model can be expressed as $p_\theta(\mathbf{x}, \mathbf{z}) = p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})$. In a VAE, $q_\phi(\mathbf{z}|\mathbf{x})$ and $p_\theta(\mathbf{x}|\mathbf{z})$ are respectively termed the encoder and decoder.

3.2. Differentiable Oblivious Decision Trees (DODTs)

An oblivious decision tree (ODT) (Lou & Obukhov, 2017) is a full, binary decision tree whose internal nodes at the same level are restricted to have the same splitting feature and splitting threshold. An ODT is less expressive than a regular decision tree, but its low variance makes it an ideal weak learner for a gradient boosting algorithm. Such an algorithm typically decreases bias at the cost of increasing variance (Bauer & Kohavi, 1997; Ganjisaffar et al., 2011). Because the variance of its component weak learner (ODT) is very low to begin with, the boosting algorithm’s overall

increase in variance is controlled, and is more than offset by the performance gains from decreasing bias. Empirically, boosted ODTs (Lou & Obukhov, 2017; Prokhorenkova et al., 2018) provide state-of-the-art results compared to regular boosted decision/regression trees.

An ODT of depth d is equivalent to a table with 2^d entries, each corresponding to a particular combination of d feature splits. An ODT is completely specified by its splitting features $\mathbf{f} \in \mathbb{R}^d$, splitting thresholds $\mathbf{b} \in \mathbb{R}^d$, and a d -dimensional *response* tensor \mathbf{R} (with 2^d entries) that maps the d decisions along a root-to-leaf path to the corresponding leaf value. Given an n -dimensional input $\mathbf{x} \in \mathbb{R}^n$, the output $h(\mathbf{x})$ of an ODT is

$$h(\mathbf{x}) = \mathbf{R}_{\mathbb{1}(\mathbf{f}_1(\mathbf{x})-\mathbf{b}_1), \dots, \mathbb{1}(\mathbf{f}_d(\mathbf{x})-\mathbf{b}_d)}$$

where $\mathbb{1}(\cdot)$ is the Heaviside function.

In a *differentiable* ODT (DODT) (Popov et al., 2020), the output of an ODT is made differentiable so that it can be trained end-to-end via backpropagation. The splitting features \mathbf{f} and Heaviside functions are replaced by their differentiable continuous approximations. Each splitting feature $\mathbf{f}_i(\mathbf{x})$ is now represented as a weighted sum of features $\hat{\mathbf{f}}_i(\mathbf{x})$, with the weights obtained via an α -entmax function (Peters et al., 2019) over a learnable feature selection vector $\mathbf{F}_i \in \mathbb{R}^n$, i.e.,

$$\hat{\mathbf{f}}_i(\mathbf{x}) = \sum_{j=1}^n \mathbf{x}_j \cdot \text{entmax}(\mathbf{F}_{i,j}).$$

The Heaviside function $\mathbb{1}(\mathbf{f}_i(\mathbf{x})-\mathbf{b}_i)$ is replaced with a two-class entmax $\begin{bmatrix} c_i(\mathbf{x}) \\ 1 - c_i(\mathbf{x}) \end{bmatrix} = \text{entmax}(\frac{\hat{\mathbf{f}}_i(\mathbf{x})-\mathbf{b}_i}{\tau_i}, 0]$ where τ_i is a learnable parameter to standardize the scales of the features. A *choice* tensor \mathbf{C} of the same size as the response tensor \mathbf{R} is obtained via an outer product over all c_i ’s and $(1 - c_i)$ ’s, i.e.,

$$\mathbf{C}(\mathbf{x}) = \begin{bmatrix} c_1(\mathbf{x}) \\ 1 - c_1(\mathbf{x}) \end{bmatrix} \otimes \begin{bmatrix} c_2(\mathbf{x}) \\ 1 - c_2(\mathbf{x}) \end{bmatrix} \otimes \dots \otimes \begin{bmatrix} c_d(\mathbf{x}) \\ 1 - c_d(\mathbf{x}) \end{bmatrix}$$

The scalar output $\hat{h}(\mathbf{x})$ of a DODT is computed as a sum over entries in the response vector \mathbf{R} weighted by the corresponding values in \mathbf{C} , i.e.,

$$\hat{h}(\mathbf{x}) = \sum_{i_1, \dots, i_d \in \{0,1\}^d} \mathbf{R}_{i_1, \dots, i_d} \cdot \mathbf{C}_{i_1, \dots, i_d}(\mathbf{x}).$$

The parameters of a DODT (i.e., \mathbf{F}_i , \mathbf{b}_i , τ_i , and \mathbf{R}) can be learned in an end-to-end fashion via stochastic gradient descent.

