

# Mathematics of deep learning

## Guest Lecture

### Leon Bottou 02.06.2018

Concatenated slides corresponding to:

- <http://leon.bottou.org/papers/tr-optml-2016>
- <http://leon.bottou.org/papers/lafond-vasilache-bottou-2017>

# Optimization Methods for Machine Learning

## Part II – The theory of SG

Leon Bottou

*Facebook AI Research*



Frank E. Curtis

*Lehigh University*



Jorge Nocedal

*Northwestern University*



See <http://arxiv.org/abs/1606.04838>

# Summary

1. Setup
2. Fundamental Lemmas
3. SG for Strongly Convex Objectives
4. SG for General Objectives
5. Work complexity for Large-Scale Learning
6. Comments

# 1- Setup

# The generic SG algorithm

The SG algorithm produces successive iterates  $w_k \in \mathbb{R}^d$  with the goal to minimize a certain function  $F : \mathbb{R}^d \rightarrow \mathbb{R}$ .

We assume that we have access to three mechanisms

1. Given an iteration number  $k$  ,  
a mechanism to generate a realization of a random variable  $\xi_k$ .  
The  $\{\xi_k\}$  form a sequence of jointly independent random variables
2. Given an iterate  $w_k$  and a realization  $\xi_k$ ,  
a mechanism to compute a stochastic vector  $g(w_k, \xi_k) \in \mathbb{R}^d$
3. Given an iteration number,  
a mechanism to compute a scalar stepsize  $\alpha_k > 0$

# The generic SG algorithm

## Algorithm 4.1 (Stochastic Gradient (SG) Method)

- 1: Choose an initial iterate  $w_1$ .
- 2: **for**  $k = 1, 2, \dots$  **do**
- 3:     Generate a realization of the random variable  $\xi_k$ .
- 4:     Compute a stochastic vector  $g(w_k, \xi_k)$ .
- 5:     Choose a stepsize  $\alpha_k > 0$ .
- 6:     Set the new iterate as  $w_{k+1} \leftarrow w_k - \alpha_k g(w_k, \xi_k)$ .
- 7: **end for**

# The generic SG algorithm

The function  $F : \mathbb{R}^d \rightarrow \mathbb{R}$  could be

$$F(w) = \begin{cases} R(w) = \mathbb{E}[f(w; \xi)] & \text{the expected risk,} \\ R_n(w) = \frac{1}{n} \sum_{\xi=1}^n f(w; \xi) & \text{the empirical risk.} \end{cases}$$

The stochastic vector could be

$$g(w_k, \xi_k) = \begin{cases} \nabla f(w_k; \xi_k) & \text{the gradient for one example,} \\ \frac{1}{n_k} \sum_{i=1}^{n_k} \nabla f(w_k; \xi_{k,i}) & \text{the gradient for a minibatch,} \\ H_k \frac{1}{n_k} \sum_{i=1}^{n_k} \nabla f(w_k; \xi_{k,i}), & \text{possibly rescaled} \end{cases}$$

# The generic SG algorithm

## Stochastic processes

- We assume that the  $\{\xi_k\}$  are jointly independent to avoid the full machinery of stochastic processes. But everything still holds if the  $\{\xi_k\}$  form an adapted stochastic process, where each  $\xi_k$  can depend on the previous ones.

## Active learning

- We can handle more complex setups by view  $\xi_k$  as a “random seed”. For instance, in active learning,  $g(w_k, \xi_k)$  firsts construct a multinomial distribution on the training examples in a manner that depends on  $w_k$ , then uses the random seed  $\xi_k$  to pick one according to that distribution.

The same mathematics cover all these cases.

## 2- Fundamental lemmas

# Smoothness

## Smoothness

- Our analysis relies on a smoothness assumption.  
We chose this path because it also gives results for the nonconvex case.  
We'll discuss other paths in the commentary section.

**Assumption 4.1 (Lipschitz-continuous gradients).** *The objective function  $F : \mathbb{R}^d \rightarrow \mathbb{R}$  is continuously differentiable and its gradient,  $\nabla F : \mathbb{R}^d \rightarrow \mathbb{R}^d$ , is Lipschitz continuous with Lipschitz constant  $L > 0$ , i.e.,*

$$\|\nabla F(w) - \nabla F(\bar{w})\|_2 \leq L\|w - \bar{w}\|_2 \quad \text{for all } \{w, \bar{w}\} \subset \mathbb{R}^d.$$

## Well known consequence

$$F(w) \leq F(\bar{w}) + \nabla F(\bar{w})^T(w - \bar{w}) + \frac{1}{2}L\|w - \bar{w}\|_2^2 \quad \text{for all } \{w, \bar{w}\} \subset \mathbb{R}^d. \quad (4.3)$$

# Smoothness

- $\mathbb{E}_{\xi_k}[\cdot]$  is the expectation with respect to the distribution of  $\xi_k$  only.
- $\mathbb{E}_{\xi_k}[F(w_{k+1})]$  is meaningful because  $w_{k+1}$  depends on  $\xi_k$  (step 6 of SG)

**Lemma 4.2.** Under Assumption 4.1, the iterates of SG (Algorithm 4.1) satisfy the following inequality for all  $k \in \mathbb{N}$ :

$$\begin{aligned} & \mathbb{E}_{\xi_k}[F(w_{k+1})] - F(w_k) \\ & \leq -\alpha_k \nabla F(w_k)^T \mathbb{E}_{\xi_k}[g(w_k, \xi_k)] + \frac{1}{2} \alpha_k^2 L \mathbb{E}_{\xi_k}[\|g(w_k, \xi_k)\|_2^2]. \quad (4.4) \end{aligned}$$

