

CS7301: Advanced Topics in Optimization for Machine Learning

Lecture 6.2: Variants of SGD for Deep Learning

Rishabh Iyer

Department of Computer Science
University of Texas, Dallas

<https://github.com/rishabhk108/AdvancedOptML>

March 5th, 2021



- Recall SGD:

$$w_{k+1} = w_k - \alpha_k \nabla f_{i_k}(w_k)$$



SGD and Momentum

- Recall SGD:

$$w_{k+1} = w_k - \alpha_k \nabla f_{i_k}(w_k)$$

- Stochastic Momentum: Improves upon SGD via Momentum (similar to the GD case):

$$w_{k+1} = w_k - \alpha_k \nabla f_{i_k}(w_k + \gamma_t(w_k - w_{k-1})) + \beta_k(w_k - w_{k-1})$$



SGD and Momentum

- Recall SGD:

$$w_{k+1} = w_k - \alpha_k \nabla f_{i_k}(w_k)$$

- Stochastic Momentum: Improves upon SGD via Momentum (similar to the GD case):

$$w_{k+1} = w_k - \alpha_k \nabla f_{i_k}(w_k + \gamma_t(w_k - w_{k-1})) + \beta_k(w_k - w_{k-1})$$

- Heavy Ball (HB) Momentum: $\gamma_k = 0$



SGD and Momentum

- Recall SGD:

$$w_{k+1} = w_k - \alpha_k \nabla f_{i_k}(w_k)$$

- Stochastic Momentum: Improves upon SGD via Momentum (similar to the GD case):

$$w_{k+1} = w_k - \alpha_k \nabla f_{i_k}(w_k + \gamma_k(w_k - w_{k-1})) + \beta_k(w_k - w_{k-1})$$

- Heavy Ball (HB) Momentum: $\gamma_k = 0$
- Nesterov's Accelerated Gradient (NAG): $\gamma_k = \beta_k$.



Issues with SGD/Momentum

- These techniques are not adaptive: Require extensive tuning for learning rate schedules.



Issues with SGD/Momentum

- These techniques are not adaptive: Require extensive tuning for learning rate schedules.
- Without the right LR schedule, convergence can be slow!



Issues with SGD/Momentum

- These techniques are not adaptive: Require extensive tuning for learning rate schedules.
- Without the right LR schedule, convergence can be slow!
- They are also less robust to initialization



Issues with SGD/Momentum

- These techniques are not adaptive: Require extensive tuning for learning rate schedules.
- Without the right LR schedule, convergence can be slow!
- They are also less robust to initialization
- Fix: Adapt learning rate based on gradient information until now.



Adaptive Gradient Descent Framework of Algorithms

- Adaptive Methods try to automatically adapt the learning rate.



Adaptive Gradient Descent Framework of Algorithms

- Adaptive Methods try to automatically adapt the learning rate.
- Define H_k as a positive semi-definite Matrix (recall the Hessian method?)



Adaptive Gradient Descent Framework of Algorithms

- Adaptive Methods try to automatically adapt the learning rate.
- Define H_k as a positive semi-definite Matrix (recall the Hessian method?)
- The simplest Adaptive Algorithm called AdaGrad (Duchi et al 2011) is:

$$w_{k+1} = w_k - \alpha_k H_k^{-1} \nabla f_{i_k}(w_k)$$

where H_k is a positive semi-definite Matrix!



Adaptive Gradient Descent Framework of Algorithms

- Adaptive Methods try to automatically adapt the learning rate.
- Define H_k as a positive semi-definite Matrix (recall the Hessian method?)
- The simplest Adaptive Algorithm called AdaGrad (Duchi et al 2011) is:

$$w_{k+1} = w_k - \alpha_k H_k^{-1} \nabla f_{i_k}(w_k)$$

where H_k is a positive semi-definite Matrix!

- The simplest definition of H_k is a diagonal matrix (recall we need SGD algorithms to be super-fast, so no matrix inversion possible!)



Adaptive Gradient Descent Framework of Algorithms

- Adaptive Methods try to automatically adapt the learning rate.
- Define H_k as a positive semi-definite Matrix (recall the Hessian method?)
- The simplest Adaptive Algorithm called AdaGrad (Duchi et al 2011) is:

$$w_{k+1} = w_k - \alpha_k H_k^{-1} \nabla f_{i_k}(w_k)$$

where H_k is a positive semi-definite Matrix!

