Adaptive Gradient Methods with Dynamic Bound of Learning Rate

Liangchen Luo!* Yuanhao Xiong2* Yan Liu³ Xu Sun!

1PKU 2ZJU 3USC *Equal Contribution

Dilema

Adam conveges faster while generalizes worse; SGD converges slower while generalizes better.

1. Non-Convergence = Extreme Learning Rate

- · Intuitive assuption (Wilson et al., 2017): The bad performance of adaptive methods may stem from UNSTABLE and EXTREME final learning rates.
- · AMSGrad (Reddi et al., 2018), a new proposed optimizer, is claimed to have smaller learning rates compared with Adam, that may help abate the imparct of huge learning rates. However, it neglects possible effects of small ones.

Still missing 1) convincing evidence of the assumption; 2) ways to address tiny learning rates

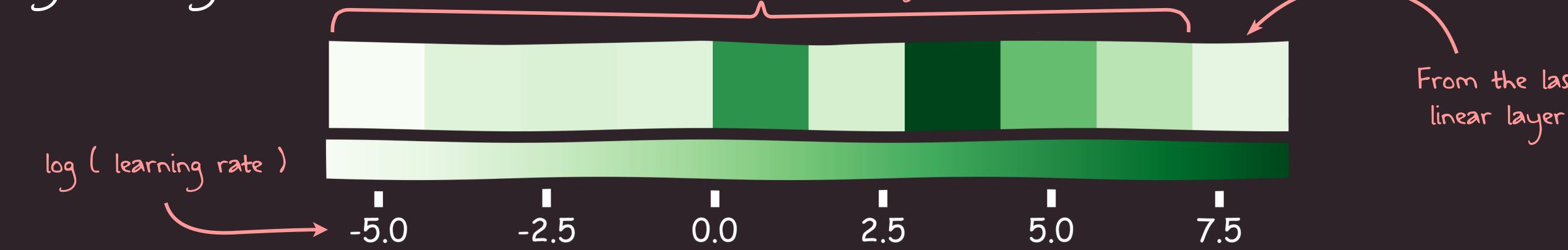
Previous

attempts &

ovservations

Evidence of the Assumption

· We sample learning rates of several weights and biases of ResNet-34 at the end of training on CIFAR-10 using Adam. Learning rates are composed of tiny ones less than 0.01 as well as huge ones greater than 1000. Each from a randomly chosen convolutional kernel



Proof of the Assumption

- Theorem 1. There is an online convex optimization problem where for any initial step size α , Adam has non-zero average regret i.e., $R_T/T \Rightarrow 0$ as $T \Rightarrow \infty$.
- Theorem 2. For any constant β_1 , $\beta_2 \in [0, 1)$ such that ${\beta_1}^2 < \beta_2$, there is an online convex optimization problem where for any initial step size α , Adam has non-zero average regret i.e., $R_T/T \Rightarrow 0$ as $T \Rightarrow \infty$.
- Theorem 3. For any constant β_1 , $\beta_2 \in [0, 1)$ such that ${\beta_1}^2 < \beta_2$, there is a stochastic convex optimization problem where for any initial step size α , Adam does not converge to the optimal solution.

#Find out why Adam fails

2. Dynamic Bound on Adam's Learning Rate

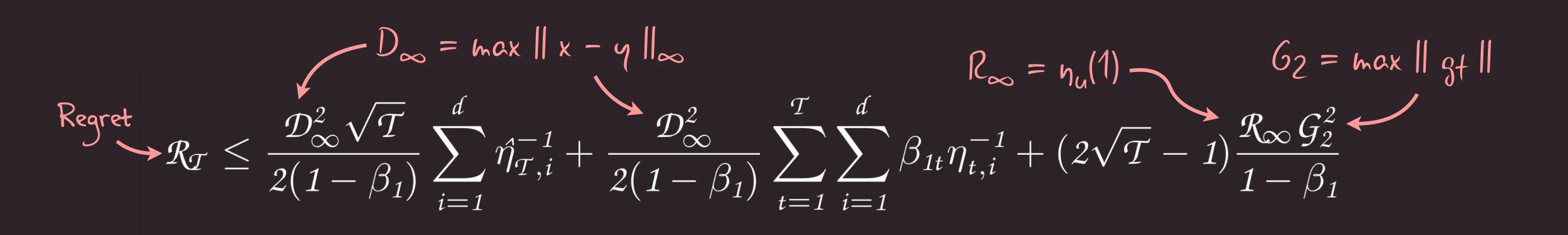
- · Our aim is to devise a strategy that combines the benefits of adaptive methods, viz. fast initial progress, and the good final generalization properties of SGD.
- · The key point is to restrict the actual learning rate of adaptive methods at the end of training. Make it neither too large nor too small.
- Inspired by Gradient Clipping, a popular technique for preventing from gradient explosion, we can also employ clipping on Adam's learning rates: An operation that clips the learning rate element-wisely such that the output is constrained to be in $[\eta_l, \eta_u]$.