Expected decrease

Noise

# Smoothness

**Lemma 4.2.** Under Assumption 4.1, the iterates of SG (Algorithm 4.1) satisfy the following inequality for all  $k \in \mathbb{N}$ :

$$\begin{aligned} & \mathbb{E}_{\xi_k}[F(w_{k+1})] - F(w_k) \\ & \leq -\alpha_k \nabla F(w_k)^T \mathbb{E}_{\xi_k}[g(w_k, \xi_k)] + \frac{1}{2}\alpha_k^2 L \mathbb{E}_{\xi_k}[\|g(w_k, \xi_k)\|_2^2]. \end{aligned} \quad (4.4)$$

*Proof.* By Assumption 4.1, the iterates generated by SG satisfy

$$\begin{aligned} F(w_{k+1}) - F(w_k) & \leq \nabla F(w_k)^T (w_{k+1} - w_k) + \frac{1}{2}L\|w_{k+1} - w_k\|_2^2 \\ & \leq -\alpha_k \nabla F(w_k)^T g(w_k, \xi_k) + \frac{1}{2}\alpha_k^2 L \|g(w_k, \xi_k)\|_2^2. \end{aligned}$$

Taking expectations in these inequalities with respect to the distribution of  $\xi_k$ , and noting that  $w_{k+1}$ —but not  $w_k$ —depends on  $\xi_k$ , we obtain the desired bound.  $\square$

# Moments

**Assumption 4.3 (First and second moment limits).** *The objective function and SG (Algorithm 4.1) satisfy the following:*

- (a) *The sequence of iterates  $\{w_k\}$  is contained in an open set over which  $F$  is bounded below by a scalar  $F_{\inf}$ .*
- (b) *There exist scalars  $\mu_G \geq \mu > 0$  such that, for all  $k \in \mathbb{N}$ ,*

$$\nabla F(w_k)^T \mathbb{E}_{\xi_k} [g(w_k, \xi_k)] \geq \mu \|\nabla F(w_k)\|_2^2 \quad \text{and} \quad (4.7a)$$

$$\|\mathbb{E}_{\xi_k} [g(w_k, \xi_k)]\|_2 \leq \mu_G \|\nabla F(w_k)\|_2. \quad (4.7b)$$

- (c) *There exist scalars  $M \geq 0$  and  $M_V \geq 0$  such that, for all  $k \in \mathbb{N}$ ,*

$$\mathbb{V}_{\xi_k} [g(w_k, \xi_k)] \leq M + M_V \|\nabla F(w_k)\|_2^2. \quad (4.8)$$

# Moments

(b) There exist scalars  $\mu_G \geq \mu > 0$  such that, for all  $k \in \mathbb{N}$ ,

$$\nabla F(w_k)^T \mathbb{E}_{\xi_k}[g(w_k, \xi_k)] \geq \mu \|\nabla F(w_k)\|_2^2 \quad \text{and} \quad (4.7a)$$

$$\|\mathbb{E}_{\xi_k}[g(w_k, \xi_k)]\|_2 \leq \mu_G \|\nabla F(w_k)\|_2. \quad (4.7b)$$

(c) There exist scalars  $M \geq 0$  and  $M_V \geq 0$  such that, for all  $k \in \mathbb{N}$ ,

$$\mathbb{V}_{\xi_k}[g(w_k, \xi_k)] \leq M + M_V \|\nabla F(w_k)\|_2^2. \quad (4.8)$$

- In expectation  $g(w_k, \xi_k)$  is a sufficient descent direction.
- True if  $\mathbb{E}_{\xi_k}[g(w_k, \xi_k)] = \nabla F(w_k)$  with  $\mu = \mu_G = 1$ .
- True if  $\mathbb{E}_{\xi_k}[g(w_k, \xi_k)] = H_k \nabla F(w_k)$  with bounded spectrum.

# Moments

(b) There exist scalars  $\mu_G \geq \mu > 0$  such that, for all  $k \in \mathbb{N}$ ,

$$\nabla F(w_k)^T \mathbb{E}_{\xi_k}[g(w_k, \xi_k)] \geq \mu \|\nabla F(w_k)\|_2^2 \quad \text{and} \quad (4.7a)$$

$$\|\mathbb{E}_{\xi_k}[g(w_k, \xi_k)]\|_2 \leq \mu_G \|\nabla F(w_k)\|_2. \quad (4.7b)$$

(c) There exist scalars  $M \geq 0$  and  $M_V \geq 0$  such that, for all  $k \in \mathbb{N}$ ,

$$\mathbb{V}_{\xi_k}[g(w_k, \xi_k)] \leq M + M_V \|\nabla F(w_k)\|_2^2. \quad (4.8)$$

- $\mathbb{V}_{\xi_k}[\cdot]$  denotes the variance w.r.t.  $\xi_k$
- Variance of the noise must be bounded in a mild manner.

# Moments

(b) There exist scalars  $\mu_G \geq \mu > 0$  such that, for all  $k \in \mathbb{N}$ ,

$$\nabla F(w_k)^T \mathbb{E}_{\xi_k}[g(w_k, \xi_k)] \geq \mu \|\nabla F(w_k)\|_2^2 \quad \text{and} \quad (4.7a)$$

$$\|\mathbb{E}_{\xi_k}[g(w_k, \xi_k)]\|_2 \leq \mu_G \|\nabla F(w_k)\|_2. \quad (4.7b)$$

(c) There exist scalars  $M \geq 0$  and  $M_V \geq 0$  such that, for all  $k \in \mathbb{N}$ ,

$$\mathbb{V}_{\xi_k}[g(w_k, \xi_k)] \leq M + M_V \|\nabla F(w_k)\|_2^2. \quad (4.8)$$

- Combining (4.7b) and (4.8) gives

$$\mathbb{E}_{\xi_k}[\|g(w_k, \xi_k)\|_2^2] \leq M + M_G \|\nabla F(w_k)\|_2^2 \quad (4.9)$$

with  $M_G := M_V + \mu_G^2 \geq \mu^2 > 0$ .

