

Stochastic Gradient Methods for Neural Networks

Chih-Jen Lin
National Taiwan University

Last updated: May 25, 2020

Outline

- 1 Gradient descent
- 2 Mini-batch SG
- 3 Adaptive learning rate
- 4 Discussion



Outline

- 1 Gradient descent
- 2 Mini-batch SG
- 3 Adaptive learning rate
- 4 Discussion



NN Optimization Problem I

- Recall that the NN optimization problem is

$$\min_{\theta} f(\theta)$$

where

$$f(\theta) = \frac{1}{2C} \theta^T \theta + \frac{1}{l} \sum_{i=1}^l \xi(z^{L+1,i}(\theta); y^i, Z^{1,i})$$

- Let's simplify the loss part a bit

$$f(\theta) = \frac{1}{2C} \theta^T \theta + \frac{1}{l} \sum_{i=1}^l \xi(\theta; y^i, Z^{1,i})$$

- The issue now is how to do the minimization



Gradient Descent I

- This is one of the most used optimization method
- First-order approximation

$$f(\boldsymbol{\theta} + \Delta\boldsymbol{\theta}) \approx f(\boldsymbol{\theta}) + \nabla f(\boldsymbol{\theta})^T \Delta\boldsymbol{\theta}$$

- Solve

$$\begin{aligned} \min_{\Delta\boldsymbol{\theta}} \quad & \nabla f(\boldsymbol{\theta})^T \Delta\boldsymbol{\theta} \\ \text{subject to} \quad & \|\Delta\boldsymbol{\theta}\| = 1 \end{aligned} \quad (1)$$

- If no constraint, the above sub-problem goes to $-\infty$

Gradient Descent II

- The solution of (1) is

$$\Delta\theta = -\frac{\nabla f(\theta)}{\|\nabla f(\theta)\|}$$

- This is called **steepest descent method**
- In general all we need is a **descent direction**

$$\nabla f(\theta)^T \Delta\theta < 0$$



Gradient Descent III

- From

$$f(\boldsymbol{\theta} + \alpha\Delta\boldsymbol{\theta}) = f(\boldsymbol{\theta}) + \alpha\nabla f(\boldsymbol{\theta})^T \Delta\boldsymbol{\theta} + \frac{1}{2}\alpha^2\Delta\boldsymbol{\theta}^T \nabla^2 f(\boldsymbol{\theta})\Delta\boldsymbol{\theta} + \dots,$$

if

$$\nabla f(\boldsymbol{\theta})^T \Delta\boldsymbol{\theta} < 0,$$

then with a small enough α ,

$$f(\boldsymbol{\theta} + \alpha\Delta\boldsymbol{\theta}) < f(\boldsymbol{\theta})$$



Line Search I

- Because we only consider an approximation

$$f(\boldsymbol{\theta} + \Delta\boldsymbol{\theta}) \approx f(\boldsymbol{\theta}) + \nabla f(\boldsymbol{\theta})^T \Delta\boldsymbol{\theta}$$

we may not have the strict decrease of the function value

- That is,

$$f(\boldsymbol{\theta}) < f(\boldsymbol{\theta} + \Delta\boldsymbol{\theta})$$

may occur

- In optimization we then need a **step selection** procedure



Line Search II

- Exact line search

$$\min_{\alpha} f(\boldsymbol{\theta} + \alpha\Delta\boldsymbol{\theta})$$

This is a one-dimensional optimization problem

- In practice, people use **backtracking line search**
- We check

$$\alpha = 1, \beta, \beta^2, \dots$$

with $\beta \in (0, 1)$ until

$$f(\boldsymbol{\theta} + \alpha\Delta\boldsymbol{\theta}) < f(\boldsymbol{\theta}) + \nu\nabla f(\boldsymbol{\theta})^T(\alpha\Delta\boldsymbol{\theta})$$



