

Deep Learning [CSL4020]

Intermediate Project Report 1

Submitted By:

Shalin Jain (B21CS070)
Shubh Goyal (B21CS073)

Shashwat Roy (B21CS071)
Sukriti Goyal (B21CS075)

Problem Statement

Finding the perfect model parameters giving the minimum loss is like searching for a needle in the loss landscape. However, optimizers automate and fasten the searching process through various techniques like gradient descent, analyzing gradients (partial derivatives of the loss w.r.t. parameters) to steer updates towards lower loss regions. Popular optimizers such as SGD, RMSProp and Adam further introduce momentum and gradient accumulation techniques for better learning rate adaptation, and parameter adjustments. But still, there remains room for making them faster and more efficient. We focus on formulating, developing and analyzing a new optimizer, by introducing further adaptive techniques in the existing optimizers.

Advancements After Adam

- **AdamW¹**: It decouples weight decay from learning rate, thus improving Adam's generalization ability, enabling performance competitive with SGD momentum for image classification tasks where vanilla Adam underperforms.
- **Nadam²**: It replaces the regular momentum in Adam by the Nesterov momentum, which speeds up convergence and improves model performance.
- **RAdam³**: Rectified Adam (RAdam) introduces a new term on the lines of learning rate warmup, rectifying the variance of the adaptive learning rate.
- **Lookahead⁴**: In this algorithm, two sets of weights are updated iteratively. Intuitively, the algorithm chooses a search direction by looking ahead at the sequence of "fast weights" generated by another optimizer. It helps improve the learning stability and lowers the variance of its inner optimizer.
- **AdaBelief⁵**: It introduces an exponentially moving average based 'belief' factor based on the magnitude and sign of the current gradient instead of directly averaging gradients for the updation of weights and biases. It provides increased stability and quality in GAN samples.
- **Dynamical Learning Rate Scheduling**: Adaptively adjusts the learning rate during training, especially across different stages. YellowFin⁶, an automatic tuner for momentum and learning rate in SGD, is one such example. YellowFin

optionally uses a negative-feedback loop to compensate for the momentum dynamics in asynchronous settings on the fly.

Shortcomings of these Optimizers

- Adaptive learning rate methods often use heuristics like moving averages or squared gradients, which might not accurately capture the true variance of the loss landscape.
 - Techniques like momentum that sometimes help in escaping the local minima rely on specific parameter choices and might not be effective in all scenarios.
 - Due to the adaptive nature of Adam, it becomes susceptible to the outliers in the data, thus giving poor performance on noisy sparse data.
-

Datasets and Architectures for Testing Improvements

As we gathered from different research papers on optimizers, the following datasets were used for the benchmarking of the optimisers so we also plan to use the following datasets along with the corresponding performant model architectures for evaluating our approaches.

- ResNet-18/ResNet-20 (for supervised image classification task) and GANs (for image generation task) training using *CIFAR-10*⁷ image dataset and *MNIST*⁸ digit dataset.
 - ResNet-18/ResNet-20 training on *ImageNet*⁹ dataset.
 - LSTM training on *One Billion Word Dataset*¹⁰ and *Penn Treebank Dataset*¹¹ to check performance for language modeling purposes.
-

References

1. <https://arxiv.org/pdf/1711.05101v3.pdf>
 2. <https://openreview.net/pdf?id=OM0jvwB8jlp57ZJjtNEZ>
 3. <https://arxiv.org/pdf/1908.03265v4.pdf>
 4. https://proceedings.neurips.cc/paper_files/paper/2019/file/90fd4f88f588ae64038134f1eeaa023f-Paper.pdf
 5. <https://arxiv.org/pdf/2010.07468.pdf>
 6. <https://arxiv.org/pdf/1706.03471.pdf>
 7. <https://www.cs.toronto.edu/%7Ekriz/cifar.html>
 8. <https://www.kaggle.com/datasets/hojjatk/mnist-dataset>
 9. <https://paperswithcode.com/dataset/imagenet>
 10. <https://www.statmt.org/lm-benchmark/>
 11. <https://catalog.ldc.upenn.edu/docs/LDC95T7/cl93.html>
-