

ADAM (Adaptive Moments) I

- The update rule (Kingma and Ba, 2015)

$$\mathbf{g} \leftarrow \frac{\boldsymbol{\theta}}{C} + \frac{1}{|S|} \nabla_{\boldsymbol{\theta}} \sum_{i:i \in S} \xi(\boldsymbol{\theta}; \mathbf{y}^i, Z^{1,i})$$

$$\mathbf{s} \leftarrow \rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{g}$$

$$\mathbf{r} \leftarrow \rho_2 \mathbf{r} + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$$

$$\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \rho_1^t}$$

$$\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_2^t}$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \frac{\epsilon}{\sqrt{\hat{\mathbf{r}} + \delta}} \odot \hat{\mathbf{s}}$$

ADAM (Adaptive Moments) II

- t is the current iteration index
- Roughly speaking, ADAM is the combination of
 - Momentum
 - RMSprop
- From Goodfellow et al. (2016),

$$\frac{\epsilon}{\sqrt{\hat{r}} + \delta} \odot \hat{\mathbf{s}}$$

(i.e., the use of momentum combined with rescaling) “does not have a clear theoretical motivation”

ADAM (Adaptive Moments) III

- The two steps

$$\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \rho_1^t}$$
$$\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_2^t}$$

are called “bias correction”

- Why “bias correction”?

ADAM (Adaptive Moments) IV

- They hope that

$$E[\mathbf{s}_t] = E[\mathbf{g}_t]$$

and

$$E[\mathbf{r}_t] = E[\mathbf{g}_t \odot \mathbf{g}_t],$$

where t is the iteration index

ADAM (Adaptive Moments) V

- For \mathbf{s}_t , we have

$$\begin{aligned}\mathbf{s}_t &= \rho_1 \mathbf{s}_{t-1} + (1 - \rho_1) \mathbf{g}_t \\ &= \rho_1 (\rho_1 \mathbf{s}_{t-2} + (1 - \rho_1) \mathbf{g}_{t-1}) + (1 - \rho_1) \mathbf{g}_t \\ &= (1 - \rho_1) \sum_{i=1}^t \rho_1^{t-i} \mathbf{g}_i\end{aligned}$$

We assume that \mathbf{s} is initialized as 0

ADAM (Adaptive Moments) VI

- Then

$$\begin{aligned} E[\mathbf{s}_t] &= E[(1 - \rho_1) \sum_{i=1}^t \rho_1^{t-i} \mathbf{g}_i] \\ &= E[\mathbf{g}_t] (1 - \rho_1) \sum_{i=1}^t \rho_1^{t-i} \end{aligned}$$

- Note that we assume

$$E[\mathbf{g}_i], \forall i \geq 1$$

are the same

ADAM (Adaptive Moments) VII

- Next,

$$\begin{aligned} & (1 - \rho_1) \sum_{i=1}^t \rho_1^{t-i} \\ &= (1 - \rho_1)(1 + \dots + \rho_1^{t-1}) \\ &= 1 - \rho_1^t \end{aligned}$$

- Thus

$$E[\mathbf{s}_t] = E[\mathbf{g}_t](1 - \rho_1^t)$$

and they do

$$\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \rho_1^t}$$

ADAM (Adaptive Moments) VIII

- The above derivation on bias correction partially follows from <https://towardsdatascience.com/adam-latest-trends-in-deep-learning-optimiz>
- The situation for $E[\mathbf{g}_t \odot \mathbf{g}_t]$ is similar
- How about ADAM's practical performance?
- From Goodfellow et al. (2016), “generally regarded as being **fairly robust to the choice of hyperparameters**, though the learning rate may need to be changed from the default”

ADAM (Adaptive Moments) IX

- However, from the web page we referred to for deriving the bias correction, “The original paper ... showing huge performance gains in terms of speed of training. However, after a while people started noticing, that in some cases Adam actually **finds worse solution than stochastic gradient**”
- One example of showing the above is Wilson et al. (2017)

Choosing Stochastic Gradient Algorithms

- From Goodfellow et al. (2016), “there is currently **no consensus**”
- Further, “the choice ... seemed to depend on the user’s familiarity with the algorithm”

Why Stochastic Gradient Widely Used? I

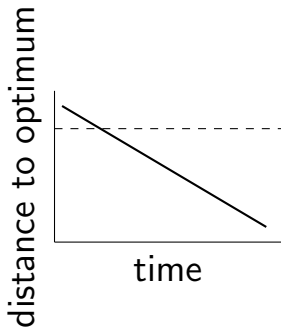
- In machine learning fast final convergence may not be important
 - An optimal solution θ^* may not lead to the best model
 - Further, we don't need a point close to θ^* . In prediction we find

$$\arg \max_k z_k^{L+1}(\theta)$$

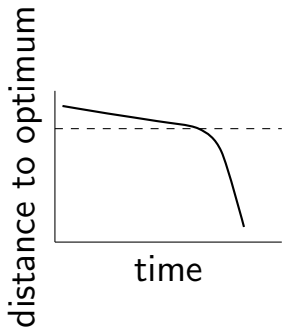
A not-so-accurate θ may be good enough

An illustration

Why Stochastic Gradient Widely Used? II



Slow final convergence



Fast final convergence

Why Stochastic Gradient Widely Used? III

- The special property of data classification is essential

$$E(\nabla_{\theta} \xi(\mathbf{z}^{L+1}; \mathbf{x}, \mathbf{y})) = \frac{1}{I} \nabla_{\theta} \sum_{i=1}^I \xi(\mathbf{z}^{L+1,i}(\theta); \mathbf{x}^i, \mathbf{y}^i)$$

- We can cheaply get a good approximation of the gradient
- Indeed stochastic gradient is less used outside machine learning

Why Stochastic Gradient Widely Used? IV

- Easy implementation. It's simpler than methods using, for example, second derivative
Now for complicated networks, (subsampled) gradient is calculated by **automatic differentiation**
- We will explain more about this
- Non-convexity plays a role
 - For convex, other methods may possess advantages to more efficiently find **the global minimum**
 - But for non-convex, efficiency to reach a **stationary point** is less useful

Why Stochastic Gradient Widely Used? V

- A global minimum usually gives a good model (as loss is minimized), but for a stationary point we are less sure
- Some variants of SG have been proposed to improve the robustness or the convergence
- All these explain why SG is popular for deep learning

References I

- I. J. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. The MIT Press, 2016.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2015.
- A. C. Wilson, R. Roelofs, M. Stern, N. Srebro, and B. Recht. The marginal value of adaptive gradient methods in machine learning. In *Advances in Neural Information Processing Systems*, pages 4148–4158, 2017.