Line Search III

- Here

$$\nu \in (0, \frac{1}{2})$$

- The convergence is well established.
- For example, under some conditions, Theorem 3.2 of Nocedal and Wright (1999) has that

$$\lim_{k \rightarrow \infty} \nabla f(\boldsymbol{\theta}^k) = 0,$$

where k is the iteration index

- This means we can reach a **stationary point** of a non-convex problem



Practical Use of Gradient Descent I

- The standard back-tracking line search is simple and useful
- However, the convergence is slow for difficult problems
- Thus in many optimization applications, methods of using second-order information (e.g., quasi Newton or Newton) are preferred

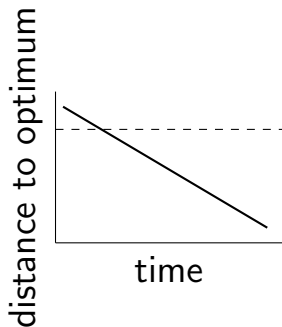
$$f(\boldsymbol{\theta} + \Delta\boldsymbol{\theta}) \approx f(\boldsymbol{\theta}) + \nabla f(\boldsymbol{\theta})^T \Delta\boldsymbol{\theta} + \frac{1}{2} \Delta\boldsymbol{\theta}^T \nabla^2 f(\boldsymbol{\theta}) \Delta\boldsymbol{\theta}$$

- These methods have fast final convergence

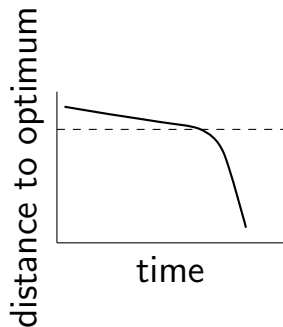


Practical Use of Gradient Descent II

- An illustration (modified from Tsai et al. (2014))



Slow final convergence



Fast final convergence

Practical Use of Gradient Descent III

- But fast final convergence may not be needed in machine learning
- The reason is that an optimal solution θ^* may not lead to the best model
- We will discuss such issues again later



Outline

- 1 Gradient descent
- 2 Mini-batch SG**
- 3 Adaptive learning rate
- 4 Discussion



Estimation of the Gradient I

- Recall the function is

$$f(\boldsymbol{\theta}) = \frac{1}{2C} \boldsymbol{\theta}^T \boldsymbol{\theta} + \frac{1}{l} \sum_{i=1}^l \xi(\boldsymbol{\theta}; \mathbf{y}^i, Z^{1,i})$$

- The gradient is

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

- Going over all data is time consuming



Estimation of the Gradient II

- What if we use a **subset** of data

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

- We may just use a subset S

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



Algorithm 1

- 1: Given an initial learning rate η .
- 2: **while do**
- 3: Choose $S \subset \{1, \dots, l\}$.
- 4: Calculate

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

- 5: May adjust the learning rate η
 - 6: **end while**
- It's known that deciding a suitable learning rate is difficult



Algorithm II

- Too small learning rate: very slow convergence
- Too large learning rate: the procedure may diverge



Stochastic Gradient “Descent” I

- In comparison with gradient descent you see that we don't do line search
- Indeed we cannot. Without the full gradient, the sufficient decrease condition may never hold.

$$f(\boldsymbol{\theta} + \alpha\Delta\boldsymbol{\theta}) < f(\boldsymbol{\theta}) + \nu\nabla f(\boldsymbol{\theta})^T(\alpha\Delta\boldsymbol{\theta})$$

- Therefore, we don't have a “descent” algorithm here
- It's possible that

$$f(\boldsymbol{\theta}^{\text{next}}) > f(\boldsymbol{\theta})$$

- Though people frequently use “SGD,” it's unclear if “D” is suitable in the name of this method



Momentum I

- This is a method to improve the convergence speed
- A new vector \mathbf{v} and a parameter $\alpha \in [0, 1)$ are introduced

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

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \mathbf{v}$$



Momentum II

- Essentially what we do is

$$\begin{aligned}\theta &\leftarrow \theta - \eta(\text{current sub-gradient}) \\ &\quad - \alpha\eta(\text{prev. sub-gradient}) \\ &\quad - \alpha^2\eta(\text{prev. prev. sub-gradient}) - \dots\end{aligned}$$

- There are some reasons why doing so can improve the convergence speed, though details are not discussed here