To mimic the collection of ODTs in a boosting algorithm, DODTs can be ensembled together, with the output of one feeding into the input of another. Such an ensemble outperforms regular boosted decision/regression trees and deep neural networks in an extensive set of empirical comparisons (Popov et al., 2020). Figure 1 illustrates an example DODT.

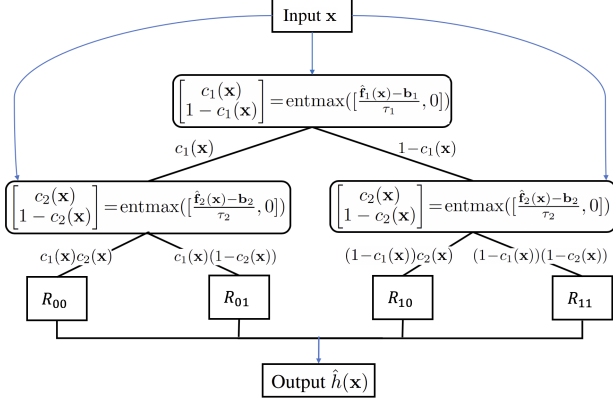


Figure 1. An example of a differentiable oblivious decision tree (DODT) of depth 2.

4. Oblivious Variational Autoencoder (OVAE)

4.1. Input/Output Representation

Our proposed OVAE model assumes that data are contained in a table containing N_c continuous columns (C_1, \dots, C_{N_c}) and N_d discrete columns (D_1, \dots, D_{N_d}), and regards each column as a random variable. In the table, each j^{th} row ($c_{1,j}, \dots, c_{N_c,j}, d_{1,j}, \dots, d_{N_d,j}$) is assumed to be a sample generated from an underlying joint distribution $\mathbb{P}(C_{1:N_c}, D_{1:N_d})$. To deal with the potential multi-modality of each continuous variable, we preprocess its continuous values using *mode-specific normalization* (MSN) (Xu et al., 2019). MSN first determines the number of modes in the distribution of each continuous variable C_i with variational Gaussian mixture models (Bishop, 2006), in which each mode m is associated with a normal distribution with mean η_m and standard deviation ψ_m . Next, for each value c_i of the continuous variable, MSN randomly samples a mode m from among the possible modes, and represents the selected mode with a one-hot encoding β_i . MSN then “normalizes” the value c_i with respect to the chosen mode’s normal distribution, i.e., $\alpha_i = \frac{c_i - \eta_m}{4\psi_m}$. Finally, each c_i is represented as the concatenation of α_i and β_i . Each row in a table is thus represented as a $(2N_c + N_d)$ -dimensional vector $\mathbf{r}_j = \alpha_{1,j} \oplus \beta_{1,j} \oplus \dots \oplus \alpha_{N_c,j} \oplus \beta_{N_c,j} \oplus \mathbf{d}_{1,j} \oplus \dots \oplus \mathbf{d}_{N_d,j}$, where \oplus is the concatenation operator and $\mathbf{d}_{i,j}$ is a one-hot encoding. Both input and generated rows share the same representation.

4.2. OVAE Encoder

We construct OVAE’s encoder by placing differentiable oblivious decision trees (DODTs) in parallel in a layer, and then stacking such DODT layers one on top of another. The outputs of the parallel trees in a layer are concatenated before being fed as input into another layer (see example in Figure 2). We include several parallel DODTs in a layer so

that each DODT can capture a different way of partitioning its input data. This is particularly useful for rich datasets that can be partitioned in multiple valid ways because it allows the DODTs to fully capture the data’s complexity. We stack one DODT layer upon another so that the latter can model the intricate dependencies among the different data partitionings in the former. $\text{DODTLayer}_{n \rightarrow k}(\mathbf{x})$ denotes a DODT layer that consists of k parallel DODTs, and it maps an n -dimensional input $\mathbf{x} \in \mathbb{R}^n$ to a k -dimensional output. The input \mathbf{x} is fed into each of the k DODTs, and each DODT outputs a scalar value.

