

# Lecture 9. Stochastic Gradient Descent

Bao Wang

Department of Mathematics

Scientific Computing and Imaging Institute

University of Utah

Math 5750/6880, Fall 2021

## Loss function

So far, we have formulated training machine learning models as

$$\min f(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}_i(\mathbf{x}) + R(\mathbf{x})$$

where  $\mathbf{x}$  is the parameter of the machine learning model,  $\mathcal{L}_i(\mathbf{x})$  is the loss of the  $i$ th training instance, and  $R(\mathbf{x})$  is the regularization term.

How to find the optimal  $\mathbf{x}^*$  if  $N$  is very large? Compute the gradient of  $f$  is very difficult.

Let us first ignore the regularization term. (Using proximal!)

## Empirical risk minimization

Let  $\{\mathbf{a}_i, y_i\}_{i=1}^n$  be  $n$  random samples, and consider

$$\min_{\mathbf{x}} F(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}; \{\mathbf{a}_i, y_i\}),$$

e.g. quadratic loss  $f(\mathbf{x}; \{\mathbf{a}_i, y_i\}) = (\mathbf{a}_i^\top \mathbf{x} - y_i)^2$ . If one draws index  $j \sim \text{Unif}(1, \dots, n)$  uniformly at random, then

$$F(\mathbf{x}) = \mathbb{E}_j[f(\mathbf{x}; \{\mathbf{a}_j, y_j\})].$$

More generally, we consider the following stochastic programming

$$\min_{\mathbf{x}} F(\mathbf{x}) = \mathbb{E}[f(\mathbf{x}; \xi)]$$

- >  $\xi$ : randomness in problem
- > suppose  $f(\cdot, \xi)$  is convex for every  $\xi$  (and hence  $F(\cdot)$  is convex)

## A natural solution

Under "mild" technical conditions, we have

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \nabla F(\mathbf{x}^t) = \mathbf{x}^t - \eta_t \nabla \mathbb{E}[f(\mathbf{x}^t; \xi)] = \mathbf{x}^t - \eta_t \mathbb{E}[\nabla_{\mathbf{x}} f(\mathbf{x}^t; \xi)]$$

Issues:

- > distribution of  $\xi$  maybe unknown.
- > even if it is known, evaluating high-dimensional expectation is often expensive.

## Stochastic gradient descent (SGD)

Stochastic approximation/stochastic gradient descent (SGD):

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \mathbf{g}(\mathbf{x}^t; \xi^t) \quad (1)$$

where  $\mathbf{g}(\mathbf{x}^t; \xi^t)$  is unbiased estimate of  $\nabla F(\mathbf{x}^t)$ , i.e.

$$\mathbb{E}[\mathbf{g}(\mathbf{x}^t; \xi^t)] = \nabla F(\mathbf{x}^t).$$

SGD is a stochastic algorithm for finding a critical point  $\mathbf{x}$  obeying  $\nabla F(\mathbf{x}) = 0$ .

## Stochastic gradient descent (SGD)

**Example.** Consider the empirical risk minimization problem

$$\min_{\mathbf{x}} F(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}; \{\mathbf{a}_i, y_i\})$$

SGD for empirical risk minimization

for  $t = 0, 1, \dots$

choose  $i_t$  uniformly at random, and run

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \nabla_{\mathbf{x}} f_{i_t}(\mathbf{x}^t; \{\mathbf{a}_i, y_i\})$$

## SGD for empirical risk minimization

**Benefits:** SGD exploits information more efficiently than batch methods

1. practical data usually involve lots of redundancy; using all data simultaneously in each iteration might be inefficient.
2. SGD is particularly efficient at the very beginning, as it achieves fast initial improvement with very low per-iteration cost.

## SGD for empirical risk minimization

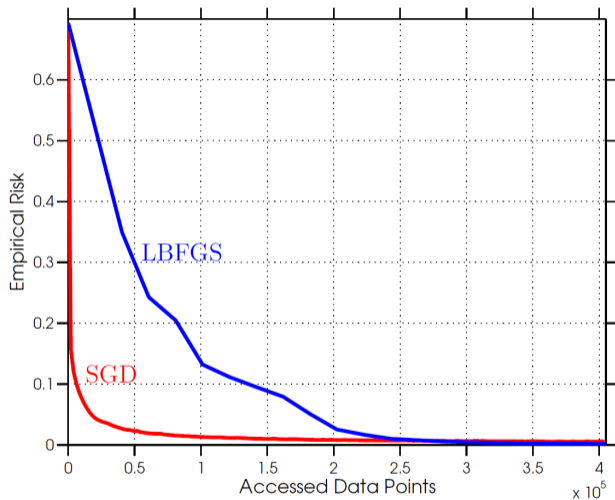


Figure: Logistic regression for RCV1 dataset ( $\eta_t \equiv 4$ ).