- The simplest definition of H_k is a diagonal matrix (recall we need SGD algorithms to be super-fast, so no matrix inversion possible!)
- Define:

$$H_k = \text{diag}(\{\sum_{i=1}^k \eta_i g_i \circ g_i\}^{1/2})$$



One Equation to Unify them All!

- One can also come up with a Momentum version of AdaGrad (called Adam).



One Equation to Unify them All!

- One can also come up with a Momentum version of AdaGrad (called Adam).
- We can unify all adaptive and non-adaptive variants into a single update equation:

$$w_{k+1} = w_k - \alpha_k H_k^{-1} \nabla f_{i_k}(w_k + \gamma_k(w_k - w_{k-1})) + \beta_k H_k^{-1} H_{k-1}(w_k - w_{k-1})$$



One Equation to Unify them All!

- One can also come up with a Momentum version of AdaGrad (called Adam).
- We can unify all adaptive and non-adaptive variants into a single update equation:

$$w_{k+1} = w_k - \alpha_k H_k^{-1} \nabla f_{i_k}(w_k + \gamma_k(w_k - w_{k-1})) + \beta_k H_k^{-1} H_{k-1}(w_k - w_{k-1})$$

- Recall $H_k = \text{diag}(\{\sum_{i=1}^k \eta_i g_i \circ g_i\}^{1/2})$. Define $G_k = H_k \circ H_k$ and $D_k = \text{diag}(g_k \circ g_k)$.



One Equation to Unify them All!

- We can unify all adaptive and non-adaptive variants into a single update equation:

$$w_{k+1} = w_k - \alpha_k H_k^{-1} \nabla f_{i_k}(w_k + \gamma_k(w_k - w_{k-1})) + \beta_k H_k^{-1} H_{k-1}(w_k - w_{k-1})$$

- Recall $H_k = \text{diag}(\{\sum_{i=1}^k \eta_i g_i \circ g_i\}^{1/2})$. Define $G_k = H_k \circ H_k$ and $D_k = \text{diag}(g_k \circ g_k)$.

	SGD	HB	NAG	AdaGrad	RMSProp	Adam
G_k	I	I	I	$G_{k-1} + D_k$	$\beta_2 G_{k-1} + (1 - \beta_2) D_k$	$\frac{\beta_2}{1 - \beta_2^k} G_{k-1} + \frac{(1 - \beta_2)}{1 - \beta_2^k} D_k$
α_k	α	α	α	α	α	$\alpha \frac{1 - \beta_1}{1 - \beta_1^k}$
β_k	0	β	β	0	0	$\frac{\beta_1(1 - \beta_1^{k-1})}{1 - \beta_1^k}$
γ	0	0	β	0	0	0



More on ADAM

- Adam is basically HB Momentum + Adaptive.
- Define $m_k = \beta_1 m_{k-1} + (1 - \beta_1)g_k$
- Define $v_k = \beta_2 v_{k-1} + (1 - \beta_2)g_k \circ g_k$
- Intuition of m_k and v_k are estimates of first moment (mean) and second moment (uncentered variance) of the gradients.
- Since m_k and v_k are initialized to 0, they are biased towards zero when the decay rates are small. To counter this, they are further normalized by $1 - \beta^k$.
- Define $\hat{m}_k = m_k / (1 - \beta_1^k)$ and $\hat{v}_k = v_k / (1 - \beta_2^k)$.
- The ADAM update is $w_{k+1} = w_k - \alpha_k \hat{m}_k \circ \hat{v}_k^{-1/2}$
- Parameters used in practice: $\beta_1 = 0.9, \beta_2 = 0.999$.



Extensions

Numerous extensions of the above techniques

- AdaMax is an extension of ADAM to use the l_{∞} norm (i.e. max) instead of square.
- NADAM applies Nesterovs momentum instead of HB Momentum to Adaptive Methods.
- ADADelta is an extension of RMSProp to use the RMS operator on the weight differences as well.
- Recent Algorithm (AMSGrad) by Reddi et al (ICLR 2018) which fixes a theoretical error in ADAM (causing it to not converge even for convex functions) simply by ensuring v_t 's remain positive!
- See more details to compare the different optimization algorithms (and also what they are) here:

<https://ruder.io/optimizing-gradient-descent/>.