Note: SGD(M) with a learning rate of $\alpha *$ can be considered as the case where $\eta_l = \eta_u = \alpha *$; as for Adam, $\eta_l = 0$ and $\eta_u = \infty$.

We propose AdaBound by employing dynamic bound rather than constant threshold:

Gradient
$$g_t = \nabla f_t(\mathbf{x}_t)$$
 Lower bound function converges from 0 to $\alpha *$;
First/Second-order $\rightarrow \mathbf{h}_t = \beta_{1t}\mathbf{h}_{t-1} + (1-\beta_{1t})g_t$ Upper bound function converges from ∞ to $\alpha *$, which makes Adam gradually transforms to SGD.
$$\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1-\beta_2)g_t^2$$
 which makes Adam gradually transforms to SGD.
$$\hat{\mathbf{h}}_t = \text{Clip}(\alpha / \sqrt{\mathbf{v}}_t, \, \mathbf{h}_t(t), \, \mathbf{h}_u(t)) \text{ and } \mathbf{h}_t = \hat{\mathbf{h}}_t / \sqrt{t}$$
 Initial step size
$$\mathbf{x}_{t+1} = \prod_{F, \, \text{diag}(\mathbf{h}_t^{-1})} (\mathbf{x}_t - \mathbf{h}_t \odot \mathbf{h}_t)$$
 Clipped learning rate

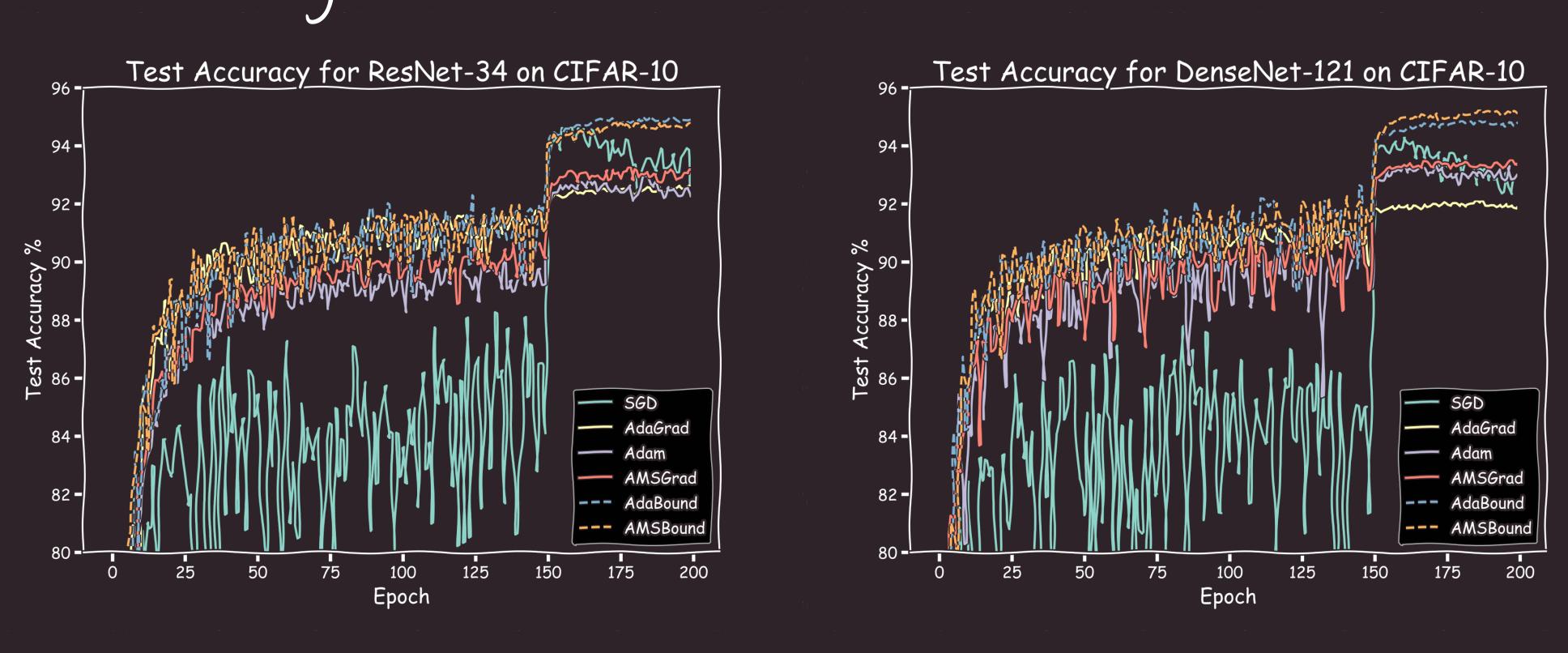
The algorithm is provable to have the following bound on regret:



