

How does the F1-score of diffusion-based tabular generative models compare to CTGAN when augmenting data for t

Assignee Research

June 10, 2026

Abstract

Class imbalance in tabular datasets poses a challenge for machine learning classification tasks, often leading to biased models that underperform in predicting minority class instances. This study presents a comparative analysis of synthetic data generation methods for addressing class imbalance in tabular data. We evaluate four augmentation approaches—Synthetic Minority Over-sampling Technique (SMOTE), Gaussian Copula, Tabular Variational Autoencoder (TVAE), and Conditional Tabular Generative Adversarial Network (CTGAN)—using the University of California Irvine (UCI) Bank Marketing dataset, w

1 Introduction

This paper examines: Synthetic Data Augmentation for Imbalanced Tabular Data: A Comparative Study of Generation Methods. Research question: How does the F1-score of diffusion-based tabular generative models compare to CTGAN when augmenting data for training LLMs on imbalanced text classification benchmarks using the HAN benchmark?.

2 Methodology

Systematic literature search across multiple databases yielded 11 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.5/10.

3 Results

11 papers retrieved. 10 claims extracted; 8 independently verified. Quality review score: 7.5/10.

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 study evaluates four synthetic data augmentation approaches: SMOTE, Gaussian Copula, TVAE, and CTGAN.	✓	0.20
The evaluation uses the University of California Irvine (UCI) Bank Marketing dataset.	✓	0.16
The UCI Bank Marketing dataset exhibits a class imbalance ratio of approximately 7.88:1.	✓	0.26
Statistical fidelity was evaluated using four metrics: marginal numerical similarity, categorical distribution similarity	✓	0.27
All augmentation methods achieved statistically significant improvements over the baseline with a p-value less than 0.05	✓	0.20
SMOTE achieved the highest recall for minority class detection at 54.2%.	✓	0.24
SMOTE's recall represents a 117.6% relative improvement over the baseline.	×	0.14
SMOTE achieved the highest F1-Score for minority class detection at 0.437.	✓	0.24
SMOTE's F1-Score represents a 22.4% improvement over the baseline.	×	0.12
Gaussian Copula provided the highest composite fidelity score of 0.930.	✓	0.26

References

- <https://doi.org/10.3390/s25165144>
- <https://doi.org/10.3390/electronics15040883>
- <https://doi.org/10.3390/electronics13173509>