Convergence analysis

## Strongly convex and smooth problems

**Assumptions.** Consider

$$\min_{\mathbf{x}} F(\mathbf{x}) := \mathbb{E}[f(\mathbf{x}; \xi)]$$

where

1.  $F$ :  $\mu$ -strongly convex,  $L$ -smooth
2.  $g(\mathbf{x}^t; \xi^t)$ : an **unbiased** estimate of  $\nabla F(\mathbf{x}^t)$  given  $\{\xi^0, \dots, \xi^{t-1}\}$
3. for all  $\mathbf{x}$ ,

$$\mathbb{E}[\|g(\mathbf{x}; \xi)\|_2^2] \leq \sigma_g^2 + c_g \|\nabla F(\mathbf{x})\|_2^2. \quad (2)$$

## Convergence: fixed stepsizes

**Theorem.** [Convergence of SGD for strongly convex problems: fixed stepsizes] Under the assumptions above, if  $\eta_t \equiv \eta \leq \frac{1}{Lc_g}$ , then SGD (1) achieves

$$\mathbb{E}[F(\mathbf{x}^t) - F(\mathbf{x}^*)] \leq \frac{\eta L \sigma_g^2}{2\mu} + (1 - \eta\mu)^t (F(\mathbf{x}^0) - F(\mathbf{x}^*)).$$

**Proof.** See L. Bottou, F. Curtis, and J. Nocedal. "Optimization Methods for Large-Scale Machine Learning", SIAM Rev, 2018.

## Convergence: fixed stepsizes

### Remarks.

1. fast (linear) convergence at the very beginning
2. converges to some neighborhood of  $\mathbf{x}^*$  — variance in gradient computation prevents further progress
3. when gradient computation is noiseless (i.e.  $\sigma_g = 0$ ), it converges linearly to optimal points
4. smaller stepsizes  $\eta$  yield better converging points

## One practical strategy

Run SGD with fixed stepsizes; whenever progress stalls, reduce stepsizes and continue SGD.

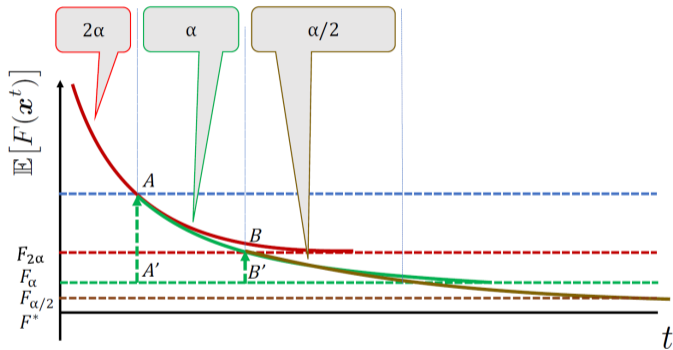


Figure: Whenever progress stalls, we half the stepsizes and repeat.

Stagewise stepsize.

## Convergence with diminishing stepsizes

**Theorem.** [Convergence of SGD for strongly convex problems: diminishing stepsizes]  
Suppose  $F$  is  $\mu$ -strongly convex, and (2) holds with  $c_g = 0$ . If  $\eta_t = \frac{\theta}{t+1}$  for some  $\theta > \frac{1}{2\mu}$ , then SGD (1) achieves

$$\mathbb{E}[\|\mathbf{x}^t - \mathbf{x}^*\|_2^2] \leq \frac{c_\theta}{t+1}$$

where  $c_\theta = \max \left\{ \frac{2\theta^2\sigma_g^2}{2\mu\theta-1}, \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 \right\}$ .