The architecture for the encoder distribution $q_\phi(\mathbf{z}_j | \mathbf{r}_j)$ is as follows (\mathbf{r}_j represents a row in a table with N_c continuous columns and N_d discrete columns as described in Section 4.1).

$$\begin{aligned} \mathbf{h}_1 &= \text{DODTLayer}_{|\mathbf{r}_j| \rightarrow k}(\mathbf{r}_j) \\ \mathbf{h}_2 &= \text{DODTLayer}_{k \rightarrow k}(\mathbf{h}_1) \\ \boldsymbol{\mu} &= \text{DODTLayer}_{k \rightarrow k}(\mathbf{h}_2) \\ \boldsymbol{\sigma} &= \exp(\text{DODTLayer}_{k \rightarrow k}(\mathbf{h}_2)) \\ q_\phi(\mathbf{z}_j | \mathbf{r}_j) &\sim \mathcal{N}(\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma})) \end{aligned}$$

4.3. OVAE Decoder

OVAE’s decoder has to model the $\alpha_{i,j}$ and $\beta_{i,j}$ values for continuous columns, and the $\mathbf{d}_{i,j}$ values for discrete columns for each j^{th} row (the symbols are described in Section 4.1). OVAE assumes that each $\alpha_{i,j}$ has a normal distribution (with a column-specific variance δ_i), and that $\beta_{i,j}$ and $\mathbf{d}_{i,j}$ each has a categorical probability mass function. The architecture for the decoder distribution $p_\theta(\mathbf{r}_j | \mathbf{z}_j)$ is as follows.

$$\begin{aligned} \mathbf{h}_1 &= \text{DODTLayer}_{k \rightarrow k}(\mathbf{z}_j) \\ \mathbf{h}_2 &= \text{DODTLayer}_{k \rightarrow k}(\mathbf{h}_1) \\ \bar{\alpha}_{i,j} &= \tanh(\text{DODTLayer}_{k \rightarrow 1}(\mathbf{h}_2)) && \text{for } 1 \leq i \leq N_c \\ \hat{\alpha}_{i,j} &\sim \mathcal{N}(\bar{\alpha}_{i,j}, \delta_i) && \text{for } 1 \leq i \leq N_c \\ \hat{\beta}_{i,j} &\sim \text{softmax}(\text{DODTLayer}_{k \rightarrow m_i}(\mathbf{h}_2)) && \text{for } 1 \leq i \leq N_c \\ \hat{\mathbf{d}}_{i,j} &\sim \text{softmax}(\text{DODTLayer}_{k \rightarrow |D_i|}(\mathbf{h}_2)) && \text{for } 1 \leq i \leq N_d \\ p_\theta(\mathbf{r}_j | \mathbf{z}_j) &= \prod_{i=1}^{N_c} \mathbb{P}(\hat{\alpha}_{i,j} = \alpha_{i,j}) \prod_{i=1}^{N_c} \mathbb{P}(\hat{\beta}_{i,j} = \beta_{i,j}) \\ &\quad \prod_{i=1}^{N_d} \mathbb{P}(\hat{\mathbf{d}}_{i,j} = \mathbf{d}_{i,j}) \end{aligned}$$

m_i is the number of modes in continuous column i , and $|D_i|$ is the number of distinct discrete values in discrete column D_i . To obtain one-hot encodings from the $\hat{\beta}_{i,j}$ and $\hat{\mathbf{d}}_{i,j}$ vectors, we simply set the maximum value in each vector to 1, and all other values to 0. The parameters of OVAE are the δ_i ’s and its constituent DODTs’ parameters. These parameters are learned via stochastic gradient descent by maximizing the evidence lower bound (Section 3.1).

Table 1. Results on classification datasets.

