

Scaling Adversarial Robustness of Tabular Foundation Models Across Structured Data Benchmarks

Assignee Research

June 11, 2026

Abstract

The development of tabular foundation models (TFMs) has accelerated in recent years, showing strong potential to outperform traditional ML methods for structured data. A key finding is that TFMs can be pretrained entirely on synthetic datasets, opening opportunities to design data generators that encourage desirable model properties. Prior work has mainly focused on crafting high-quality priors over generators to improve overall pretraining performance. Our insight is that parameterizing the generator distribution enables an adversarial robustness perspective: during training, we can adapt the

1 Introduction

This paper examines: Robust Tabular Foundation Models. Research question: How does the adversarial robustness of tabular foundation models scale with model size when evaluated on the TabTime benchmark and other structured data benchmarks like TabMWP or TabFew?.

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 7.3/10.

3 Results

13 papers retrieved. 12 claims extracted; 9 independently verified. Quality review score: 7.3/10.

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
Tabular foundation models (TFMs) rely on in-context learning (ICL) for classification and regression tasks with structure	✓	0.21
TFMs can produce high-quality predictions on new datasets in milliseconds when GPU-accelerated.	✓	0.17
Training TFMs relies on generating diverse synthetic datasets constructed from structural causal models (SCMs).	✓	0.21
All current publicly available, competitive TFMs have been pretrained on datasets generated from a fixed prior distribution	✓	0.20
Fixed priors in TFM training underrepresent certain regions of the parameter space, potentially degrading performance on	✓	0.22
State-of-the-art TFMs lag behind tree-based methods on some benchmarks.	×	0.14
The proposed RTFM algorithm is a model-agnostic two-stage adversarial training algorithm for TFMs.	✓	0.16
Applying RTFM to TabPFN V2 using only 90k additional training datasets significantly improves its ranking on several real	✓	0.19
The maximization stage of the proposed method uses a black-box optimization algorithm to search the SCM parameter space	✓	0.24
With $n_{ds}=20$, $e=7$, and sufficient CPU cores, the estimated optimality gap can be computed in a matter of seconds.	✓	0.21
The benchmark table includes synthetic dataset configurations with feature counts ranging from 5 to 128.	×	0.02
The benchmark table includes synthetic dataset configurations using activation functions such as tanh, identity, elu, an	×	0.03

References

- <http://arxiv.org/abs/2007.08428v4>

- <http://arxiv.org/abs/2103.15670v3>
- <http://arxiv.org/abs/2512.03307v1>