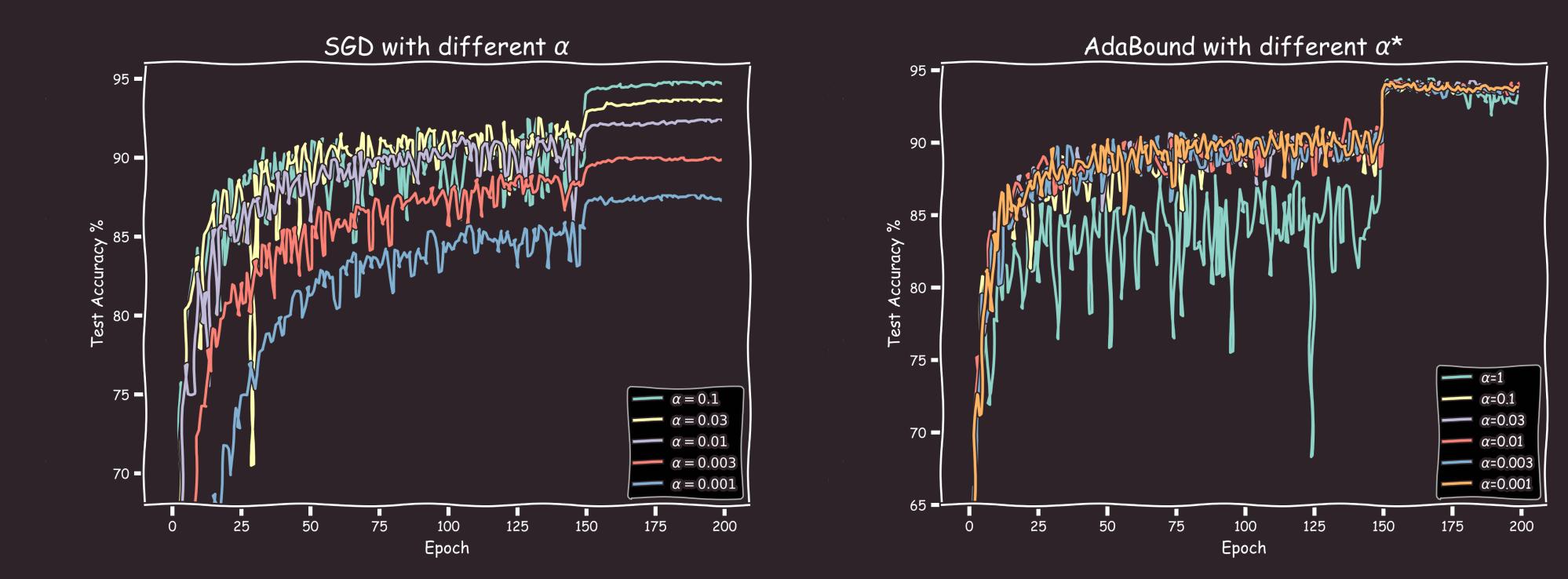
AdaBound: Transforms from Adam to SGD by applying dynamic bound on learning rate.

3. Experiments

Test accuracy for ResNet-34 and DenseNet-121 on CIFAR-10:



More rubust with different hyperparameters:





code

References

[1] S. Reddi, S. Kale, and S. Kumar. On the convergence of adam and beyond. In Proc. of ICLR, 2018.

[2] A. Wilson, R. Roelofs, M. Stern, N. Srebro, and B. Recht. The marginal value of adaptive gradient methods in machine learning. In Proc. of NeurIPS, 2017.