# Moments

**Lemma 4.4.** Under Assumptions 4.1 and 4.3, the iterates of SG (Algorithm 4.1) satisfy the following inequalities for all  $k \in \mathbb{N}$ :

$$\begin{aligned}\mathbb{E}_{\xi_k}[F(w_{k+1})] - F(w_k) \\ \leq -\mu\alpha_k \|\nabla F(w_k)\|_2^2 + \frac{1}{2}\alpha_k^2 L \mathbb{E}_{\xi_k}[\|g(w_k, \xi_k)\|_2^2] \quad (4.10a)\end{aligned}$$

$$\leq -(\mu - \frac{1}{2}\alpha_k L M_G)\alpha_k \|\nabla F(w_k)\|_2^2 + \frac{1}{2}\alpha_k^2 L M. \quad (4.10b)$$

Expected decrease

Noise

- The convergence of SG depends on the balance between these two terms.

# Moments

**Lemma 4.4.** Under Assumptions 4.1 and 4.3, the iterates of SG (Algorithm 4.1) satisfy the following inequalities for all  $k \in \mathbb{N}$ :

$$\begin{aligned}\mathbb{E}_{\xi_k}[F(w_{k+1})] - F(w_k) &\leq -\mu\alpha_k \|\nabla F(w_k)\|_2^2 + \frac{1}{2}\alpha_k^2 L \mathbb{E}_{\xi_k}[\|g(w_k, \xi_k)\|_2^2] \quad (4.10a)\\ &\leq -(\mu - \frac{1}{2}\alpha_k L M_G)\alpha_k \|\nabla F(w_k)\|_2^2 + \frac{1}{2}\alpha_k^2 L M.\end{aligned}$$

*Proof.* By Lemma 4.2 and (4.7a), it follows that

$$\begin{aligned}\mathbb{E}_{\xi_k}[F(w_{k+1})] - F(w_k) &\leq -\alpha_k \nabla F(w_k)^T \mathbb{E}_{\xi_k}[g(w_k, \xi_k)] + \frac{1}{2}\alpha_k^2 L \mathbb{E}_{\xi_k}[\|g(w_k, \xi_k)\|_2^2]\\ &\leq -\mu\alpha_k \|\nabla F(w_k)\|_2^2 + \frac{1}{2}\alpha_k^2 L \mathbb{E}_{\xi_k}[\|g(w_k, \xi_k)\|_2^2],\end{aligned}$$

which is (4.10a). Assumption 4.3, giving (4.9), then yields (4.10b). □

## 3- SG for Strongly Convex Objectives

# Strong convexity

**Assumption 4.5 (Strong convexity).** *The objective function  $F : \mathbb{R}^d \rightarrow \mathbb{R}$  is strongly convex in that there exists a constant  $c > 0$  such that for all  $(\bar{w}, w) \in \mathbb{R}^d \times \mathbb{R}^d$*

$$F(\bar{w}) \geq F(w) + \nabla F(w)^T(\bar{w} - w) + \frac{1}{2}c\|\bar{w} - w\|_2^2. \quad (4.11)$$

*Hence,  $F$  has a unique minimizer, denoted as  $w_* \in \mathbb{R}^d$  with  $F_* := F(w_*)$ .*

## Known consequence

$$2c(F(w) - F_*) \leq \|\nabla F(w)\|_2^2 \text{ for all } w \in \mathbb{R}^d. \quad (4.12)$$

## Why does strong convexity matter?

- It gives the strongest results.
- It often happens in practice (one regularizes to facilitate optimization!)
- It describes any smooth function near a strong local minimum.

# Total expectation

## Different expectations

- $\mathbb{E}_{\xi_k}[\quad]$  is the expectation with respect to the distribution of  $\xi_k$  only.
- $\mathbb{E}[\quad]$  is the total expectation w.r.t. the joint distribution of all  $\xi_k$ .

For instance, since  $w_k$  depends only on  $\xi_1, \xi_2, \dots, \xi_{k-1}$ ,

$$\mathbb{E}[F(w_k)] = \mathbb{E}_{\xi_1} \mathbb{E}_{\xi_2} \dots \mathbb{E}_{\xi_{k-1}}[F(w_k)]$$

## Results in expectation

- We focus on results that characterize the properties of SG in expectation.
- The stochastic approximation literature usually relies on rather complex martingale techniques to establish almost sure convergence results. We avoid them because they do not give much additional insight.

# SG with fixed stepsize

**Theorem 4.6 (Strongly Convex Objective, Fixed Stepsize).** Under Assumptions 4.1, 4.3, and 4.5 (with  $F_{\inf} = F_*$ ), suppose that the SG method (Algorithm 4.1) is run with a fixed stepsize,  $\alpha_k = \bar{\alpha}$  for all  $k \in \mathbb{N}$ , satisfying

$$0 < \bar{\alpha} \leq \frac{\mu}{LM_G}. \quad (4.13)$$

Then, for all  $k \in \mathbb{N}$  the expected optimality gap satisfies :

$$\begin{aligned} \mathbb{E}[F(w_k) - F_*] &\leq \frac{\bar{\alpha}LM}{2c\mu} + (1 - \bar{\alpha}c\mu)^{k-1} \left( F(w_1) - F_* - \frac{\bar{\alpha}LM}{2c\mu} \right) \\ &\xrightarrow{k \rightarrow \infty} \frac{\bar{\alpha}LM}{2c\mu}. \end{aligned} \quad (4.14)$$

- Only converges to a neighborhood of the optimal value.
- Both (4.13) and (4.14) describe well the actual behavior.