| | adult | census | credit | cover. | intru. | mnist12/28 | | Average |
|-----------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | F1 | F1 | F1 | Macro-F1 | Macro-F1 | Micro-F1 | Micro-F1 | F1 |
| Identity | 0.67 | 0.49 | 0.72 | 0.65 | 0.86 | 0.89 | 0.92 | 0.74 |
| CLBN | 0.33 | 0.31 | 0.41 | 0.32 | 0.38 | 0.74 | 0.18 | 0.38 |
| PrivBayes | 0.41 | 0.12 | 0.19 | 0.27 | 0.38 | 0.12 | 0.08 | 0.23 |
| medGAN | 0.38 | 0.00 | 0.00 | 0.09 | 0.30 | 0.09 | 0.10 | 0.14 |
| VEEGAN | 0.24 | 0.09 | 0.00 | 0.08 | 0.26 | 0.19 | 0.14 | 0.14 |
| tableGAN | 0.49 | 0.36 | 0.18 | 0.00 | 0.00 | 0.10 | 0.00 | 0.16 |
| CTGAN | <u>0.60</u> | 0.39 | 0.67 | 0.32 | 0.53 | 0.39 | 0.37 | 0.47 |
| TVAE | 0.63 | 0.38 | 0.10 | 0.43 | <u>0.51</u> | <u>0.79</u> | <u>0.79</u> | <u>0.52</u> |
| OVAE | <u>0.60</u> | <u>0.38</u> | <u>0.51</u> | 0.45 | 0.53 | 0.83 | 0.84 | 0.59 |

5. Experiments

5.1. Datasets

For our experiments, we use 5 real-world regression datasets, 7 real-world classification datasets (12 tabular datasets in total).

All 5 real-world regression datasets are from the UCI machine learning repository (Dua & Graff, 2017) (bike-sharing (bike), GPU kernel performance (gpu), wine quality (wine), power plant (power), and online news popularity (news)). Among the 7 real-world classification datasets, 4 are from the UCI machine learning repository (adult, census, covertype, and intrusion), and 1 is from Kaggle (credit). The remaining 2 classification datasets mnist28 and mnist12 are respectively obtained by binarizing 28×28 and 12×12 MNIST images (LeCun & Cortes, 2010) into feature vectors (with an additional label column indicating the target digit). These two datasets allow us to investigate the performances of OVAE and its comparison systems on high dimensional binary tabular data. All datasets are tabular in nature, and each is divided into a training set \mathbf{T}_{train} and a test set \mathbf{T}_{test} .

5.2. Methodology

We compare our OVAE model against 7 other models (see Section 2). Two of these models are based on Bayesian networks: CLBN (Chow & Liu, 1968) and PrivBayes (Zhang et al., 2017). The remaining models are state-of-the-art deep learning ones: medGAN (Choi et al., 2017), VEEGAN (Srivastava et al., 2017), tableGAN (Park et al., 2018), CTGAN (Xu et al., 2019), and TVAE (Xu et al., 2019).

We follow the experimental methodology adopted by (Xu et al., 2019). For every real-world dataset, we train each model on the training tabular data \mathbf{T}_{train} , and use the trained model to generate synthetic tabular data \mathbf{T}_{syn} . We then train a set of standard regressors or classifiers (e.g., (boosted) regression/decision tree, linear regression, and multilayer perceptron) on \mathbf{T}_{syn} , and evaluate the set of

regressors/classifiers on the test tabular data \mathbf{T}_{test} . We run this process thrice for each model per dataset, and report the average result of the set of regressors/classifiers. For regression tasks, we report the average R^2 score of a set of regressors on \mathbf{T}_{test} . R^2 reflects the proportion of the variance in a dependent variable that is successfully modeled. R^2 is (negatively) related to mean-squared error (MSE) such that a higher R^2 leads to a correspondingly lower MSE (and vice versa). Thus the higher the R^2 score, the better the performance of a model. For classification tasks, we report variants of the F1 score depending on the skew of class labels in a dataset. adult, credit and census are binary-class datasets, and their class-label distributions are skewed towards one of two classes. To measure a model’s performance on the (more) difficult task of predicting the minority class, we consider the minority class as the positive class, and report the standard F1 score for these three datasets. covertype and intrusion are multi-class datasets, and they also exhibit highly imbalanced class labels. Thus we report the macro-F1 score for them. mnist12 and mnist28 are multi-class datasets with class labels that are well-balanced. Hence, we use micro-F1 (accuracy) as a suitable metric for these two datasets. The higher the F1 score (or its variants), the better the performance of a model.

