

Sythetic Data Generation

School of Data science

The Chinese University of Hong Kong, Shenzhen

August 26, 2022





- 1 Background of tabular data generation
- 2 Overview of generative model for tabular data



Part 1. Background of tabular data generation



- Heterogeneous Property: mixed type data.
 - Different from image and language data, tabular data has dense numerical and sparse categorical features.
e.g., categorical, ordinal, continuous

- Ubiquitous in many crucial applications:
 - medical diagnosis based on patient history
 - predictive analytics for financial applications,
e.g., risk analysis, estimation of creditworthiness, the recommendation of investment strategies, and portfolio management



- Low-quality training data,
 - e.g., missing values, class-imbalanced
- Complex or irregular dependencies between different columns,
 - e.g., a change of a categorical feature can entirely flip a prediction on tabular data
 - Many features are uninformative
- Handling the categorical features remains particularly challenging

Therefore, for classification and regression problems with tabular data, using tree ensemble models can outperform deep learning methods.¹²

¹Grinsztajn, et al. Why do tree-based models still outperform deep learning on tabular data? NIPS 2022 workshop track.

²Shwartz et al. Tabular Data: Deep Learning is Not All You Need. ICML 2021 workshop track.



- Utility: Incorporate more training data to enhance the performance
- Privacy: Sensitive data from users,
 - e.g., Information Leakage of medical diagnosis, Membership Inference Attacks
- Control generation,
 - e.g., class conditional generation, imputation



Part 2. Overview of generative model for tabular data

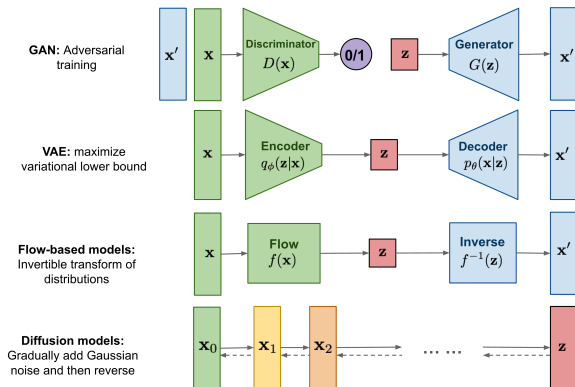


Method	Based upon	Application
medGAN [46]	Autoencoder+GAN	Medical Records
TableGAN [145] DCGAN	General	
Mottini et al. [149]	Cramér GAN	Passenger Records
Camino et al. [150]	medGAN, ARAE	General
medBGAN, medWGAN [151]	WGAN-GP, Boundary seeking GAN	Medical Records
ITS-GAN [124]	GAN with AE for constraints	General
CTGAN, TVAE [144]	Wasserstein GAN, VAE	General
actGAN [126]	WGAN-GP	Health Data
VAEM [143]	VAE (Hierarchical)	General
OVAE [152]	Oblivious VAE	General
TAEI [44]	AE+SMOTE (in multiple setups)	General
Causal-TGAN [153]	Causal-Model, WGAN-GP	General
Copula-Flow [45]	Invertible Flows	General

TABLE III: Generation of tabular data using deep neural network models (in chronological order).

³Borisov, et al. (2021). Deep Neural Networks and Tabular Data: A Survey. ArXiv, abs/2110.01889.

Overview of generative model:⁴



GFlowNet: A sampling method for discrete type data training by reinforcement criterion

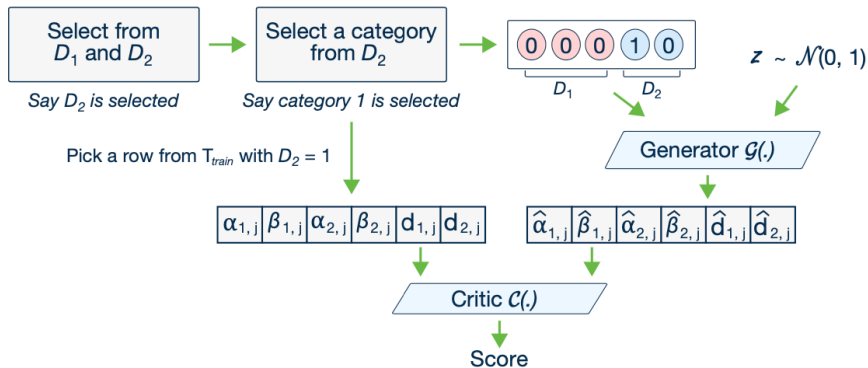
⁴ <https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>