**Remark.** convergence rate  $O(1/t)$  with diminishing stepsize  $\eta_t \sim 1/t$ .

**Proof.** Using the SGD update rule, we have Using the SGD update rule, we have

$$\begin{aligned}\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2 &= \|\mathbf{x}^t - \eta_t \mathbf{g}(\mathbf{x}^t; \xi^t) - \mathbf{x}^*\|_2^2 \\ &= \|\mathbf{x}^t - \mathbf{x}^*\|_2^2 - 2\eta_t (\mathbf{x}^t - \mathbf{x}^*)^\top \mathbf{g}(\mathbf{x}^t; \xi^t) + \eta_t^2 \|\mathbf{g}(\mathbf{x}^t; \xi^t)\|_2^2.\end{aligned}\tag{3}$$

Since  $\mathbf{x}^t$  is independent of  $\xi_t$ , apply the law of total expectation to obtain

$$\begin{aligned}\mathbb{E}[(\mathbf{x}^t - \mathbf{x}^*)^\top \mathbf{g}(\mathbf{x}^t; \xi^t)] &= \mathbb{E}[\mathbb{E}[(\mathbf{x}^t - \mathbf{x}^*)^\top \mathbf{g}(\mathbf{x}^t; \xi^t) | \xi_1, \dots, \xi_{t-1}]] \\ &= \mathbb{E}[(\mathbf{x}^t - \mathbf{x}^*)^\top \mathbb{E}[\mathbf{g}(\mathbf{x}^t; \xi^t) | \xi_1, \dots, \xi_{t-1}]] \\ &= \mathbb{E}[(\mathbf{x}^t - \mathbf{x}^*)^\top \nabla F(\mathbf{x}^t)]\end{aligned}\tag{4}$$

Furthermore, strong convexity gives

$$\begin{aligned}\langle \nabla F(\mathbf{x}^t), \mathbf{x}^t - \mathbf{x}^* \rangle &= \langle \nabla F(\mathbf{x}^t) - \underbrace{\nabla F(\mathbf{x}^*)}_{=0}, \mathbf{x}^t - \mathbf{x}^* \rangle \geq \mu \|\mathbf{x}^t - \mathbf{x}^*\|_2^2 \\ \Rightarrow \mathbb{E}[\langle \nabla F(\mathbf{x}^t), \mathbf{x}^t - \mathbf{x}^* \rangle] &\geq \mu \mathbb{E}[\|\mathbf{x}^t - \mathbf{x}^*\|_2^2].\end{aligned}\tag{5}$$

Combine (3), (4), (5) and (2) (with  $c_g = 0$ ) to obtain

$$\mathbb{E}[\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2] \leq (1 - 2\mu\eta_t)\mathbb{E}[\|\mathbf{x}^t - \mathbf{x}^*\|_2^2] + \underbrace{\eta_t^2 \sigma_g^2}_{\text{does not vanish unless } \eta_t \rightarrow 0}.\tag{6}$$

Take  $\eta_t = \frac{\theta}{t+1}$  and use induction to conclude the proof.



## Optimality

Informally, when minimizing strongly convex functions, no algorithm performing  $t$  queries to noisy first-order oracles can achieve an accuracy better than the order of  $1/t$   
 $\Rightarrow$  SGD with step sizes  $\eta_t \sim 1/t$  is optimal.

## Optimality

[Nemirovski & Yudin 1983] More precisely, consider a class of problems in which  $f$  is  $\mu$ -strongly convex and  $L$ -smooth, and  $\text{Var}(\|\mathbf{g}(\mathbf{x}^t; \xi^t)\|_2) \leq \sigma^2$ . Then the worst-case iteration complexity for (stochastic) first-order methods:

$$\sqrt{\frac{L}{\mu}} \log\left(\frac{L\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{\epsilon}\right) + \frac{\sigma^2}{\mu\epsilon}.$$

> For deterministic case:  $\sigma = 0$ , and hence the lower bound is

$$\sqrt{\frac{L}{\mu}} \log\left(\frac{L\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{\epsilon}\right) \quad (\text{achievable by Nesterov's method})$$

> For noisy case with large  $\sigma$ , the lower bound is dominated by (1/t convergence rate)

$$\frac{\sigma^2}{\mu} \cdot \frac{1}{\epsilon}$$

## Recap

Nesterov's method for strongly convex problems:

$$\begin{aligned}\mathbf{x}^{t+1} &= \text{prox}_{\eta_t h}(\mathbf{y}^t - \eta_t \nabla f(\mathbf{x}^t)) \\ \mathbf{y}^{t+1} &= \mathbf{x}^{t+1} + \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}(\mathbf{x}^{t+1} - \mathbf{x}^t)\end{aligned}$$

