

Stochastic Gradient Descent

Zou Lexiao

Harbin Institute of Technology, Shenzhen

January 2023

Table of Contents

1 Choosing the step size (learning rate)

2 Iterate averaging

3 Variance reduction

- SVRG
- SAGA
- Application to deep learning

4 Preconditioned SGD

- AdaGrad
- RMSprop and AdaDelta
- Adam
- Issues with adaptive learning rates
- Non-diagonal preconditioning matrices

Motivation

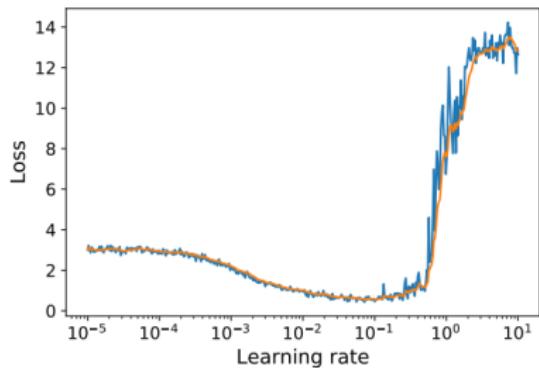


Figure 8.17: Loss vs learning rate (horizontal axis). Training loss vs learning rate for a small MLP fit to FashionMNIST using vanilla SGD. (Raw loss in blue, EWMA smoothed version in orange). Generated by `lrschedule_tf.ipynb`.

One heuristic

- Start with a small learning rate and gradually increase it, evaluating performance using a small number of minibatches.
- In practice, it is better to pick a rate that is slightly smaller than (i.e., to the left of) the one with the lowest loss, to ensure stability.

Learning rate schedule

- Definition:

Learning rate schedule: Adjust the step size over time.

- Sufficient Condition for convergence:

Robbins-Monro conditions¹:

$$\eta_t \rightarrow 0, \frac{\sum_{t=1}^{\infty} \eta_t^2}{\sum_{t=1}^{\infty} \eta_t} \rightarrow 0$$

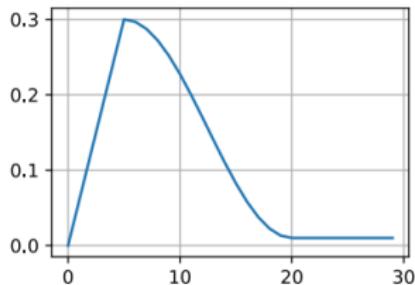
- Examples

- piecewise constant $\eta_t = \eta_i$ if $t_i \leq t \leq t_{i+1}$
 - Step decay
 - Reduce-on-plateau
- exponential decay $\eta_t = \eta_0 e^{-\lambda t}$
- polynomial decay $\eta_t = \eta_0 (\beta t + 1)^{-\alpha}$
 - square-root schedule: $\alpha = 0.5, \beta = 1$

¹Herbert Robbins and Sutton Monro. "A Stochastic Approximation Method". In: *The Annals of Mathematical Statistics* 22.3 (1951), pp. 400–407. DOI: 10.1214/aoms/1177729586. URL: <https://doi.org/10.1214/aoms/1177729586>.

Learning rate schedule: In deep learning community

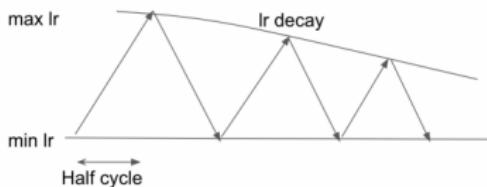
Learning rate warmup(one-cycle learning rate schedule):



(a)

- Notice:
To ensure convergence to a point, we must reduce the learning rate to 0.

Cyclical learning rate:



- Motivation: To escape local minima.
- Methods: The minimum and maximum learning rates can be found based on the initial “dry run” described above, and the half-cycle can be chosen based on how many restarts you want to do with your training budget.
- Related works: **stochastic gradient descent with warm restarts(SGDR)**
storing all the checkpoints visited after each cool down, and using all of them as members of a **model ensemble**.

Table of Contents

1 Choosing the step size (learning rate)

2 Iterate averaging

3 Variance reduction

- SVRG
- SAGA
- Application to deep learning

4 Preconditioned SGD

- AdaGrad
- RMSprop and AdaDelta
- Adam
- Issues with adaptive learning rates
- Non-diagonal preconditioning matrices

Iterate averaging

- Motivation:

The parameter estimates produced by SGD can be very unstable over time.

- Methods:

$$\bar{\theta}_t = \frac{1}{t} \sum_{i=1}^t \theta_i = \frac{1}{t} \theta_t + \frac{t-1}{t} \bar{\theta}_{t-1}$$

- Advantages:

- Achieves the best possible asymptotic convergence rate among SGD algorithms, matching that of variants using second-order information².
- Statistical benefits: in the case of linear regression, this method is equivalent to ℓ_2 regularization (i.e., ridge regression)

- Related work: **Stochastic Weight Averaging (SWA)**³

SWA exploits the flatness in objectives used to train deep neural networks, to find solutions which provide better generalization.

