

# Does flow-matching generated synthetic data preserve feature correlations better than VAE and GAN baselines wh

Assignee Research

June 10, 2026

## Abstract

The rising use of machine learning in various fields requires robust methods to create synthetic tabular data. Data should preserve key characteristics while addressing data scarcity challenges. Current approaches based on Generative Adversarial Networks, such as the state-of-the-art CTGAN model, struggle with the complex structures inherent in tabular data. These data often contain both continuous and discrete features with non-Gaussian distributions. Therefore, we propose a novel Variational Autoencoder (VAE)-based model that addresses these limitations. Inspired by the TVAE model, our appro

## 1 Introduction

This paper examines: An improved tabular data generator with VAE-GMM integration. Research question: Does flow-matching generated synthetic data preserve feature correlations better than VAE and GAN baselines when evaluated on high-dimensional imbalanced tabular benchmarks?.

## 2 Methodology

Systematic literature search across multiple databases yielded 10 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.6/10.

## 3 Results

10 papers retrieved. 16 claims extracted; 0 independently verified. Quality review score: 4.6/10.

## 4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

## 5 Extracted Claims

Claim	Verified	Confidence
The accuracy of CTGAN on the Adult dataset is 0.80 with a confidence interval of (0.79, 0.82).	×	0.02
The accuracy of TVAE on the Adult dataset is 0.79 with a confidence interval of (0.77, 0.80).	×	0.01
The accuracy of the proposed approach on the Adult dataset is 0.80 with a confidence interval of (0.77, 0.82).	×	0.03
The C-index of CTGAN on the Metabric dataset is 0.58 with a confidence interval of (0.53, 0.63).	×	0.02
The C-index of TVAE on the Metabric dataset is 0.57 with a confidence interval of (0.52, 0.62).	×	0.01
The C-index of the proposed approach on the Metabric dataset is 0.60 with a confidence interval of (0.54, 0.65).	×	0.03
The C-index of CTGAN on the STD dataset is 0.64 with a confidence interval of (0.57, 0.70).	×	0.02
The C-index of TVAE on the STD dataset is 0.54 with a confidence interval of (0.47, 0.61).	×	0.01
The C-index of the proposed approach on the STD dataset is 0.54 with a confidence interval of (0.46, 0.60).	×	0.03
The performance metrics obtained using data generated by each model (TVAE and the proposed approach) are comparable to t	×	0.09
The VAE was first introduced by [11] as a probabilistic generative model to perform Bayesian inference on datasets consi	×	0.05
The core principle of VAE lies in learning a probabilistic generative model to capture the underlying latent structure w	×	0.07
VAEs employ variational methods to estimate the true posterior density.	×	0.02
The optimization of VAE parameters typically involves maximizing the Evidence Lower Bound (ELBO).	×	0.03
The ELBO provides a lower bound on the data’s marginal log-likelihood.	×	0.03
The reparameterization trick is used to address the challenge of computing the gradient of the ELBO with respect to .	×	0.02

## References

- <http://arxiv.org/abs/2404.08434v2>
- <http://arxiv.org/abs/2512.21798v2>
- <http://arxiv.org/abs/2502.17119v2>