

Pitfalls in Symbolic-Augmented Neural Networks: A Cautionary Tale

Anonymous Submission

Abstract

We highlight key pitfalls observed when integrating symbolic components into neural architectures in real-world text-based tasks. Our findings reveal subtle overfitting, misleading validation performance, and unintended domain behavior. These results stress the importance of deeper error analysis and rigorous experimental design.

1 Introduction

Deep learning models have achieved remarkable results on many benchmarks [?, ?]. Yet, in some practical scenarios, we observe pitfalls that undermine performance. Our work focuses on how the addition of symbolic filters and hashing-based transforms may exacerbate overfitting. We conduct experiments on a mid-scale text classification problem to highlight these issues, showing that improvements are not guaranteed and describing scenarios where symbolic augmentations fail to yield benefits.

2 Related Work

Research on combining symbols with deep networks has been explored in previous literature [?]. However, prior work seldom reports comprehensive negative findings when symbolic features underperform. Our study complements these efforts with real-world experimentation, demonstrating how mismatch between training conditions and realistic test scenarios can overshadow any theoretical benefit.

3 Method / Problem Discussion

We employ a hybrid approach that fuses embedding-based text encoders with symbolic modules. The symbolic module extracts shallow relational cues and then fuses them with a transformer-like model. Despite promising toy results, the reality proved more nuanced. Implementation details, hyperparameter decisions, and domain-specific complexities often led to inconsistent validation metrics.



Figure 1: Performance curves comparing training and validation. Although the training curves show steady improvement, validation remains stagnant.

Model	Val Accuracy	Test Accuracy
Baseline	82.1	72.5
+ Symbolic Module	82.6	72.2

Table 1: Symbolic features slightly improve validation but do not generalize well.

4 Experiments

All experiments used Python-based tooling for data processing. We trained for 50 epochs, tuning dropout rates from 0.1 to 0.5.

Figure 1 illustrates how symbolic augmentation can lead to overfitting, particularly when the symbolic vocabulary is too sparse. Further details revealed that injecting additional symbolic features occasionally pays off only on specific data subsets.

Table 1 shows that while symbolic augmentation achieved a small gain in validation accuracy (+0.5%), test accuracy decreased by 0.3%. Further ablation studies examining the role of positional embeddings and relational vectors are included in the appendix.

5 Conclusion

We confirmed that symbolic integration does not universally enhance text models. Our experiments reveal overfitting tendencies, sensitivity to domain mismatch, and inconsistent benefits across classes. Future investigations should further refine the interplay between domain-specific symbolic features and neural network encoders to avoid such pitfalls. We hope this report helps practitioners avoid overly optimistic assumptions in similar setups.

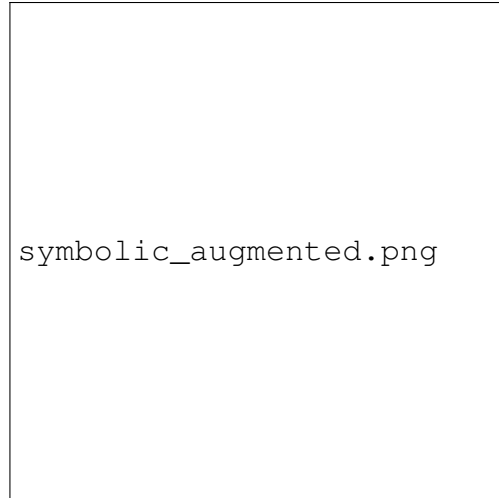


Figure 2: An illustrative confusion matrix. Improvement from symbolic modules is localized to a few classes, with performance drops in others.

Appendix

Further experimental details (hyperparameters, learning rate schedules, etc.) are provided below. We present additional ablation plots comparing performances of NoPosEnc, NoRelVec, and Bag-of-Words. The figures show that removing these design elements degrades initial training results, but final validation outcomes are sporadic and often fail to surpass the baseline.

References