²Boris Polyak and Anatoli Juditsky. "Acceleration of Stochastic Approximation by Averaging". In: *SIAM Journal on Control and Optimization* 30 (July 1992), pp. 838–855. DOI: [10.1137/0330046](https://doi.org/10.1137/0330046).

³Pavel Izmailov et al. "Averaging Weights Leads to Wider Optima and Better Generalization". In: *arXiv e-prints*, arXiv:1803.05407 (Mar. 2018), arXiv:1803.05407. arXiv: 1803.05407 [cs.LG].

Table of Contents

- 1 Choosing the step size (learning rate)
- 2 Iterate averaging
- 3 Variance reduction
 - SVRG
 - SAGA
 - Application to deep learning
- 4 Preconditioned SGD
 - AdaGrad
 - RMSprop and AdaDelta
 - Adam
 - Issues with adaptive learning rates
 - Non-diagonal preconditioning matrices

stochastic variance reduced gradient (SVRG)

Ever so often (e.g., once per epoch), we compute the full gradient at a "snapshot" of the model parameters $\tilde{\theta}$; the corresponding "exact" gradient is therefore $\nabla \mathcal{L}(\tilde{\theta})$. At step t , we compute the usual stochastic gradient at the current parameters, $\nabla \mathcal{L}_t(\theta_t)$, but also at the snapshot parameters, $\nabla \mathcal{L}_t(\tilde{\theta})$, which we use as a baseline. We can then use the following improved gradient estimate

$$\mathbf{g}_t = \nabla \mathcal{L}_t(\theta_t) - \nabla \mathcal{L}_t(\tilde{\theta}) + \nabla \mathcal{L}(\tilde{\theta})$$

stochastic averaged gradient accelerated (SAGA)

We first initialize by computing $\mathbf{g}_n^{\text{local}} = \nabla \mathcal{L}_n(\boldsymbol{\theta}_0)$ for all n , and the average, $\mathbf{g}^{\text{avg}} = \frac{1}{N} \sum_{n=1}^N \mathbf{g}_n^{\text{local}}$. Then, at iteration t , we use the gradient estimate

$$\mathbf{g}_t = \nabla \mathcal{L}_n(\boldsymbol{\theta}_t) - \mathbf{g}_n^{\text{local}} + \mathbf{g}^{\text{avg}}$$

where $n \sim \text{Unif}\{1, \dots, N\}$ is the example index sampled at iteration t . We then update $\mathbf{g}_n^{\text{local}} = \nabla \mathcal{L}_n(\boldsymbol{\theta}_t)$ and \mathbf{g}^{avg} by replacing the old $\mathbf{g}_n^{\text{local}}$ by its new value.

- Unlike SVRG, it only requires one full batch gradient computation, at the start of the algorithm. While the downside is the extra memory cost.
- It is recommended for use in the sklearn logistic regression code when N is large and \mathbf{x} is sparse.

Application to deep learning

For example, the use of batch normalization , data augmentation and dropout all break the assumptions of SVRG, since the loss will differ randomly in ways that depend not just on the parameters and the data index n .

Table of Contents

1 Choosing the step size (learning rate)

2 Iterate averaging

3 Variance reduction

- SVRG
- SAGA
- Application to deep learning

4 Preconditioned SGD

- AdaGrad
- RMSprop and AdaDelta
- Adam
- Issues with adaptive learning rates
- Non-diagonal preconditioning matrices

Preconditioned SGD

$$\theta_{t+1} = \theta_t - \eta_t \mathbf{M}_t^{-1} \mathbf{g}_t$$

- The noise in the gradient estimates make it difficult to reliably estimate the Hessian.
- Most practitioners use a diagonal preconditioner

AdaGrad

$$\theta_{t+1,d} = \theta_{t,d} - \eta_t \frac{1}{\sqrt{s_{t,d} + \epsilon}} g_{t,d}$$

where $d = 1 : D$ indexes the dimensions of the parameter vector, and

$$s_{t,d} = \sum_{i=1}^t g_{i,d}^2$$

is the sum of the squared gradients and $\epsilon > 0$ is a small term to avoid dividing by zero.

RMSprop and AdaDelta

Motivation: the denominator get larger over time. However, it might hurt performance as the denominator gets large too fast.

RMSprop

$$s_{t+1,d} = \beta s_{t,d} + (1 - \beta) g_{t,d}^2$$

- We apply the EWMA of the past squared gradients rather than their sum.
- In practice we usually use $\beta \sim 0.9$, which puts more weight on recent examples.

RMSprop and AdaDelta

AdaDelta

$$\Delta\theta_t = -\eta_t \frac{\sqrt{\delta_{t-1} + \epsilon}}{\sqrt{s_t + \epsilon}} g_t$$

where

$$\delta_t = \beta\delta_{t-1} + (1 - \beta)(\Delta\theta_t)^2$$

- Eliminate the need to tune learning rate η_t , since the “units” of the numerator and denominator cancel, so we are just elementwise-multiplying the gradient by a scalar.⁴

However, since these adaptive learning rates need not decrease with time (unless we choose η_t to explicitly do so), these methods are not guaranteed to converge to a solution.