**Theorem 4.** [Convergence of accelerated proximal gradient methods for strongly convex case] Suppose  $f$  is  $\mu$ -strongly convex and  $L$ -smooth. If  $\eta_t \equiv 1/L$ , then

$$F(\mathbf{x}^t) - F^{opt} \leq \left(1 - \frac{1}{\sqrt{\kappa}}\right)^t \left(F(\mathbf{x}^0) - F^{opt} + \frac{\mu \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2}{2}\right).$$

## Comparison with batch GD

	iteration complexity	per-iteration cost	total comput. cost
batch GD	$\log \frac{1}{\epsilon}$	$n$	$n \log \frac{1}{\epsilon}$
SGD	$\frac{1}{\epsilon}$	1	$\frac{1}{\epsilon}$

**Figure:** Empirical risk minimization with  $n$  samples.

SGD is more appealing for large  $n$  and moderate accuracy  $\epsilon$  (in which case  $\frac{1}{\epsilon} < n \log \frac{1}{\epsilon}$ ) – which often arise in the big data regime! ( $n$  is very big.)

## Convex problems

What if we lose strong convexity?

$$\min_{\mathbf{x}} F(\mathbf{x}) := \mathbb{E}[f(\mathbf{x}; \xi)]$$

1.  $F$ : convex
2.  $\mathbb{E}[\|\mathbf{g}(\mathbf{x}; \xi)\|_2^2] \leq \sigma_g^2$  for all  $\mathbf{x}$
3.  $\mathbf{g}(\mathbf{x}^t; \xi^t)$  is an unbiased estimate of  $\nabla F(\mathbf{x}^t)$  given  $\{\xi^0, \dots, \xi^{t-1}\}$ .

Suppose we run SGD and return a weighted average

$$\tilde{\mathbf{x}}^t := \sum_{k=0}^t \frac{\eta_k}{\sum_{j=0}^t \eta_j} \mathbf{x}^k.$$

**Theorem.** Under the assumption above, one has

$$\mathbb{E}[F(\tilde{\mathbf{x}}^t) - F(\mathbf{x}^*)] \leq \frac{\frac{1}{2}\mathbb{E}[\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2] + \frac{1}{2}\sigma_g^2 \sum_{k=0}^t \eta_k^2}{\sum_{k=0}^t \eta_k}.$$

In particular, if  $\eta_t \sim 1/\sqrt{t}$ , then

$$\mathbb{E}[F(\tilde{\mathbf{x}}^t) - F(\mathbf{x}^*)] \lesssim \frac{\log t}{\sqrt{t}}.$$

GD for cvx:  $\eta_t = 1/L$  and convergence rate is  $1/t$ .

**Proof.** By convexity of  $F$ , we have

$$F(\mathbf{x}) \geq F(\mathbf{x}^t) + (\mathbf{x} - \mathbf{x}^t)^\top \nabla F(\mathbf{x}^t) \Rightarrow \mathbb{E}[(\mathbf{x}^t - \mathbf{x}^*)^\top \nabla F(\mathbf{x}^t)] \geq \mathbb{E}[F(\mathbf{x}^t) - F(\mathbf{x}^*)].$$

This together with (3) and (4) implies

$$2\eta_k \mathbb{E}[F(\mathbf{x}^k) - F(\mathbf{x}^*)] \leq \mathbb{E}[\|\mathbf{x}^k - \mathbf{x}^*\|_2^2] - \mathbb{E}[\|\mathbf{x}^{k+1} - \mathbf{x}^*\|_2^2] + \eta_k^2 \sigma_g^2.$$

Sum over  $k = 0, \dots, t$  to obtain

$$\begin{aligned} \sum_{k=0}^t 2\eta_k \mathbb{E}[F(\mathbf{x}^k) - F(\mathbf{x}^*)] &\leq \mathbb{E}[\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2] - \mathbb{E}[\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2] + \sigma_g^2 \sum_{k=0}^t \eta_k^2 \\ &\leq \mathbb{E}[\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2] + \sigma_g^2 \sum_{k=0}^t \eta_k^2. \end{aligned}$$

Setting  $v_t = \frac{\eta_t}{\sum_{k=0}^t \eta_k}$  yields

$$\sum_{k=0}^t v_k \mathbb{E}[F(\mathbf{x}^k) - F(\mathbf{x}^*)] \leq \frac{\frac{1}{2} \mathbb{E}[\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2] + \frac{1}{2} \sigma_g^2 \sum_{k=0}^t \eta_k^2}{\sum_{k=0}^t \eta_k}.$$

By convexity of  $F$ , we arrive at

$$\mathbb{E}[F(\tilde{\mathbf{x}}^t) - F(\mathbf{x}^*)] \leq \frac{\frac{1}{2} \mathbb{E}[\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2] + \frac{1}{2} \sigma_g^2 \sum_{k=0}^t \eta_k^2}{\sum_{k=0}^t \eta_k}.$$



Reducing variance via iterate averaging

Stepsize choice  $O(1/t)$ ?