# SG with fixed stepsize (proof)

*Proof.* Using Lemma 4.4 with (4.13) and (4.12), we have for all  $k \in \mathbb{N}$  that

$$\begin{aligned}\mathbb{E}_{\xi_k}[F(w_{k+1})] - F(w_k) &\leq -(\mu - \frac{1}{2}\bar{\alpha}LM_G)\bar{\alpha}\|\nabla F(w_k)\|_2^2 + \frac{1}{2}\bar{\alpha}^2 LM \\ &\leq -\frac{1}{2}\bar{\alpha}\mu\|\nabla F(w_k)\|_2^2 + \frac{1}{2}\bar{\alpha}^2 LM \\ &\leq -\bar{\alpha}c\mu(F(w_k) - F_*) + \frac{1}{2}\bar{\alpha}^2 LM.\end{aligned}$$

Subtracting  $F_*$  from both sides and taking total expectations,

$$\mathbb{E}[F(w_{k+1}) - F_*] \leq (1 - \bar{\alpha}c\mu)\mathbb{E}[F(w_k) - F_*] + \frac{1}{2}\bar{\alpha}^2 LM.$$

Subtracting the constant  $\bar{\alpha}LM/(2c\mu)$  from both sides, one obtains

$$\mathbb{E}[F(w_{k+1}) - F_*] - \frac{\bar{\alpha}LM}{2c\mu} \leq (1 - \bar{\alpha}c\mu) \left( \mathbb{E}[F(w_k) - F_*] - \frac{\bar{\alpha}LM}{2c\mu} \right). \quad (4.15)$$

Observe that (4.15) is a contraction inequality since, by (4.13) and (4.9),

$$0 < \bar{\alpha}c\mu \leq \frac{c\mu^2}{LM_G} \leq \frac{c\mu^2}{L\mu^2} = \frac{c}{L} \leq 1. \quad (4.16)$$

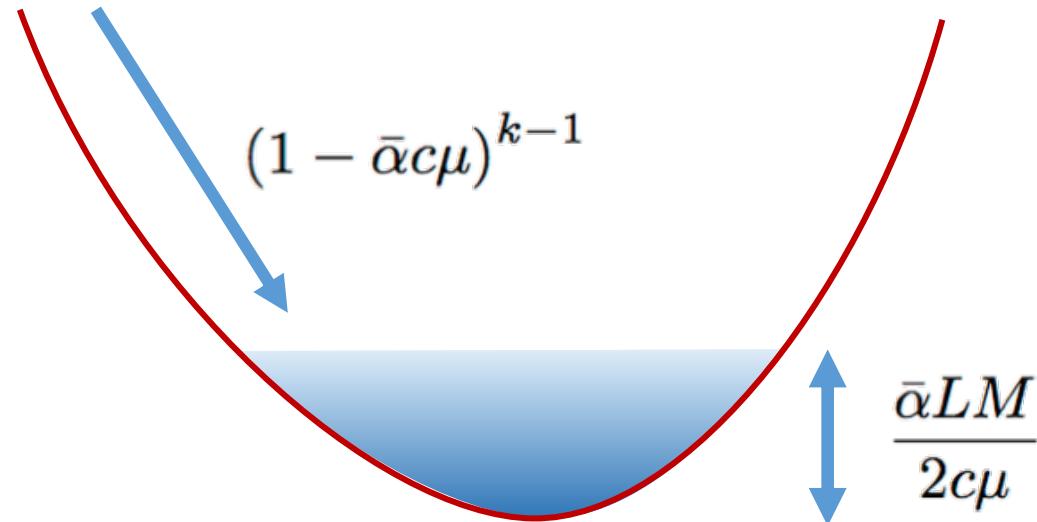
The result thus follows by applying (4.15) repeatedly. □

