

Synthetic Data Generation Techniques and Their Impact on Tabular Foundation Model Accuracy

Assignee Research

June 7, 2026

Abstract

This report synthesises findings from 13 peer-reviewed papers addressing the following research question: How does the choice of synthetic data generation technique (e.g., GANs vs. VAEs vs. diffusion models) impact the downstream task accuracy of tabular foundation models when evaluated on heterogeneous. 14 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.2/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Diffusion and Flow Matching Models for Tabular Data: A Survey. Research question: How does the choice of synthetic data generation technique (e.g., GANs vs. VAEs vs. diffusion models) impact the downstream task accuracy of tabular foundation models when evaluated on heterogeneous benchmarks like TabularBench and OpenML?.

2 Methodology

Systematic literature search across multiple databases yielded 13 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.2/10.

3 Results

13 papers retrieved. 14 claims extracted; 0 independently verified. Quality review score: 4.2/10.

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
Future benchmarks should include missing-data settings, imbalanced classes, high-cardinality categorical features, relat	×	0.05
Formal differential privacy has begun to appear in tabular diffusion models [116], [117].	×	0.05
Recent memorization studies [118], [119] indicate that a small subset of records may dominate memorized generations.	×	0.01
Anomaly detection methods such as TCCM [71] show that feature-level residuals can support interpretability.	×	0.05
Healthcare models such as FlexGen-EHR [85] and PatientFlow [41] point toward multimodal and longitudinal tabular generat	×	0.04
Several recent methods combine diffusion or flow matching with autoencoders, transformers, tree models, or feature-token	×	0.12
Data augmentation for tabular data can be divided into two different tasks: 1) data synthesis and 2) over-sampling.	×	0.12
Over-sampling can be considered as a special case of single table synthesis where we only generate a part of the table.	×	0.02
Diffusion SOS [64] 2022 KDD Synthesis (single, generic) uses SDEs and evaluates Utility.	×	0.02
STaSy [32] 2023 ICLR Synthesis (single, generic) uses SDEs and evaluates Fidelity, Utility, Diversity.	×	0.03
TabDDPM [8] 2023 ICML Synthesis (single, generic) uses DDPM+MLD and evaluates Fidelity, Utility, Privacy.	×	0.03
CoDi [34] 2023 ICML Synthesis (single, generic) uses DDPM+MLD and evaluates Utility, Diversity.	×	0.02
AutoDiff [33] 2023 NeurIPSW Synthesis (single, generic) uses Any and evaluates Fidelity, Utility, Privacy.	×	0.03
MissDiff [76] 2023 ICMLW Synthesis (single, generic) uses SDEs and evaluates Fidelity.	×	0.01

References

- <http://arxiv.org/abs/2411.15497v3>
- <http://arxiv.org/abs/2504.20900v1>
- <http://arxiv.org/abs/2502.17119v2>