Symbolic Reasoning Gone Awry: Unexpected Pitfalls in Deep Model Clustering

Anonymous Submission

Abstract

Deep clustering methods have increasingly been used to tackle symbolic or compositional tasks, yet they often exhibit surprising drawbacks when exposed to complexities such as varying K-values or incomplete supervision. In a real-world deployment context, these pitfalls can lead to unreliable or misleading results. We investigate these vulnerabilities through controlled experiments and highlight lessons that can help practitioners avoid similar pitfalls.

1 Introduction

Numerous contemporary approaches for structured data rely on deep neural networks to identify underlying clusters or compositional factors. Although such methods have shown promise on curated benchmarks, real-world conditions often expose unexpected behaviors, such as high overfitting and brittle performance across different parameter settings. These issues can have implications for practical deployments [??].

We systematically explore the impact of varying cluster sizes and partial supervision schemes on symbolic tasks. Our findings include: (1) evidence that some model configurations fail to adapt to even moderate domain shifts, (2) partial improvements that highlight trade-offs between complexity-weighted accuracy and validation loss, and (3) non-trivial anomalies in the training logs resulting in empty or invalid records. We hope these negative or inconclusive results will promote further discussion and improvements.

2 Related Work

Deep clustering, factorized representations, and symbolic analysis have garnered notable attention [???]. Despite reported successes, concern has grown over robustness to real-world data [?]. Unlike prior works that present strong empirical results, we document cases in which factor inference degrades when assumptions about data structure are violated, aligning with vulnerabilities noted by ?.

3 Method Discussion

We experimented with a baseline histogram-based clustering pipeline and a more flexible, autoencoder-driven model. K-values vary to test how subtask granularity affects training. We tracked validation loss, classification accuracy, and a complexity-weighted metric for capturing compositional nuances. Data processing used synthetic shapes with discrete color and form labels, mimicking realistic symbolic scenarios in robotics or vision.

4 Experiments

All runs used three random seeds. Anomalies arose in certain seeds where logs improperly recorded zero values, illustrating how subtle mistakes can mask true performance. Divergence across seeds further underscores how small hyperparameter changes can destabilize training.

The results suggest that data preprocessing and logging protocols must be carefully monitored. Even small interruptions or poorly tuned hyperparameters can derail the learning process. Our partial successes include a sharper performance increase when combining shape and color labels, though the method remains sensitive to initialization.



(a) Validation loss over epochs for varying K.

(b) Complexity-weighted metric for baseline.

Figure 1: Performance metrics for the baseline approach. Higher K-values improve factor granularity but can destabilize training.

5 Conclusion

We uncovered concrete problems that arise in deep clustering for symbolic tasks, including anomalous training logs and inconsistent performance across seeds. We encourage other researchers to test similar conditions in their pipelines and remain alert to potential pitfalls in K-based factorization. Future work includes refining pipeline components, improving logging resilience, and better uncertainty quantification.

References

A Supplementary Material

Here, we provide additional figures, per-seed lines, and tables to expand on the primary findings. Extended experiments further highlight the fragility of clustering methods under varied hyperparameters and partial label noise. Robust logging infrastructure is critical for reproducible deep clustering research.

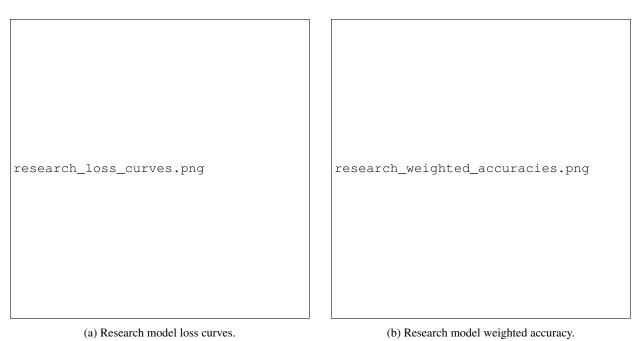


Figure 2: In the proposed method, higher complexity weighting leads to more stable performance, yet some trials produce anomalous logs.