# SG with fixed stepsize

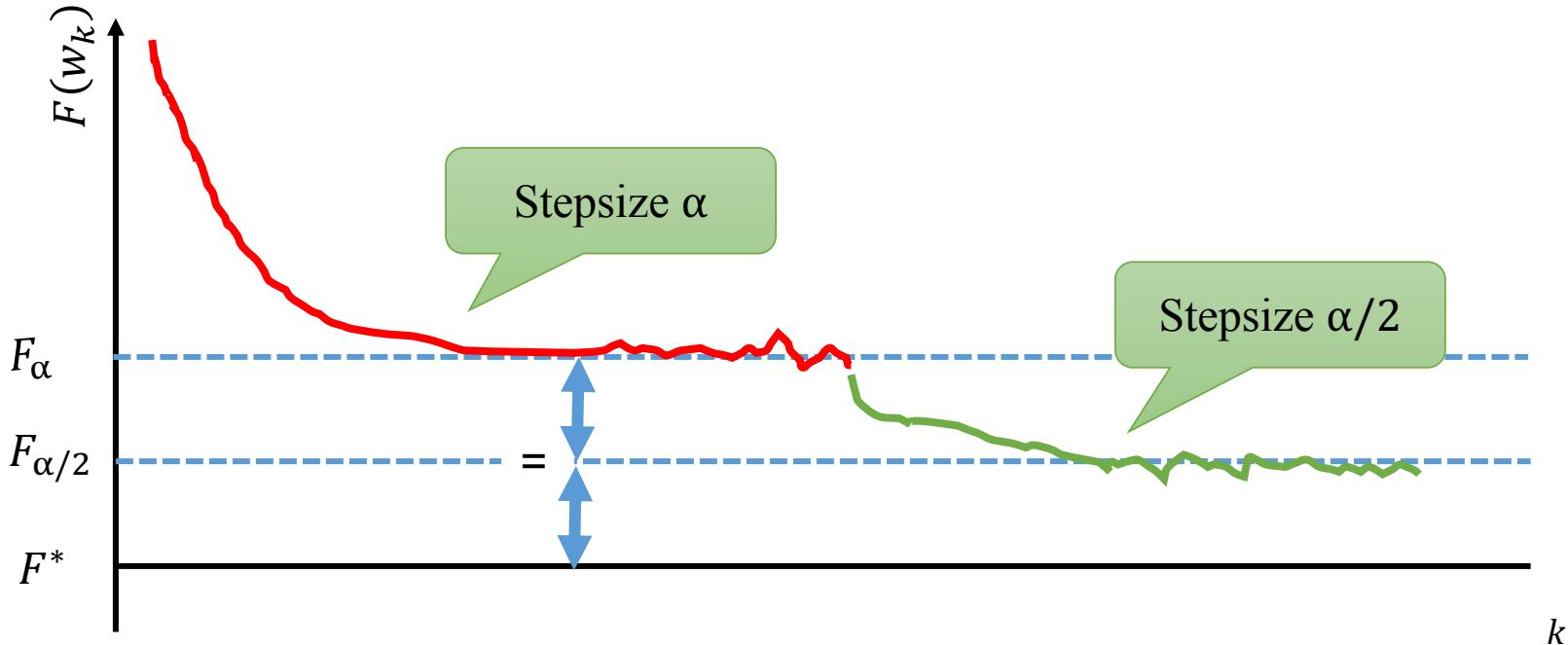
$$\mathbb{E}[F(w_k) - F_*] \leq \frac{\bar{\alpha}LM}{2c\mu} + (1 - \bar{\alpha}c\mu)^{k-1} \left( F(w_1) - F_* - \frac{\bar{\alpha}LM}{2c\mu} \right) \quad (4.14)$$

Note the interplay between the stepsize  $\bar{\alpha}$  and the variance bound  $M$ .

- If  $M = 0$ , one recovers the linear convergence of batch gradient descent.
- If  $M > 0$ , one reaches a point where the noise prevents further progress.

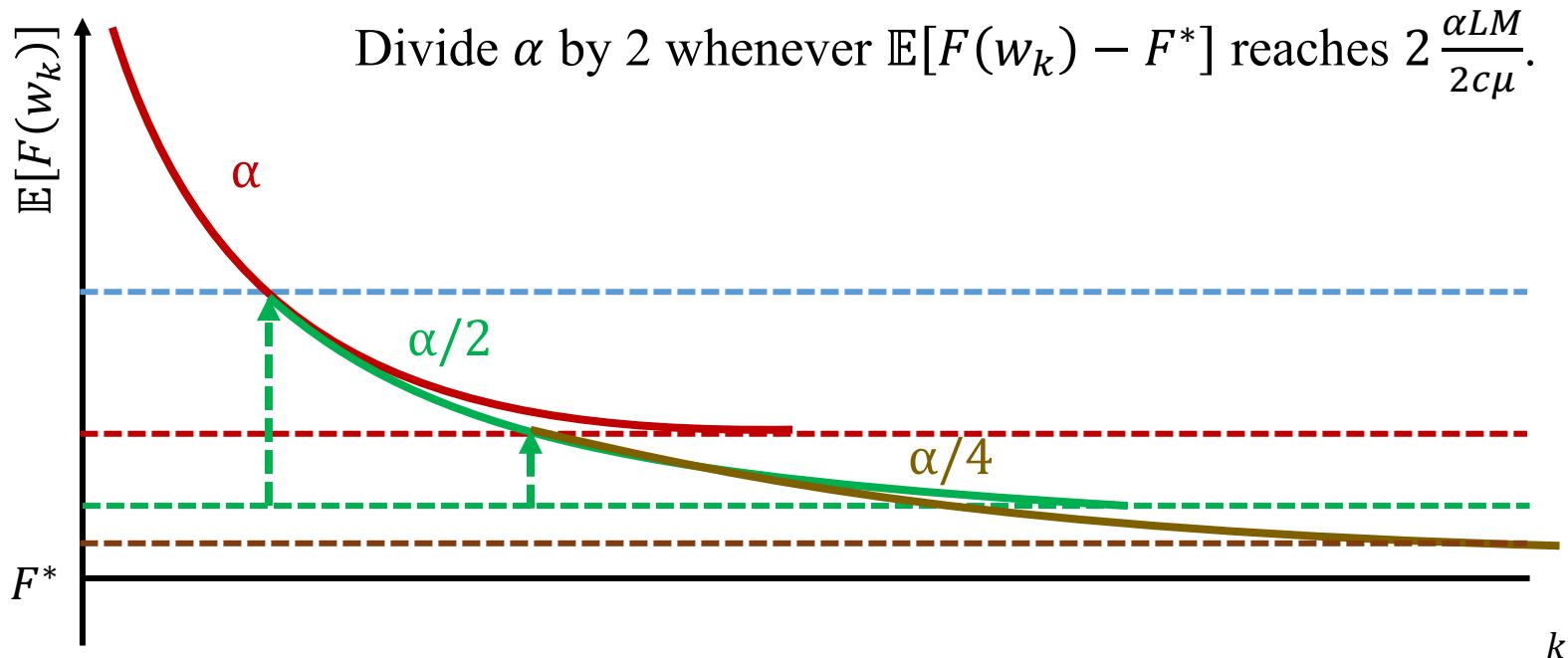


# Diminishing the stepsizes



- If we wait long enough, halving the stepsize  $\alpha$  eventually halves  $F(w_k) - F^*$ .
- We can even estimate  $F^* \approx 2F_{\alpha/2} - F_\alpha$

# Diminishing the stepsizes faster



- Divide  $\alpha$  by 2 whenever  $\mathbb{E}[F(w_k)]$  reaches  $\alpha LM / c\mu$ .
- Time  $\tau_\alpha$  between changes :  $(1 - \alpha c\mu)^{\tau_\alpha} = 1/3$  means  $\tau_\alpha \propto 1/\alpha$ .
- Whenever we halve  $\alpha$  we must wait twice as long to halve  $F(w) - F^*$ .
- Overall convergence rate in  $\mathcal{O}(1/k)$ .

