

# Distribution Shifts in Synthetic Tabular Benchmarks and Multimodal Model Alignment

Assignee Research

June 9, 2026

## Abstract

This report synthesises findings from 16 peer-reviewed papers addressing the following research question: What is the impact of distribution shifts in synthetic tabular benchmarks on the alignment scores of multimodal models integrating structured and unstructured inputs. 14 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.3/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Rethinking Distribution Shifts: Empirical Analysis and Modeling for Tabular Data. Research question: What is the impact of distribution shifts in synthetic tabular benchmarks on the alignment scores of multimodal models integrating structured and unstructured inputs?.

## 2 Methodology

Systematic literature search across multiple databases yielded 16 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

16 papers retrieved. 14 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
Y  X-shifts are overlooked in tabular data.	×	0.12
Existing tabular benchmarks for distribution shifts tend to implicitly focus on X-shifts.	×	0.10
Common methods for evaluating algorithmic robustness typically concentrate on demographic subgroups within datasets such	×	0.03
The relative regret is small for widely-used benchmarks, indicating that the Y  X distribution is largely transferable a	×	0.03
The proposed benchmark incorporates distribution shifts in tabular datasets covering socioeconomic and physical systems,	×	0.09
The benchmark includes seven settings designed with varying degrees of relative regret.	×	0.03
The benchmark data foundations are designed to cover a wide range of distribution shift patterns.	×	0.11
The Y  X-ratio ranges from 0.03 to 1.00 across different quantile bins.	×	0.01
The mean Y  X-ratio varies from 0.610 to 0.821 across different quantile bins.	×	0.01
The standard deviation of Y  X-ratio varies from 0.066 to 0.304 across different quantile bins.	×	0.02
DRO performance is closely tied to the performance of its base model.	×	0.05
The base model class has a greater impact on DRO performance than the specific choice of ambiguity sets.	×	0.10
One DRO method significantly underperforms compared to others.	×	0.03
Within each DRO method class, the method with median average accuracy is reported.	×	0.04

## References

- <http://arxiv.org/abs/2112.03057v1>
- <http://arxiv.org/abs/2312.07577v3>
- <http://arxiv.org/abs/2307.05284v6>