AttentiveLayerAdam: Analysis of Orthogonal Constraints in Transformer Optimization

Aardvark

November 6, 2025

Abstract

This paper presents an investigation of orthogonalization constraints in transformer optimization through AttentiveLayerAdam, a modified Adam optimizer with layer-specific learning rates and attention weight orthogonalization. Our method showed progressive improvement across ablations but ultimately underperformed the AdamW baseline (4.927) with a final validation loss of 9.853. We analyze the computational overhead (23% slower than AdamW) and compare to recent orthogonal optimization approaches. While demonstrating stable orthogonal constraints are feasible, results suggest they require refinement to compete with standard approaches.

1 Introduction

We introduce AttentiveLayerAdam, combining: (1) layer-specific learning rates, (2) controlled orthogonalization of attention weights, and (3) gradient clipping. Our evaluation compares to AdamW and analyzes failure modes.

2 Method

AttentiveLayerAdam modifies Adam with:

$$m_{t} = \beta_{1} m_{t-1} + (1 - \beta_{1}) g_{t}$$

$$v_{t} = \beta_{2} v_{t-1} + (1 - \beta_{2}) g_{t}^{2}$$

$$\hat{g}_{t} = (I - \lambda W_{t} W_{t}^{T}) m_{t} / (\sqrt{v_{t}} + \epsilon)$$

3 Results

Key results:

• AdamW baseline: 4.927

• AttentiveLayerAdam: 9.853

 \bullet Training time: +23% vs AdamW

4 Conclusions

While demonstrating stable optimization, Attentive LayerAdam trails AdamW, suggesting:

- Orthogonal constraints may need adjustment
- Implementation overhead is significant
- \bullet Complementary techniques may be needed

References

[1] Loshchilov, I., Hutter, F. Decoupled weight decay regularization. arXiv:1711.05101 (2017).