⁴Matthew Zeiler. “ADADELTA: An adaptive learning rate method”. In: 1212 (Dec 2012). A set of small, light-blue navigation icons typically used in Beamer presentations for navigating between slides and sections.

Adam=RMSProp + momentum

$$\begin{aligned}\mathbf{m}_t &= \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t \\ \mathbf{s}_t &= \beta_2 \mathbf{s}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2\end{aligned}$$

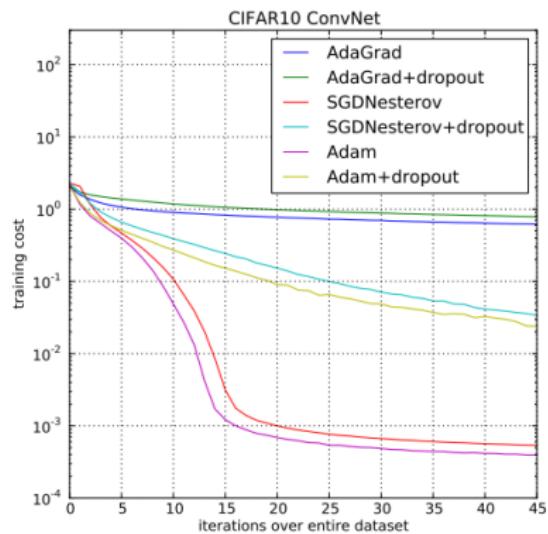
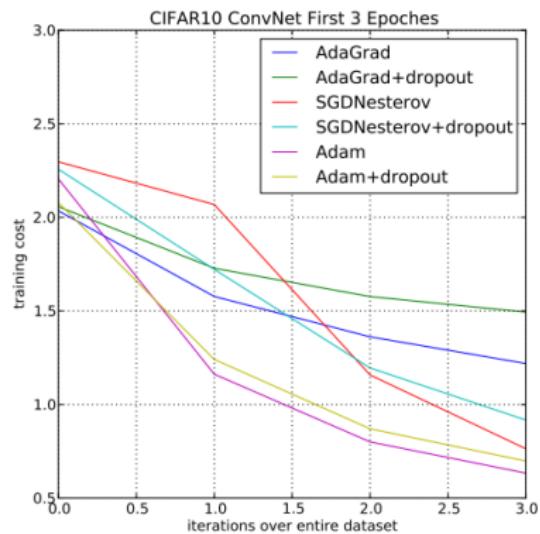
We then perform the update

$$\Delta \theta_t = -\eta_t \frac{1}{\sqrt{\mathbf{s}_t} + \epsilon} \mathbf{m}_t$$

If we initialize with $\mathbf{m}_0 = \mathbf{s}_0 = \mathbf{0}$, then initial estimates will be biased towards small values. The authors therefore recommend using the bias-corrected moments, which increase the values early in the optimization process. These estimates are given by

$$\begin{aligned}\hat{\mathbf{m}}_t &= \mathbf{m}_t / (1 - \beta_1^t) \\ \hat{\mathbf{s}}_t &= \mathbf{s}_t / (1 - \beta_2^t)\end{aligned}$$

Adam



Issues with adaptive learning rates

- Although it is called adaptive learning rate methods, it still needs to set the base learning rate η_0
- Since the EWMA methods are typically used in the stochastic setting where the gradient estimates are noisy, their learning rate adaptation can result in non-convergence even on convex problems.

Non-diagonal preconditioning matrices

Motivation: Although the methods we have discussed above can adapt the learning rate of each parameter, they do not solve the more fundamental problem of ill-conditioning due to correlation of the parameters, and hence do not always provide as much of a speed boost over vanilla SGD as one may hope.⁵

⁵Rohan Anil et al. "Scalable Second Order Optimization for Deep Learning". In: (Mar 2021).

References

- [1] Herbert Robbins and Sutton Monro. "A Stochastic Approximation Method". In: *The Annals of Mathematical Statistics* 22.3 (1951), pp. 400–407. DOI: 10.1214/aoms/1177729586. URL: <https://doi.org/10.1214/aoms/1177729586>.
- [2] Boris Polyak and Anatoli Juditsky. "Acceleration of Stochastic Approximation by Averaging". In: *SIAM Journal on Control and Optimization* 30 (July 1992), pp. 838–855. DOI: 10.1137/0330046.
- [3] Pavel Izmailov et al. "Averaging Weights Leads to Wider Optima and Better Generalization". In: *arXiv e-prints*, arXiv:1803.05407 (Mar. 2018), arXiv:1803.05407. arXiv: 1803.05407 [cs.LG].
- [4] Matthew Zeiler. "ADADELTA: An adaptive learning rate method". In: 1212 (Dec. 2012).
- [5] Rohan Anil et al. "Scalable Second Order Optimization for Deep Learning". In: (Mar. 2021).