Outline

- 1 Gradient descent
- 2 Mini-batch SG
- 3 Adaptive learning rate**
- 4 Discussion



AdaGrad I

- Scaling learning rates inversely proportional to the square root of sum of past gradient squares (Duchi et al., 2011)
- Update rule:

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

$$\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$$

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

- \mathbf{r} : sum of past gradient squares



AdaGrad II

ϵ and δ are given constants

- \odot : Hadamard product (element-wise product of two vectors/matrices)
- A large \mathbf{g} component
 - \Rightarrow a larger \mathbf{r} component
 - \Rightarrow fast decrease of the learning rate
- Conceptual explanation from Duchi et al. (2011):
 - frequently occurring features \Rightarrow low learning rates
 - infrequent features \Rightarrow high learning rates



AdaGrad III

“the intuition is that each time an infrequent feature is seen, the learner should **take notice**.”

- But how is this explanation related to \mathbf{g} components?
- Let's consider **linear** classification. Recall our optimization problem is

$$\frac{\mathbf{w}^T \mathbf{w}}{2} + C \sum_{i=1}^l \xi(\mathbf{w}; y_i, \mathbf{x}_i)$$



AdaGrad IV

- For methods such as SVM or logistic regression, the loss function can be written as a function of $\mathbf{w}^T \mathbf{x}$

$$\xi(\mathbf{w}; y, \mathbf{x}) = \hat{\epsilon}(\mathbf{w}^T \mathbf{x})$$

Then the gradient is

$$\mathbf{w} + C \sum_{i=1}^I \hat{\epsilon}'(\mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i$$

- Thus the gradient is related to the **density of features**



AdaGrad V

- The above analysis is for linear classification
- But now we have a **non-convex** neural network!
- **Empirically**, people find that the sum of squared gradient since the beginning causes **too fast decrease of the learning rate**



RMSPprop I

- The original reference seems to be the lecture slides at https://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf
- Idea: they think AdaGrad's learning rate may be too small before reaching a locally convex region
- That is, OK to sum all past gradient squares in convex, but not non-convex
- Thus they do “**exponentially weighted moving average**”



RMSPprop II

- Update rule

$$\begin{aligned} \mathbf{r} &\leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g} \\ \boldsymbol{\theta} &\leftarrow \boldsymbol{\theta} - \frac{\epsilon}{\sqrt{\delta + \mathbf{r}}} \odot \mathbf{g} \end{aligned}$$

- AdaGrad:

$$\begin{aligned} \mathbf{r} &\leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g} \\ \boldsymbol{\theta} &\leftarrow \boldsymbol{\theta} - \frac{\epsilon}{\sqrt{\mathbf{r} + \delta}} \odot \mathbf{g} \end{aligned}$$



RMSPprop III

- Somehow the setting is a bit heuristic and the reason behind the change (from AdaGrad to RMSPprop) is not really that strong



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

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

- 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)
- We may do some experiments later



Outline

- 1 Gradient descent
- 2 Mini-batch SG
- 3 Adaptive learning rate
- 4 Discussion



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”
- This isn’t very good. Can we have some systematic investigation?



Why Stochastic Gradient Widely Used? I

- The special property of data classification is essential

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

Indeed stochastic gradient is less used outside machine learning

- Easy implementation. It's simpler than methods using, for example, second derivative
- Non-convexity plays a role



Why Stochastic Gradient Widely Used? II

- For convex, other methods are efficient to find **the global minimum**
- But for non-convex, efficiency to reach a **stationary point** is less useful
- A global minimum usually gives a good model (loss minimized), but for a stationary point we are less sure
- All these explain why SG is popular for deep learning
- What are your opinions? Any other reasons you can think of



Issues of Stochastic Gradient I

- We have shown several variants
- Don't you think some settings are a bit ad hoc?
- There are reasons behind each change. But some are just heuristic
- Can we try a paradigm completely different?
- But before that we need some first-hand experiences and know implementation details.



References I

- J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12:2121–2159, 2011.
- 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.
- J. Nocedal and S. J. Wright. *Numerical Optimization*. Springer-Verlag, New York, NY, 1999.
- 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.