# SG with diminishing stepsizes

**Theorem 4.7 (Strongly Convex Objective, Diminishing Stepsizes).** *Under Assumptions 4.1, 4.3, and 4.5 (with  $F_{\inf} = F_*$ ), suppose that SG (Algorithm 4.1) is run with a stepsize sequence such that, for all  $k \in \mathbb{N}$ ,*

$$\alpha_k = \frac{\beta}{\gamma + k} \text{ for some } \beta > \frac{1}{c\mu} \text{ and } \gamma > 0 \text{ s.t. } \alpha_1 \leq \frac{\mu}{LM_G}. \quad (4.18)$$

*Then, for all  $k \in \mathbb{N}$ , the expected optimality gap satisfies*

$$\mathbb{E}[F(w_k) - F_*] \leq \frac{\nu}{\gamma + k}, \quad (4.19)$$

*where*

$$\nu := \max \left\{ \frac{\beta^2 LM}{2(\beta c\mu - 1)}, (\gamma + 1)(F(w_1) - F_*) \right\}. \quad (4.20)$$

# SG with diminishing stepsizes

## Theorem 4.7 (Strongly Convex Objective)

Stepsize decreases in  $1/k$

Same maximal stepsize (as in 4.1, 4.3, and 4.5 (with  $r_{\inf} = r_*$ ), subject to the condition that SG converges in a finite number of iterations).  
In other words, there exists a stepsize sequence such that, for all  $k \in \mathbb{N}$ ,

$$\alpha_k = \frac{\beta}{\gamma + k} \text{ for some } \beta > \frac{1}{c\mu} \text{ and } \gamma > 0 \text{ s.t. } \alpha_1 \leq \frac{\mu}{LM_G}. \quad (4.18)$$

For  $k \in \mathbb{N}$ , the expected optimality gap satisfies

$$\mathbb{E}[F(w_k) - F_*] \leq \frac{\nu}{\gamma + k}, \quad (4.19)$$

where

$$\nu := \max \left\{ \frac{\beta^2 LM}{2(\beta c\mu - 1)}, (\gamma + 1)(F(w_1) - F_*) \right\}. \quad (4.20)$$

...otherwise

gap  $\propto$  stepsize

# SG with diminishing stepsizes (proof)

*Proof.* Proceeding as in the proof of Theorem 4.6, one gets

$$\mathbb{E}[F(w_{k+1}) - F_*] \leq (1 - \alpha_k c\mu) \mathbb{E}[F(w_k) - F_*] + \frac{1}{2} \alpha_k^2 LM. \quad (4.21)$$

We now prove (4.19) by induction. First, the definition of  $\nu$  ensures that it holds for  $k = 1$ . Then, assuming (4.19) holds for some  $k \geq 1$ , it follows from (4.21) that

$$\begin{aligned} \mathbb{E}[F(w_{k+1}) - F_*] &\leq \left(1 - \frac{\beta c\mu}{\hat{k}}\right) \frac{\nu}{\hat{k}} + \frac{\beta^2 LM}{2\hat{k}^2} \quad (\text{with } \hat{k} := \gamma + k) \\ &= \left(\frac{\hat{k}-1}{\hat{k}^2}\right) \nu - \underbrace{\left(\frac{\beta c\mu - 1}{\hat{k}^2}\right) \nu}_{\text{nonpositive by the definition of } \nu} + \frac{\beta^2 LM}{2\hat{k}^2} \leq \frac{\nu}{\hat{k}+1}, \end{aligned}$$

where the last inequality follows because  $\hat{k}^2 \geq (\hat{k}+1)(\hat{k}-1)$ .  $\square$

# Mini batching

	Computation	Noise
$g(w_k, \xi_k) = \begin{cases} \nabla f(w_k; \xi_k) \\ \frac{1}{n_{\text{mb}}} \sum_{i=1}^{n_{\text{mb}}} \nabla f(w_k; \xi_{k,i}) \end{cases}$	1	$M$
	$n_{\text{mb}}$	$M/n_{\text{mb}}$

Using minibatches with stepsize  $\bar{\alpha}$  :

$$\mathbb{E}[F(w_k) - F_*] \leq \frac{\bar{\alpha}LM}{2c\mu n_{\text{mb}}} + [1 - \bar{\alpha}c\mu]^{k-1} \left( F(w_1) - F_* - \frac{\bar{\alpha}LM}{2c\mu n_{\text{mb}}} \right).$$

Using single example with stepsize  $\bar{\alpha} / n_{\text{mb}}$  :

$$\mathbb{E}[F(w_k) - F_*] \leq \frac{\bar{\alpha}LM}{2c\mu n_{\text{mb}}} + \left[ 1 - \frac{\bar{\alpha}c\mu}{n_{\text{mb}}} \right]^{k-1} \left( F(w_1) - F_* - \frac{\bar{\alpha}LM}{2c\mu n_{\text{mb}}} \right).$$

$n_{\text{mb}}$  times more iterations that are  $n_{\text{mb}}$  times cheaper.

same

# Minibatching

## Ignoring implementation issues

- We can match minibatch SG with stepsize  $\bar{\alpha}$  using single example SG with stepsize  $\bar{\alpha} / n_{mb}$  .
- We can match single example SG with stepsize  $\bar{\alpha}$  using minibatch SG with stepsize  $\bar{\alpha} \times n_{mb}$  provided  $\bar{\alpha} \times n_{mb}$  is smaller than the max stepsize.