We also have an *Identity* system that simply copies \mathbf{T}_{train} ,

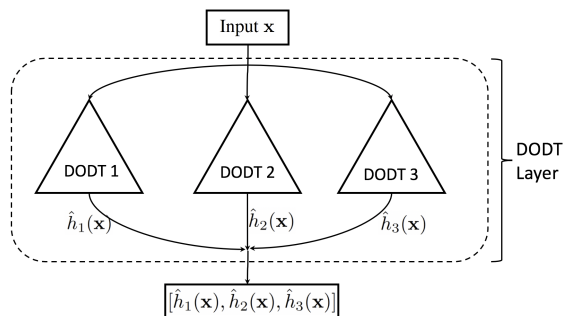


Figure 2. A DODT layer with 3 parallel DODTs (each triangle represents a DODT). The outputs of the DODTs are concatenated.

Table 2. Results on regression datasets.

| | bike R^2 | gpu R^2 | wine R^2 | power R^2 | news R^2 | Average R^2 |
|-----------|----------------------|-------------------|---------------|-------------------|-------------------|----------------------|
| Identity | 0.86 | 0.75 | 0.28 | 0.88 | 0.14 | 0.58 |
| CLBN | -49×10^1 | -18 | -27 | -22×10^3 | -6.3 | -45×10^2 |
| PrivBayes | 0.71 | -29×10^2 | -10 | 0.88 | -4.5 | -58×10 |
| medGAN | -50×10^{16} | -54×10^2 | -42 | -11.82 | -8.80 | -10×10^{16} |
| VEEGAN | -0.43 | -16×10^2 | -11 | -0.67 | -65×10^5 | -13×10^5 |
| tableGAN | 0.74 | -82×10^2 | -0.10 | 0.77 | -3.09 | -16×10^2 |
| CTGAN | <u>0.78</u> | 0.60 | <u>0.02</u> | -0.19 | -0.43 | 0.16 |
| TVAE | 0.70 | <u>0.68</u> | 0.22 | 0.72 | -0.20 | <u>0.42</u> |
| OVAE | 0.82 | 0.70 | 0.22 | <u>0.81</u> | <u>-0.30</u> | 0.45 |

and treats it as \mathbf{T}_{syn} in the aforementioned methodology (rather than doing the hard work of learning a model to generate \mathbf{T}_{syn}). The R^2 , F1, and L_{test} scores associated with the Identity system serve as upper bounds for the scores of OVAE and its comparison models. OVAE is trained with stochastic gradient descent using quasi-hyperbolic ADAM (Ma & Yarats, 2018) as the optimizer (we use the recommended parameters in its paper). In each DODT layer in the OVAE model, we use either 128, 256 or 512 parallel differentiable oblivious decision trees, (i.e., $k \in \{128, 256, 512\}$ in Section 4.2 and 4.3). The depth of each DODT is set to 6. The α in α -entmax activation function used in Differentiable Oblivious Decision Tree (DODT) is 1.5. The parameters are chosen using preliminary experiments, and depend on whether the resultant models can fit into our GPU’s memory (Nvidia Geforce RTX 2080 Ti; 11GB).

The regression results are shown in Table 2 (best results are boldfaced; second best results are underlined). The numbers in the *news* are as reported in (Xu et al., 2019) (modulo the OVAE results). All other numbers are from our experiments. Our OVAE model outperforms TVAE on 3 real-world regression datasets (*bike*, *gpu*, and *power*), ties on one, and loses on another. OVAE differs from TVAE primarily in using layers of differentiable oblivious decision trees (DODTs) in place of the standard feedforward neural networks in TVAE’s encoders and decoders. The results give credence to our hypothesis that DODTs provide a useful inductive bias for improving tabular data generation. Note that where OVAE is not the best model (*power* and *news*), it is second best. On average, OVAE is the best performer on the real-world regression datasets. (We were unable to replicate CTGAN’s and TVAE’s results on *news* as given in (Xu et al., 2019) using its provided code, but decided to still report those numbers in Table 2 to cast the comparison systems in the best light. From our experiments, CTGAN’s and TVAE’s R^2 scores on *news* are respectively -0.07 and -0.64, and their resulting average R^2 scores are 0.23 and 0.34. The relative ordering of the systems remains unchanged, with OVAE being the best, but its R^2 gap from

TVAE is larger: 0.11 as opposed to the current 0.03.)

The classification results are shown in Table 1. The numbers in the OVAE row are obtained from our experiments; all other numbers are as reported in (Xu et al., 2019). OVAE outperforms TVAE on 5 real-world classification datasets (*credit*, *covertypes*, *intrusion*, *mnist12*, and *mnist28*), ties on one, and loses on another; on average, OVAE outperforms TVAE on the real-world classification datasets. Again, this shows that using DODTs in OVAE leads to better results vis-à-vis TVAE. Like on the real-world regression datasets, OVAE is consistently the best model (on *covertypes*, *intrusion*, *mnist12*, and *mnist28*) or the second best performer (*adult*, *census*, and *credit*). In aggregate, OVAE is the best performing system.