Theoretical Results

Numerous extensions of the above techniques

- The first theoretical result was shown for AdaGrad. The convergence result there is a *Regret* bound which is common for online algorithms.
- As mentioned above, the paper introducing ADAM actually had a bug in its analysis. The same also holds for RMSProp, AdaDelta and NADAM etc. They do not have **theoretical Regret bounds** backing them.
- Paper introducing AMSGrad showed regret bounds with a modified version of ADAM (and correspondingly RMSProp, NADAM, ...)
- All this only holds for convex functions. No results known for Non-Convex Functions.



Comparison of the different methods

- AdaGrad (Duchi et al 2011) is one of the most influential papers of the last decade!



Comparison of the different methods

- AdaGrad (Duchi et al 2011) is one of the most influential papers of the last decade!
- The starting point of numerous new techniques for adaptive methods.



Comparison of the different methods

- AdaGrad (Duchi et al 2011) is one of the most influential papers of the last decade!
- The starting point of numerous new techniques for adaptive methods.
- There is really no one technique that is provably better than the other. Each technique has its own pros and cons!

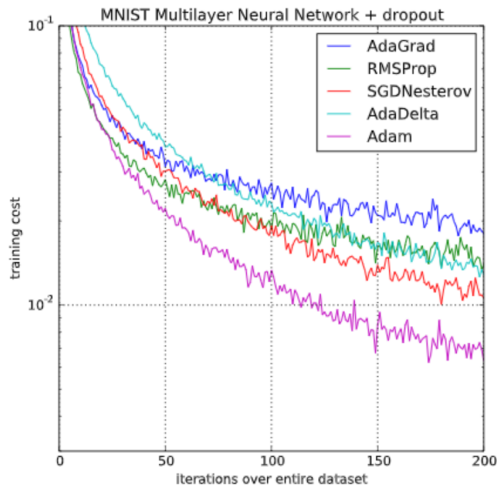


Comparison of the different methods

- AdaGrad (Duchi et al 2011) is one of the most influential papers of the last decade!
- The starting point of numerous new techniques for adaptive methods.
- There is really no one technique that is provably better than the other. Each technique has its own pros and cons!
- In the next few slides, I'll try to put together a few takeaways from some recent papers which have studied this specifically for non-convex optimization.



Kingma et al, ICLR 2015 – Original ADAM paper



Adaptive vs Non Adaptive Techniques: Comparisons

- Benefits of AdaGrad: AdaGrad can significantly improve upon SGD in sparse feature sets! It automatically sets the learning rate, and secondly, automatically updates the learning rates with a decay schedule! Also, it has a per coordinate learning rate!



Adaptive vs Non Adaptive Techniques: Comparisons

- Benefits of AdaGrad: AdaGrad can significantly improve upon SGD in sparse feature sets! It automatically sets the learning rate, and secondly, automatically updates the learning rates with a decay schedule! Also, it has a per coordinate learning rate!
- In dense settings and particularly in deep models, Adagrad works very poorly because of rapid decay in learning rates



Adaptive vs Non Adaptive Techniques: Comparisons

- Benefits of AdaGrad: AdaGrad can significantly improve upon SGD in sparse feature sets! It automatically sets the learning rate, and secondly, automatically updates the learning rates with a decay schedule! Also, it has a per coordinate learning rate!
- In dense settings and particularly in deep models, Adagrad works very poorly because of rapid decay in learning rates
- ADAM, RMSProp, AdaDelta, ... all try to fix this issue!



Adaptive vs Non Adaptive Techniques: Comparisons

- Benefits of AdaGrad: AdaGrad can significantly improve upon SGD in sparse feature sets! It automatically sets the learning rate, and secondly, automatically updates the learning rates with a decay schedule! Also, it has a per coordinate learning rate!
- In dense settings and particularly in deep models, Adagrad works very poorly because of rapid decay in learning rates
- ADAM, RMSProp, AdaDelta, ... all try to fix this issue!
- In many cases, these adaptive algorithms improve upon SGD in terms of training loss and better/faster convergence!



