

Novel Evaluation Metrics for Tabular Generative Models and Their Human Judgment Correlation

Assignee Research

June 7, 2026

Abstract

This report synthesises findings from 15 peer-reviewed papers addressing the following research question: How do novel evaluation metrics for tabular generative models compare to traditional metrics (e.g., FID, Inception Score) in terms of correlation with human judgment when applied to mixed-data-type. 10 claims were extracted from source literature; 4 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 5.9/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Evaluating Generative Models for Tabular Data: Novel Metrics and Benchmarking. Research question: How do novel evaluation metrics for tabular generative models compare to traditional metrics (e.g., FID, Inception Score) in terms of correlation with human judgment when applied to mixed-data-type benchmarks?.

2 Methodology

Systematic literature search across multiple databases yielded 15 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 5.9/10.

3 Results

15 papers retrieved. 10 claims extracted; 4 independently verified. Quality review score: 5.9/10.

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.	✓	0.29
FPCAD exhibits promising performance but requires further refinements to enhance its reliability.	✓	0.19
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
FAED and FPCAD are novel and robust metrics for evaluating generative models in the tabular data domain.	✓	0.29
Quality Decrease, Mode Drop, and Mode Collapse issues were embedded into the datasets to simulate real-world generative	×	0.05
TSTR (Train on Synthetic, Test on Real) accuracy suggests that synthetic data effectively approximates real-world distri	×	0.04
TRTS (Train on Real, Test on Synthetic) score indicates that the synthetic data retains key characteristics of the real	×	0.04
TSTR is particularly useful for detecting cases where synthetic data only partially represents real data.	×	0.03
TRTS assesses whether synthetic samples introduce patterns absent in real data.	×	0.04

References

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

- <http://arxiv.org/abs/2303.04707v2>
- <http://arxiv.org/abs/2502.17119v2>