6. Conclusion and Future Work

We presented OVAE, a new model for generating synthetic tabular data. OVAE combines differentiable oblivious decision trees (DODTs) with variational autoencoders (VAEs), thereby incorporating a strong inductive bias for tabular data into VAEs. To fully capture the richness in tabular data, OVAE gathers several parallel DODTs into a layer, and stacks such layers one upon another. Empirical comparisons with seven systems on 12 real-world datasets show the promise of our approach. Our proposed OVAE model advances the line of research that obviates privacy restrictions on sensitive tabular data by generating high-fidelity synthetic data. Our OVAE model could be used in a variety of fields (e.g., finance and healthcare) to expedite the development of AI systems for social good (e.g., credit scoring for micro-finance, and epidemiological forecasting) preserving privacy. Our work is particularly pertinent in the current Covid-19 pandemic, during which a constant stream of Covid-19 patient data are collected by medical institutions, but are only accessible to a select group of researchers who are associated with the institutions. By using OVAE to generate synthetic data, we could preserve privacy of the patients and make high-fidelity fake patient data publicly available, allowing a wider span of intellectual

resources and human ingenuity to be brought to bear on Covid-19 problems. As future work, we want to incorporate domain knowledge as additional tabular constraints into OVAE, utilize normalizing flows to improve OVAE.

7. Acknowledgements

This research is partly funded by an MOE AcRF Tier 1 grant (R -253-000-146-133) and an MOH NIC grant (MOH/NIC/CDM1/2018) to Stanley Kok.

References

- Aviñó, L., Ruffini, M., and Gavaldà, R. Generating synthetic but plausible healthcare record datasets. *arXiv preprint arXiv:1807.01514*, 2018.
- Bauer, E. and Kohavi, R. An empirical comparison of voting classification algorithms: Bagging, boosting, and variants. *Machine Learning*, 36:105–139, 1997.
- Bishop, C. M. *Pattern Recognition and Machine Learning*. Springer, 2006.
- Camino, R., Hammerschmidt, C., and State, R. Generating multi-categorical samples with generative adversarial networks. In *ICML Workshop on Theoretical Foundations and Applications of Deep Generative Models*, 2018.
- Che, Z., Cheng, Y., Zhai, S., Sun, Z., and Liu, Y. Boosting deep learning risk prediction with generative adversarial networks for electronic health records. *2017 IEEE International Conference on Data Mining (ICDM)*, pp. 787–792, 2017.
- Choi, E., Biswal, S., Malin, B., Duke, J., Stewart, W. F., and Sun, J. Generating multi-label discrete patient records using generative adversarial networks. In Doshi-Velez, F., Fackler, J., Kale, D., Ranganath, R., Wallace, B., and Wiens, J. (eds.), *Proceedings of the 2nd Machine Learning for Healthcare Conference*, volume 68 of *Proceedings of Machine Learning Research*, pp. 286–305, Boston, Massachusetts, 18–19 Aug 2017. PMLR. URL <http://proceedings.mlr.press/v68/choi17a.html>.
- Chow, C. and Liu, C. Approximating discrete probability distributions with dependence trees. *IEEE Transactions on Information Theory*, 14(3):462–467, 1968.
- Dai, H., Tian, Y., Dai, B., Skiena, S., and Song, L. Syntax-directed variational autoencoder for structured data. In *International Conference on Learning Representations*, 2018.
- Dua, D. and Graff, C. UCI machine learning repository, 2017.
- Emam, K. E., Jonker, E., Arbuckle, L., and Malin, B. A systematic review of re-identification attacks on health data. *PLoS ONE*, 6, 2011.
- Ganjisaffar, Y., Caruana, R., and Lopes, C. V. Bagging gradient-boosted trees for high precision, low variance ranking models. In *SIGIR*, 2011.
- Gelman, A. and Su, Y.-S. *Data Analysis Using Regression and Multilevel/Hierarchical Models*. Cambridge University Press, 2007.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial nets. In Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N. D., and Weinberger, K. Q. (eds.), *Advances in Neural Information Processing Systems 27*, pp. 2672–2680. Curran Associates, Inc., 2014. URL <http://papers.nips.cc/paper/5423-generative-adversarial-nets.pdf>.
- Guo, J., Lu, S., Cai, H., Zhang, W., Yu, Y., and Wang, J. Long text generation via adversarial training with leaked information. *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018.
- Hu, Z., Yang, Z., Salakhutdinov, R., Liang, X., Qin, L., Dong, H., and Xing, E. Deep generative models with learnable knowledge constraints. *NeurIPS*, 2018.
- Kingma, D. P. and Welling, M. Auto-encoding variational Bayes. In *International Conference on Learning Representations*, 2014.
- Kusner, M. J., Paige, B., and Hernández-Lobato, J. M. Grammar variational autoencoder. In *International Conference on Machine Learning*, 2017.
- LeCun, Y. and Cortes, C. The MNIST database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 2010.
- Lin, C. H., Chang, C.-C., Chen, Y.-S., Juan, D.-C., Wei, W., and Chen, H.-T. COCO-GAN: Generation by parts via conditional coordinating. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 4511–4520, 2019.
- Lou, Y. and Obukhov, M. Bdt: Gradient boosted decision tables for high accuracy and scoring efficiency. *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2017.
- Ma, J. and Yarats, D. Quasi-hyperbolic momentum and ADAM for deep learning. *arXiv preprint arXiv:1810.06801*, 2018.
- Paassen, B., Koprinska, I., and Yacef, K. Tree echo state autoencoders with grammars. In *International Joint Conference on Neural Networks*, 2020.