Two conflicting regimes:

1. The noiseless case (i.e.  $\mathbf{g}(\mathbf{x}; \xi) = \nabla F(\mathbf{x})$ ): stepsizes  $\eta_t \sim 1/t$  are way too conservative.
2. The general noisy case: longer stepsizes ( $\eta_t \gg 1/t$ ) might fail to suppress noise (and hence slow down convergence).

Can we modify SGD so as to allow for larger stepsizes without compromising convergence rates?

## Acceleration by averaging

**Motivation for iteration averaging.** SGD with long stepsizes poorly suppresses noise, which tends to oscillate around the global minimizers due to the noisy nature of gradient computation. One may, however, average iterates to mitigate oscillation and reduce variance.

Run SGD and return

$$\bar{\mathbf{x}}^t := \frac{1}{t} \sum_{i=0}^{t-1} \mathbf{x}^i, \quad (7)$$

with larger stepsizes  $\eta_t \sim t^{-\alpha}$  with  $\alpha < 1$ .

**Key idea:** average the iterates to reduce variance and improve convergence.

## Example: a toy quadratic problem

Consider  $\min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{x}\|_2^2$  using constant stepsizes  $\eta_t \equiv \eta < 1$ . Assume the unbiased stochastic gradient satisfies

$$\mathbf{g}(\mathbf{x}^t; \xi^t) = \mathbf{x}^t + \xi^t \quad \text{with} \quad \mathbb{E}[\xi^t | \xi^0, \dots, \xi^{t-1}] = 0 \quad \text{and} \quad \mathbb{E}[\xi^t \xi^{t\top} | \xi^0, \dots, \xi^{t-1}] = \mathbf{I}.$$

SGD iterates:

$$\mathbf{x}^1 = \mathbf{x}^0 - \eta(\mathbf{x}^0 + \xi^0) = (1 - \eta)\mathbf{x}^0 - \eta\xi^0$$

$$\mathbf{x}^2 = \mathbf{x}^1 - \eta(\mathbf{x}^1 + \xi^1) = (1 - \eta)^2\mathbf{x}^0 - \eta(1 - \eta)\xi^0 - \eta\xi^1$$

$\vdots$

$$\mathbf{x}^t = (1 - \eta)^t\mathbf{x}^0 - \eta(1 - \eta)^{t-1}\xi^0 - \eta(1 - \eta)^{t-2}\xi^1 - \dots$$

## Example: a toy quadratic problem

i.e.,

$$\begin{aligned}\bar{\mathbf{x}}^t &\approx \underbrace{\frac{1}{t} \sum_{k=0}^{t-1} (1-\eta)^k \mathbf{x}^0}_{=\frac{1}{t} \frac{1-(1-\eta)^t}{\eta} \mathbf{x}^0 \rightarrow 0} - \underbrace{\eta \{1 + (1-\eta) + \dots\}}_{\text{imprecise: but close enough for large } t} \frac{1}{t} \sum_{k=0}^{t-1} \xi^k \\ &\approx -\frac{1}{t} \sum_{k=0}^{t-1} \xi^k \quad (\text{since } 1 + (1-\eta) + \dots = \eta^{-1}) \\ &\rightarrow \frac{1}{\sqrt{t}} \mathcal{N}(0, I) \text{ as } t \rightarrow \infty \quad (\text{the central limit theorem for martingale})\end{aligned}$$

Other popular SGD variants: SGD with momentum

$$\begin{aligned}\mathbf{v}^{t+1} &= \gamma \mathbf{v}_t + \eta \mathbf{g}(\mathbf{x}^t; \xi^t) \\ \mathbf{x}^{t+1} &= \mathbf{x}^t - \mathbf{v}^{t+1}.\end{aligned}\tag{8}$$

How it is related to the heavy-ball method?

