

Can novel evaluation metrics for generative tabular models predict the robustness of aligned multimodal system

Assignee Research

June 10, 2026

Abstract

Generative models have revolutionized multiple domains, yet their application to tabular data remains underexplored. Evaluating generative models for tabular data presents unique challenges due to structural complexity, large-scale variability, and mixed data types, making it difficult to intuitively capture intricate patterns. Existing evaluation metrics offer only partial insights, lacking a comprehensive measure of generative performance. To address this limitation, we propose three novel evaluation metrics: FAED, FPCAD, and RFIS. Our extensive experimental analysis, conducted on three stan

1 Introduction

This paper examines: Evaluating Generative Models for Tabular Data: Novel Metrics and Benchmarking. Research question: Can novel evaluation metrics for generative tabular models predict the robustness of aligned multimodal systems against distribution shifts in numerical and text features?.

2 Methodology

Systematic literature search across multiple databases yielded 7 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 6.5/10.

3 Results

7 papers retrieved. 6 claims extracted; 3 independently verified. Quality review score: 6.5/10.

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
FAED effectively captures generative modeling issues overlooked by existing metrics.	✓	0.30
FPCAD exhibits promising performance but requires further refinements to enhance its reliability.	✓	0.18
FAED successfully detects all synthesized problems (Quality Decrease, Mode Drop, and Mode Collapse) in the experimental	×	0.11
Existing metrics (SDV Fidelity, Utility, TSTR, and TRTS) fail to identify key issues in generative modeling for tabular	✓	0.20
TSTR (Train on Synthetic, Test on Real) is useful for detecting cases where synthetic data only partially represents rea	×	0.05
TRTS (Train on Real, Test on Synthetic) assesses whether synthetic samples introduce patterns absent in real data.	×	0.03

References

- <http://arxiv.org/abs/2504.20900v1>
- <http://arxiv.org/abs/2502.17119v2>
- <http://arxiv.org/abs/2112.03057v1>