Adaptive vs Non Adaptive Techniques: Comparisons

- Benefits of AdaGrad: AdaGrad can significantly improve upon SGD in sparse feature sets! It automatically sets the learning rate, and secondly, automatically updates the learning rates with a decay schedule! Also, it has a per coordinate learning rate!
- In dense settings and particularly in deep models, Adagrad works very poorly because of rapid decay in learning rates
- ADAM, RMSProp, AdaDelta, ... all try to fix this issue!
- In many cases, these adaptive algorithms improve upon SGD in terms of training loss and better/faster convergence!
- Also, momentum generally improves upon the non-momentum variants!



Adaptive vs Non Adaptive Techniques: Comparisons

- Benefits of AdaGrad: AdaGrad can significantly improve upon SGD in sparse feature sets! It automatically sets the learning rate, and secondly, automatically updates the learning rates with a decay schedule! Also, it has a per coordinate learning rate!
- In dense settings and particularly in deep models, Adagrad works very poorly because of rapid decay in learning rates
- ADAM, RMSProp, AdaDelta, ... all try to fix this issue!
- In many cases, these adaptive algorithms improve upon SGD in terms of training loss and better/faster convergence!
- Also, momentum generally improves upon the non-momentum variants!
- But....



Back to SGD! Generalization....

- But, in Machine Learning, Generalization is more important compared to just Training Loss!



Back to SGD! Generalization....

- But, in Machine Learning, Generalization is more important compared to just Training Loss!
- Recent works (for example, Wilson et al, The Marginal Value of Adaptive Gradient Methods in Machine Learning) showed a very surprising result!



Back to SGD! Generalization....

- But, in Machine Learning, Generalization is more important compared to just Training Loss!
- Recent works (for example, Wilson et al, The Marginal Value of Adaptive Gradient Methods in Machine Learning) showed a very surprising result!
- Though adaptive gradient methods tend to minimize training loss better, they do so by obtaining more complex and less generalizable solutions!



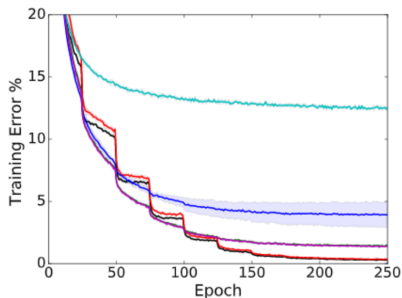
Back to SGD! Generalization....

- But, in Machine Learning, Generalization is more important compared to just Training Loss!
- Recent works (for example, Wilson et al, The Marginal Value of Adaptive Gradient Methods in Machine Learning) showed a very surprising result!
- Though adaptive gradient methods tend to minimize training loss better, they do so by obtaining more complex and less generalizable solutions!
- They gave a few synthetic examples (particularly in overparameterized scenarios) where SGD and its variants obtain the less complex solutions but Adaptive variants obtain solutions which do not generalize well!

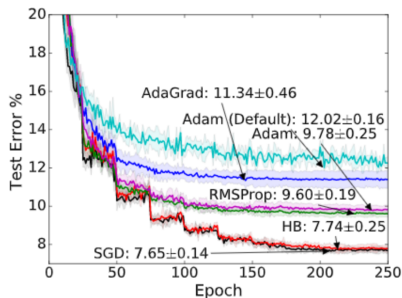


Back to SGD! Generalization....

See Wilson et al, The Marginal Value of Adaptive Gradient Methods in Machine Learning, NeurIPS 2017



(a) CIFAR-10 (Train)



(b) CIFAR-10 (Test)



Additional Reading

- Wilson et al, The Marginal Value of Adaptive Gradient Methods in Machine Learning, NeurIPS 2017
- Reddi et al, On the Convergence of ADAM and Beyond, ICLR 2018.
- Kingma and Ba, ADAM: A Method for Stochastic Optimization, ICLR 2015
- Duchi et al, Adaptive subgradient methods for online learning and stochastic optimization, Journal of Machine Learning Research 2011.
- Zeiler. ADADELTA: An Adaptive Learning Rate Method, ArXiv 2012
- <https://ruder.io/optimizing-gradient-descent/>