#connection-with-noise-conditioned-score-networks-ncsn



Origin GAN: $\min_G \max_D \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} [\log(D(\mathbf{x}))] + \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g} [\log(1 - D(\tilde{\mathbf{x}}))]$

WGAN: $\min_G \max_{D \in \mathcal{D}} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} [D(\mathbf{x})] - \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g} [D(\tilde{\mathbf{x}})]$



⁵Xu, Lei et al. Modeling Tabular data using Conditional GAN. NeurIPS (2019).



Table 3: Ablation study results on mode-specific normalization, conditional generator and training-by-sampling module, as well as the network architecture. The absolute performance change on real classification datasets (excluding MNIST) is reported.

Model	Mode-specific Normalization			Generator		Network Architecture		
	GMM5	GMM10	MinMax	w/o S.	w/o C.	GAN	WGANGP	GAN+PacGAN
Performance	-4.1%	-8.6%	-25.7%	-17.8%	-36.5%	-6.5%	+1.75%	-5.2%

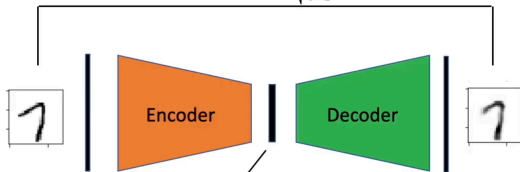
Rule of thumb:

- Conditional generation
- mode-specific normalization
- WGAN+gradient penalty



1) Minimize squared error loss: (ensures good reconstruction)

$$\mathcal{L}_1 = \|\mathbf{x} - \text{Dec}(\text{Enc}(\mathbf{x}))\|_2^2 = \sqrt{\sum_{i=1}^d (x_i - x'_i)^2}$$



2) Minimize KL divergence:

$$\mathcal{L}_2 = D_{KL} [N(\mu, \sigma) \| N(0, 1)] = -\frac{1}{2} \sum (1 + \log(\sigma^2) - \mu^2 - \sigma^2)$$

⁶ Comments:

- The generator in GANs does not have access to real data during the entire training process; thus, we can make CTGAN achieve differential privacy easier than TVAE.

⁶ Image credits to Sebastian Raschka

⁷ Xu, Lei et al. Modeling Tabular data using Conditional GAN. NeurIPS 2019.



VAE can generate discrete type data

- For categorical columns, we can use the softmax over all categories.
- For continuous-discrete columns (like salary), we can model it as a continuous variable and discretize it in the end
- For ordinal-discrete columns (like ratings), we can use ordinal regression likelihood⁸.

⁸Paquet, et al. A hierarchical model for ordinal matrix factorization. *Statistics and Computing*, 22, 945-957.



Basic idea: VAEM uses a hierarchy of latent variables, which fits in two stages.

- In the first stage, learn one type-specific VAE for each dimension. These initial one-dimensional VAEs capture marginal distribution properties and provide a latent representation that is more homogeneous across dimensions.
- In the second stage, another VAE is used to capture dependencies among the one-dimensional latent representations from the first stage.

⁹Ma, et al. VAEM: a Deep Generative Model for Heterogeneous Mixed Type Data. NIPS 2020.



