

Causal Mixture Fine-Tuning Enhances Alignment Stability Under Data Scarcity

Assignee Research

June 9, 2026

Abstract

This report synthesises findings from 11 peer-reviewed papers addressing the following research question: Does causal mixture fine-tuning improve alignment stability in language models compared to standard data augmentation under scarce validation conditions. 9 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.3/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Causal Data Augmentation for Robust Fine-Tuning of Tabular Foundation Models. Research question: Does causal mixture fine-tuning improve alignment stability in language models compared to standard data augmentation under scarce validation conditions?.

2 Methodology

Systematic literature search across multiple databases yielded 11 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.3/10.

3 Results

11 papers retrieved. 9 claims extracted; 0 independently verified. Quality review score: 4.3/10.

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
CausalMixFT achieves the highest median improvement of $(+0.12 \pm 0.63)$ over the pre-trained model on 33 classification data	×	0.09
Default fine-tuning has a variability of ± 0.98 , while CausalMixFT has a variability of ± 0.63 , indicating greater instability	×	0.08
CausalMixFT ranks first overall in average ranks across datasets, followed by the default fine-tuning baseline, while pu	×	0.07
The normalization strategy used to compare performance across different data generators is based on the zero-shot perfor	×	0.05
CausalMixFT extends the fine-tuning framework by mixing real and causally grounded synthetic samples into the fine-tunin	×	0.13
SCM-Based Synthetic Augmentation (CausalMixFT) uses SCMs fitted to the target dataset to generate synthetic data that re	×	0.12
The PC and FCI algorithms are used to estimate the structural relations between the features in CausalMixFT.	×	0.03
DoWhy’s SCM framework with additive noise models is used to sample and fit DAGs in CausalMixFT.	×	0.02
Numerical features are modeled with regressors, and categorical features with classifiers in CausalMixFT.	×	0.01

References

- <http://arxiv.org/abs/2601.04110v2>

- <http://arxiv.org/abs/2312.10793v3>
- <http://arxiv.org/abs/2110.06500v2>