## With implementation issues

- Minibatch implementations use the hardware better.
- Especially on GPU.

## 4- SG for General Objectives

# Nonconvex objectives

**Nonconvex training objectives are pervasive in deep learning.**

**Nonconvex landscape in high dimension can be very complex.**

- Critical points can be local minima or saddle points.
- Critical points can be first order of high order.
- Critical points can be part of critical manifolds.
- A critical manifold can contain both local minima and saddle points.

**We describe meaningful (but weak) guarantees**

- Essentially, SG goes to critical points.

**The SG noise plays an important role in practice**

- It seems to help navigating local minima and saddle points.
- More noise has been found to sometimes help optimization.
- But the theoretical understanding of these facts is weak.

# Nonconvex SG with fixed stepsize

**Theorem 4.8 (Nonconvex Objective, Fixed Stepsize).** *Under Assumptions 4.1 and 4.3, suppose that the SG method (Algorithm 4.1) is run with a fixed stepsize,  $\alpha_k = \bar{\alpha}$  for all  $k \in \mathbb{N}$ , satisfying*

$$0 < \bar{\alpha} \leq \frac{\mu}{LM_G}. \quad (4.25)$$

*Then, the expected sum-of-squares and average-squared gradients of  $F$  corresponding to the SG iterates satisfy the following inequalities for all  $K \in \mathbb{N}$ :*

$$\mathbb{E} \left[ \sum_{k=1}^K \|\nabla F(w_k)\|_2^2 \right] \leq \frac{K\bar{\alpha}LM}{\mu} + \frac{2(F(w_1) - F_{\inf})}{\mu\bar{\alpha}} \quad (4.26a)$$

and therefore  $\mathbb{E} \left[ \frac{1}{K} \sum_{k=1}^K \|\nabla F(w_k)\|_2^2 \right] \leq \frac{\bar{\alpha}LM}{\mu} + \frac{2(F(w_1) - F_{\inf})}{K\mu\bar{\alpha}} \quad (4.26b)$

# Nonconvex SG with fixed stepsize

**Theorem 4.8 (Nonconvex Objective Function)** Under Assumptions 4.1 and 4.3, suppose that the SG run with a fixed stepsize,  $\alpha_k = \bar{\alpha}$  for all  $k \in \mathbb{N}$ , satisfies

Same max stepsize

$$0 < \bar{\alpha} \leq \frac{\mu}{LM_G}. \quad (4.25)$$

If the average norm of the gradient is small, then the norm of the gradient cannot be often large...

of-squares and averages satisfy the following:

This goes to zero like  $1/K$  even if  $K$ :

This does not

$$\mathbb{E} \left[ \sum_{k=1}^K \|\nabla F(w_k)\|_2^2 \right] \leq \frac{2(F(w_1) - F_{\inf})}{\mu \bar{\alpha}} \quad (4.26a)$$

$$\text{and therefore } \mathbb{E} \left[ \frac{1}{K} \sum_{k=1}^K \|\nabla F(w_k)\|_2^2 \right] \leq \frac{\bar{\alpha} LM}{\mu} + \frac{2(F(w_1) - F_{\inf})}{K \mu \bar{\alpha}} \quad (4.26b)$$

# Nonconvex SG with fixed stepsize (proof)

*Proof.* Taking the total expectation of (4.10b) and from (4.25),

$$\begin{aligned}\mathbb{E}[F(w_{k+1})] - \mathbb{E}[F(w_k)] &\leq -(\mu - \frac{1}{2}\bar{\alpha}LM_G)\bar{\alpha}\mathbb{E}[\|\nabla F(w_k)\|_2^2] + \frac{1}{2}\bar{\alpha}^2LM \\ &\leq -\frac{1}{2}\mu\bar{\alpha}\mathbb{E}[\|\nabla F(w_k)\|_2^2] + \frac{1}{2}\bar{\alpha}^2LM.\end{aligned}$$

Summing both sides of this inequality for  $k \in \{1, \dots, K\}$  and recalling Assumption 4.3(a) gives

$$F_{\inf} - F(w_1) \leq \mathbb{E}[F(w_{K+1})] - F(w_1) \leq -\frac{1}{2}\mu\bar{\alpha} \sum_{k=1}^K \mathbb{E}[\|\nabla F(w_k)\|_2^2] + \frac{1}{2}K\bar{\alpha}^2LM.$$

Rearranging yields (4.26a), and dividing further by  $K$  yields (4.26b). □

# Nonconvex SG with diminishing step sizes

**Theorem 4.10 (Nonconvex Objective, Diminishing Stepsizes).** *Under Assumptions 4.1 and 4.3, suppose that the SG method (Algorithm 4.1) is run with a stepsize sequence satisfying*

$$\sum_{k=1}^{\infty} \alpha_k = \infty \quad \sum_{k=1}^{\infty} \alpha_k^2 < \infty ,$$

*then*

$$\mathbb{E} \left[ \sum_{k=1}^K \alpha_k \|\nabla F(w_k)\|_2^2 \right] < \infty$$

**Corollary 4.12.** *Under the conditions of Theorem 4.10, if we further assume that the objective function  $F$  is twice differentiable, and that the mapping  $w \mapsto \|\nabla F(w)\|_2^2$  has Lipschitz-continuous derivatives, then*

$$\lim_{k \rightarrow \infty} \mathbb{E}[\|\nabla F(w_k)\|_2^2] = 0.$$

## 5- Work complexity for Large-Scale Learning

# Large-Scale Learning

**Assume that we are in the large data regime**

- Training data is essentially unlimited.
- Computation time is limited.

