Adam (Adaptive Moments) I

• The update rule (Kingma and Ba, 2015)

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

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

$$r \leftarrow \rho_{2} r + (1 - \rho_{2}) g \odot g$$

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

$$\hat{r} \leftarrow \frac{r}{1 - \rho_{2}^{t}}$$

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

Adam (Adaptive Moments) II

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

$$rac{\epsilon}{\sqrt{\hat{\pmb{r}}}+\delta}\odot \hat{\pmb{s}}$$

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

Adam (Adaptive Moments) III

- How about Adam's practical performance?
- From Goodfellow et al. (2016), "generally regarded as being fairly robust to the choice of hyperparmeters, though the learning rate may need to be changed from the default"
- 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"

イロト イヨト イヨト ・

Adam (Adaptive Moments) IV

• One example of showing the above is Wilson et al. (2017)

イロト 不得 トイヨト イヨト

Bias Correction in Adam I

• The two steps in Adam

$$\hat{m{s}} \leftarrow rac{m{s}}{1-
ho_1^t} \ \hat{m{r}} \leftarrow rac{m{r}}{1-
ho_2^t}$$

are called "bias correction"

- Why do we need this "bias correction"?
- Note that s is the direction used to update θ .

< □ > < □ > < □ > < □ > < □ > < □ >

Bias Correction in Adam II

• We hope that its expectation is similar to the expected gradient

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

and

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

where t is the iteration index

- The problem is that due to the moving average, the vector is biased toward the initial value
- Note that our initial **s** is **0**

< □ > < 同 > < 回 > < 回 > < 回 >

Bias Correction in Adam III

• For s_t , we have

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

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Bias Correction in Adam IV

• Then

$$E[\mathbf{s}_{t}] = E[(1 - \rho_{1}) \sum_{i=1}^{t} \rho_{1}^{t-i} \mathbf{g}_{i}]$$
$$= E[\mathbf{g}_{t}](1 - \rho_{1}) \sum_{i=1}^{t} \rho_{1}^{t-i}$$

• Note that we assume

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

are the same

Chih-Jen Lin (National Taiwan Univ.)

э

イロン イ理 とくほとう ほんし

Bias Correction in Adam V

• Next,

$$(1 -
ho_1) \sum_{i=1}^t
ho_1^{t-i} = (1 -
ho_1)(1 + \dots +
ho_1^{t-1}) = 1 -
ho_1^t$$

Thus

$$E[\boldsymbol{s}_t] = E[\boldsymbol{g}_t](1-
ho_1^t)$$

and they do

$$\hat{\boldsymbol{s}} \leftarrow \frac{\boldsymbol{s}}{1 - \rho_1^t} + \boldsymbol{r} + \boldsymbol{s} + \boldsymbol{s}$$

Bias Correction in Adam VI

- The above derivation on bias correction partially follows from https://towardsdatascience.com/ adam-latest-trends-in-deep-learning-optimiz
- The situation for $E[\boldsymbol{g}_t \odot \boldsymbol{g}_t]$ is similar

イロト イヨト イヨト ・

The Importance of Bias Correction I

- An interesting story is that BERT (Devlin et al., 2019), an important NLP technique using Adam, forgot to do bias correction
- This seems to cause lengthy iterations
- See Zhang et al. (2021) for discussing this issue

< □ > < □ > < □ > < □ > < □ > < □ >

Weight Decay I

Recall in our earlier description, the simple stochastic gradient update is

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

 $\frac{\theta}{C}$

• In this rule,

comes from the regularization term $\theta^T \theta/(2C)$ in $f(\theta)$

Weight Decay II

- The use of regularization follows from standard machine learning settings
- However, in the area of neural networks, this term may come from a setting called weight decay (Hanson and Pratt, 1988)

$$oldsymbol{ heta} \leftarrow (1-\lambda)oldsymbol{ heta} - \eta(rac{1}{|oldsymbol{S}|}
abla_{oldsymbol{ heta}}\sum_{i:i\inoldsymbol{S}}\xi(oldsymbol{ heta};oldsymbol{y}^i,Z^{1,i}))$$

where λ is the rate of weight decay

• In fact, Hanson and Pratt (1988) did not give good reasons for decaying the weight of θ