From the second equation, we have  $\mathbf{v}^{t+1} = \mathbf{x}^t - \mathbf{x}^{t+1}$ , and therefore we have

$$\mathbf{x}^t - \mathbf{x}^{t+1} = \gamma(\mathbf{x}^{t-1} - \mathbf{x}^t) + \eta \mathbf{g}(\mathbf{x}^t; \xi^t),$$

i.e.

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \mathbf{g}(\mathbf{x}^t; \xi^t) + \gamma(\mathbf{x}^t - \mathbf{x}^{t-1}).$$

Other popular SGD variants: SGD with Nesterov momentum

$$\begin{aligned}\mathbf{v}^{t+1} &= \gamma \mathbf{v}_t + \eta \mathbf{g}(\mathbf{x}^t - \gamma \mathbf{v}_t; \xi^t) \\ \mathbf{x}^{t+1} &= \mathbf{x}^t - \mathbf{v}^{t+1}.\end{aligned}\tag{9}$$

In practice, SGD with Nesterov momentum performs similar to SGD with momentum.

## Other popular SGD variants: Adagrad

Let  $\mathbf{g}_i^t = \mathbf{g}(x_i^t)$ , SGD update for every parameter  $x_i$  at each time step  $t$  then becomes

$$x_i^{t+1} = x_i^t - \eta g_i^t.$$

Adagrad modifies the general learning rate  $\eta$  at each time step  $t$  for every parameter  $x_i$  based on the past gradients that have been computed for  $x_i$ :

$$x_i^{t+1} = x_i^t - \frac{\eta}{\sqrt{G_{ii}^t + \epsilon}} \cdot g_i^t,$$

where  $G^t \in \mathbb{R}^{d \times d}$  is a diagonal matrix where each diagonal element  $i, i$  is the sum of the squares of the gradients w.r.t.  $x_i$  up to time step  $t$ , while  $\epsilon$  (e.g.  $10^{-8}$ ) is a smoothing term that avoids division by zero. In vector form, we have

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \frac{\eta}{\sqrt{G^t + \epsilon}} \odot \mathbf{g}^t.$$

Adagrad eliminates the tuning of learning rate, but  $G^t + \epsilon$  can go to infinity.



## Other popular SGD variants: Adadelta or RMSprop

Instead of taking the sum of the squares of the gradients w.r.t.  $x_i$  up to time step  $t$ , Adadelta or RMSprop takes the following moving average

$$E[\mathbf{g}^t]^2 = \gamma E[\mathbf{g}^{t-1}]^2 + (1 - \gamma)[\mathbf{g}^t]^2,$$

where  $0 < \gamma < 1$  is a constant, and then we update  $\mathbf{x}^t$  as follows:

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \frac{\eta}{\sqrt{E[\mathbf{g}^t]^2 + \epsilon}} \mathbf{g}^t.$$

## Other popular SGD variants: Adam

Adam further integrates momentum into RMSprop. In particular, we compute the decaying averages of past and past squared gradients  $\mathbf{m}^t$  and  $\mathbf{v}^t$  as follows:

$$\mathbf{m}^t = \beta_1 \mathbf{m}^{t-1} + (1 - \beta_1) \mathbf{g}^t; \quad \mathbf{v}^t = \beta_2 \mathbf{v}^{t-1} + (1 - \beta_2) [\mathbf{g}^t]^2,$$

where we typically set  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . Next, we do the following bias-correction

$$\hat{\mathbf{m}}^t = \frac{\mathbf{m}^t}{1 - \beta_1^t}; \quad \hat{\mathbf{v}}^t = \frac{\mathbf{v}^t}{1 - \beta_2^t}.$$

Finally, we update the parameter as follows:

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \frac{\eta}{\sqrt{\hat{\mathbf{v}}^t} + \epsilon} \hat{\mathbf{m}}^t.$$

## Batch Normalization

To facilitate learning, we typically normalize the initial values of our parameters by initializing them with zero mean and unit variance. As training progresses and we update parameters to different extents, we lose this normalization, which slows down training and amplifies changes as the network becomes deeper.

Batch normalization reestablishes these normalizations for every mini-batch and changes are back-propagated through the operation as well. By making normalization part of the model architecture, we are able to use higher learning rates and pay less attention to the initialization parameters. Batch normalization additionally acts as a regularizer.

## Early Stopping

Monitor error on a validation set during training and stop (with some patience) if your validation error does not improve enough.