Modeling:

$$\mathbf{z} \sim \pi(\mathbf{z}), \mathbf{x} = f_{\theta}(\mathbf{z}), \mathbf{z} = f_{\theta}^{-1}(\mathbf{x})$$
$$p_{\theta}(\mathbf{x}) = \pi(\mathbf{z}) \left| \det \frac{d\mathbf{z}}{d\mathbf{x}} \right| = \pi(f_{\theta}^{-1}(\mathbf{x})) \left| \det \frac{df_{\theta}^{-1}}{d\mathbf{x}} \right|$$

Training objective:

$$\min_{\theta} \mathcal{L}(\theta|\mathcal{D}) = -\min_{\theta} \frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \log p_{\theta}(\mathbf{x})$$

Key problem:

- Normalizing Flow cannot model discrete type data

Solution:

- Dequantization, i.e., adding real-valued noise to the discrete data.
 - Uniform dequantization
 - variational dequantization

¹⁰ Lee, et al. Differentially Private Normalizing Flows for Synthetic Tabular Data Generation. AAAI 2022.



Modeling:

$$p(\mathbf{x}) = \prod_{i=1}^D p(x_i | x_1, \dots, x_{i-1}) = \prod_{i=1}^D p(x_i | x_{1:i-1})$$

Key problem:

- Tabular data is not sequential data like images or language

Solution:

- Using Masked Autoencoder, e.g., MADE (Masked Autoencoder for Distribution Estimation)

Comments:

- No work using Autoregressive model for tabular data
- A library¹¹ using traditional machine learning methods is available.
- self-supervision loss is attractive/somehow promising, but with high limitations in tabular data

¹¹Mahiou, et al. dpart: Differentially Private Autoregressive Tabular, a General Framework for Synthetic Data Generation. ICML 2022 workshop track.

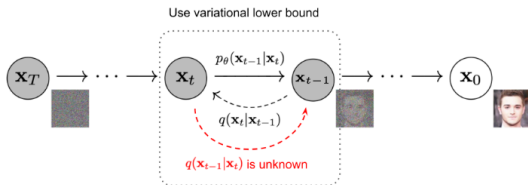


Fig. 2. The Markov chain of forward (reverse) diffusion process of generating a sample by slowly adding (removing) noise. (Image source: [Ho et al., 2020](#) with a few additional annotations)

Training objective:

$$\min_{\theta} L_{\text{VLB}} = \min_{\theta} \mathbb{E}_q \left[\underbrace{D_{\text{KL}}(q(\mathbf{x}_T | \mathbf{x}_0) \| p(\mathbf{x}_T))}_{L_T} + \sum_{t>1} \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t))}_{L_{t-1}} - \log p_{\theta}(\mathbf{x}_0 | \mathbf{x}_1) \right]$$



- Diffusion model seems to be a promising research direction on tabular data generation
- Some key points to be solved:
 - How to deal with discrete columns?
 - Discrete diffusion ¹²
 - How to deal with the relationship between discrete columns and continuous columns?
 - Conditional diffusion ¹³

¹²Austin, et al. Structured Denoising Diffusion Models in Discrete State-Spaces. NIPS 2021.

¹³Batzolis, et al. Conditional Image Generation with Score-Based Diffusion Models.



GFlowNet: A sampling method for discrete type data training by reinforcement objective.

Flow consistency equations:

$$\sum_{s,a:T(s,a)=s'} F(s,a) = R(s') + \sum_{a' \in \mathcal{A}(s')} F(s',a').$$

Training objective:

$$\mathcal{L}_{\theta,\epsilon}(\tau) = \sum_{s' \in \tau \neq s_0} \left(\log \left[\epsilon + \sum_{s,a:T(s,a)=s'} \exp F_{\theta}^{\log}(s,a) \right] - \log \left[\epsilon + R(s') + \sum_{a' \in \mathcal{A}(s')} \exp F_{\theta}^{\log}(s',a') \right] \right)^2$$

¹⁴Bengio, et al. Flow Network based Generative Models for Non-Iterative Diverse Candidate Generation. NIPS 2021.



- GFlowNet explicitly models the relationship between discrete columns
- Some key points to be solved:
 - GFlowNet can only deal with discrete columns, we need to deal with the relationship between discrete columns and continuous columns.
 - Conditional diffusion/ Conditional GAN



- We get familiar with the properties of tabular data and the difficulties for modelling tabular data
- We get familiar with all types of generative model, especially for synthesizing tabular data
- some open questions