

Multimodal Tabular Foundation Models Alignment with Human Preferences on TabMWP: Synthetic-Real Data Pretraining Effects

Assignee Research

June 8, 2026

Abstract

This report synthesises findings from 8 peer-reviewed papers addressing the following research question: How does the alignment of multimodal tabular foundation models with human preferences on TabMWP change when using mixed synthetic and real data for pretraining compared to real-data-only pretraining. 10 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.3/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Robust Tabular Foundation Models. Research question: How does the alignment of multimodal tabular foundation models with human preferences on TabMWP change when using mixed synthetic and real data for pretraining compared to real-data-only pretraining?.

2 Methodology

Systematic literature search across multiple databases yielded 8 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 3.3/10.

3 Results

8 papers retrieved. 10 claims extracted; 0 independently verified. Quality review score: 3.3/10.

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 structured	×	0.12
TFMs can produce high-quality predictions on new datasets in milliseconds when GPU-accelerated.	×	0.07
Training TFMs relies on generating large amounts of diverse synthetic datasets constructed from structural causal models	×	0.07
All current publicly available, competitive TFMs have been pretrained on datasets generated from a fixed prior distribution	×	0.05
Fixed priors in TFM training underrepresent certain regions of the parameter space, potentially degrading performance on	×	0.05
State-of-the-art TFMs lag behind tree-based methods on some benchmarks.	×	0.06
The proposed RTFM algorithm is a model-agnostic two-stage adversarial training algorithm for TFMs.	×	0.12
Applying RTFM to TabPFN V2 with only 90k additional training datasets significantly improves its ranking on several real	×	0.12
The optimality gap estimation step in the proposed method can be computed in a matter of seconds when parallelized with	×	0.02
In the described implementation, the optimality gap estimation used $nds=20$, $e=7$, and $ncores=140$.	×	0.03

References

- <http://arxiv.org/abs/2503.14504v2>
- <http://arxiv.org/abs/2512.03307v1>
- <http://arxiv.org/abs/2407.14477v4>