Weight Decay III

• Clearly, if

$$\lambda = \frac{\eta}{C}$$

then weight decay is the same as regularization

 However, as pointed out in Loshchilov and Hutter (2019), the equivalence does not hold if adaptive learning rate is used

イロト 不得下 イヨト イヨト

Weight Decay IV

• For example, in AdaGrad, the update rule is

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \frac{\epsilon}{\sqrt{\boldsymbol{r}} + \delta} \odot \left(\frac{1}{|S|} \nabla_{\boldsymbol{\theta}} \sum_{i:i \in S} \xi(\boldsymbol{\theta}; \boldsymbol{y}^{i}, Z^{1,i})\right) \\ - \frac{\epsilon}{\sqrt{\boldsymbol{r}} + \delta} \odot \frac{\boldsymbol{\theta}}{C}$$

so the regularization term is scaled in a component-wise way

• Loshchilov and Hutter (2019) advocate to decouple the weight decay step

Weight Decay V

• For example, for the momentum algorithm

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

they prefer the following equivalent form

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



 Based on this, Loshchilov and Hutter (2019) proposed AdamW

э

イロト 不得 トイヨト イヨト

AdamW I

$$\mathbf{g} \leftarrow \frac{1}{|S|} \nabla_{\boldsymbol{\theta}} \sum_{i:i \in S} \xi(\boldsymbol{\theta}; \mathbf{y}^{i}, Z^{1,i}) \\
 \mathbf{s} \leftarrow \rho_{1} \mathbf{s} + (1 - \rho_{1}) \mathbf{g} \\
 \mathbf{r} \leftarrow \rho_{2} \mathbf{r} + (1 - \rho_{2}) \mathbf{g} \odot \mathbf{g} \\
 \hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \rho_{1}^{t}} \\
 \hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_{2}^{t}} \\
 \boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \frac{\epsilon}{\sqrt{\hat{\mathbf{r}}} + \delta} \odot \hat{\mathbf{s}} - \epsilon \frac{\boldsymbol{\theta}}{\zeta}$$

2

AdamW II

- This is not equivalent to Adam because in Adam, θ/C has been used in calculating g and then scaled after
- Why is the decoupled setting better? Some discussions are in Section 3 of Loshchilov and Hutter (2019)

< □ > < □ > < □ > < □ > < □ > < □ >

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"

< ロト < 同ト < ヨト < ヨト

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

$$rg\max_k z_k^{L+1}(oldsymbol{ heta})$$

A not-so-accurate θ may be good enough An illustration (modified from Tsai et al. (2014))

イロト イヨト イヨト ・

Why Stochastic Gradient Widely Used? II



Slow final convergence Fast final convergence

< ロト < 同ト < ヨト < ヨト

Why Stochastic Gradient Widely Used? III

• The special property of data classification is essential

$$E(\nabla_{\boldsymbol{\theta}}\xi(\boldsymbol{\theta}; \mathbf{y}^{i}, Z^{1,i})) = \frac{1}{l} \nabla_{\boldsymbol{\theta}} \sum_{i=1}^{l} \xi(\boldsymbol{\theta}; \mathbf{y}^{i}, Z^{1,i})$$

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

< □ > < 同 > < 回 > < 回 > < 回 >

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

(日) (周) (三) (三)

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

< ロト < 同ト < ヨト < ヨト

References I

- J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In J. Burstein, C. Doran, and T. Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT*, pages 4171–4186, 2019. doi: 10.18653/v1/n19-1423.
- I. J. Goodfellow, Y. Bengio, and A. Courville. Deep Learning. The MIT Press, 2016.
- S. Hanson and L. Pratt. Comparing biases for minimal network construction with back-propagation. In Advances in Neural Information Processing Systems, volume 1, 1988.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In Proceedings of International Conference on Learning Representations (ICLR), 2015.
- I. Loshchilov and F. Hutter. Decoupled weight decay regularization. In *Proceedings of International Conference on Learning Representations*, 2019.
- C.-H. Tsai, C.-Y. Lin, and C.-J. Lin. Incremental and decremental training for linear classification. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2014. URL http://www.csie.ntu.edu.tw/~cjlin/papers/ws/inc-dec.pdf.
- 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.
- T. Zhang, F. Wu, A. Katiyar, K. Q. Weinberger, and Y. Artzi. Revisiting few-sample BERT fine-tuning. In *Proceedings of International Conference on Learning Representations*, 2021, 2011