- Papernot, N., Abadi, M., Erlingsson, Ú., Goodfellow, I. J., and Talwar, K. PATE-GAN: Generating synthetic data with differential privacy guarantees. In *International Conference on Learning Representations*, 2017.
- Park, N., Mohammadi, M., Gorde, K., Jajodia, S., Park, H., and Kim, Y. Data synthesis based on generative adversarial networks. *ArXiv*, abs/1806.03384, 2018.
- Patki, N., Wedge, R., and Veeramachaneni, K. The synthetic data vault. In *2016 IEEE International Conference on Data Science and Advanced Analytics (DSAA)*, pp. 399–410, 2016.
- Peters, B., Niculae, V., and Martins, A. F. T. Sparse sequence-to-sequence models. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 1504–1519, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1146.
- Popov, S., Morozov, S., and Babenko, A. Neural oblivious decision ensembles for deep learning on tabular data. In *International Conference on Learning Representations*, 2020.
- Prokhorenkova, L., Gusev, G., Vorobev, A., Dorogush, A. V., and Gulin, A. Catboost: unbiased boosting with categorical features. In *Advances in neural information processing systems*, pp. 6638–6648, 2018.
- Srivastava, A., Valkov, L., Russell, C., Gutmann, M. U., and Sutton, C. VEEGAN: Reducing mode collapse in GANs using implicit variational learning. In *Advances in Neural Information Processing Systems*, pp. 3308–3318, 2017.
- Sun, Y., Cuesta-Infante, A., and Veeramachaneni, K. Learning vine copula models for synthetic data generation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33:5049–5057, 07 2019. doi: 10.1609/aaai.v33i01.33015049.
- Xu, L., Skoularidou, M., Cuesta-Infante, A., and Veeramachaneni, K. Modeling tabular data using conditional GAN. In Wallach, H., Larochelle, H., Beygelzimer, A., d’Alché Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 32*, pp. 7335–7345. Curran Associates, Inc., 2019.
- Yahi, A., Vanguri, R., Elhadad, N., and Tatonetti, N. P. Generative adversarial networks for electronic health records: a framework for exploring and evaluating methods for predicting drug-induced laboratory test trajectories. *arXiv preprint arXiv:1712.00164*, 2017.
- Yoon, J., Jordon, J., and van der Schaar, M. PATE-GAN: Generating synthetic data with differential privacy guarantees. In *International Conference on Learning Representations*, 2019.
- Zhang, J., Cormode, G., Procopiuc, C. M., Srivastava, D., and Xiao, X. PrivBayes: Private data release via bayesian networks. *ACM Trans. Database Syst.*, 42(4), October 2017. ISSN 0362-5915. doi: 10.1145/3134428. URL <https://doi.org/10.1145/3134428>.