**The good**

- More training data  $\Rightarrow$  less overfitting
- Less overfitting  $\Rightarrow$  richer models.

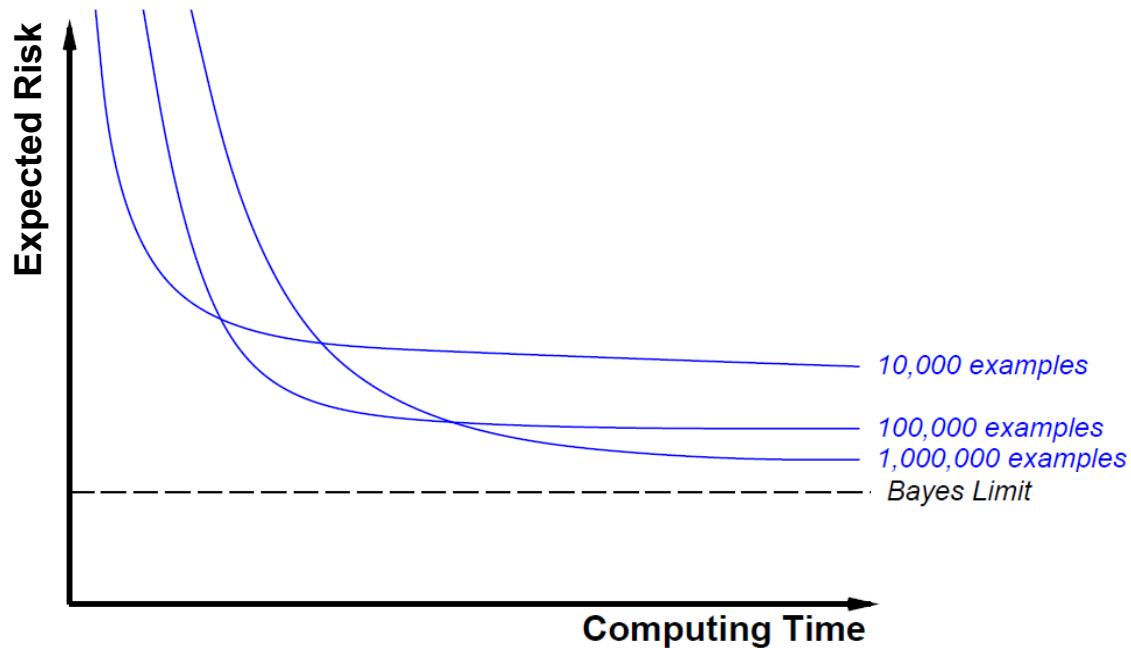
**The bad**

- Using more training data or rich models quickly exhausts the time budget.

**The hope**

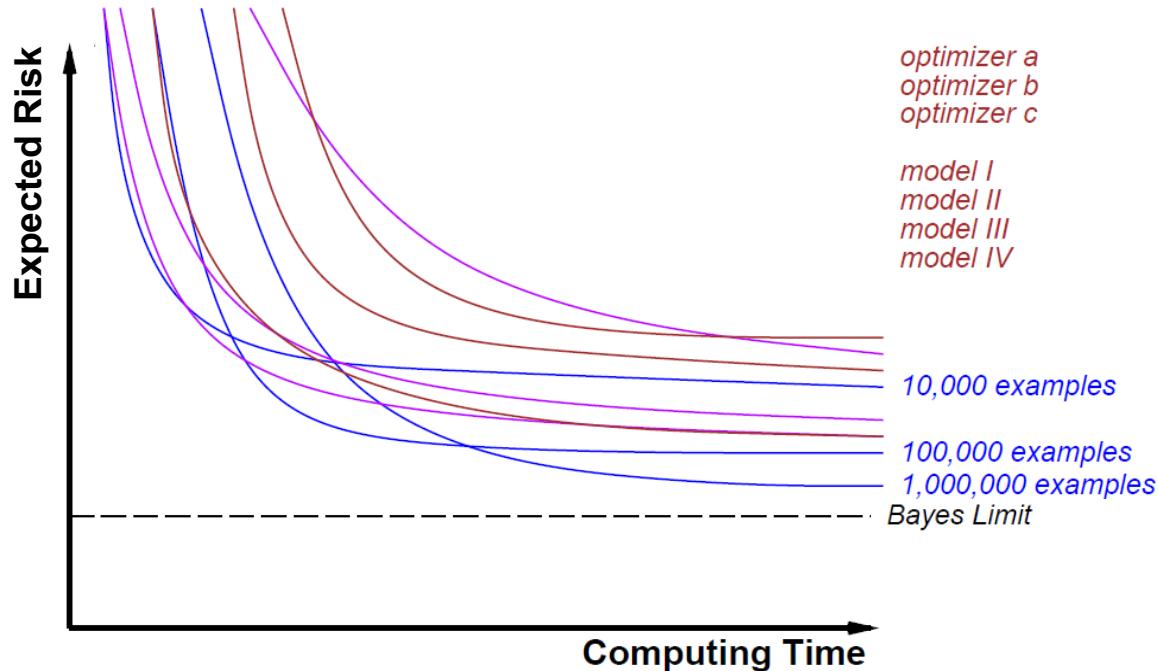
- How thoroughly do we need to optimize  $R_n(w)$  when we actually want another function  $R(w)$  to be small ?

# Expected risk versus training time



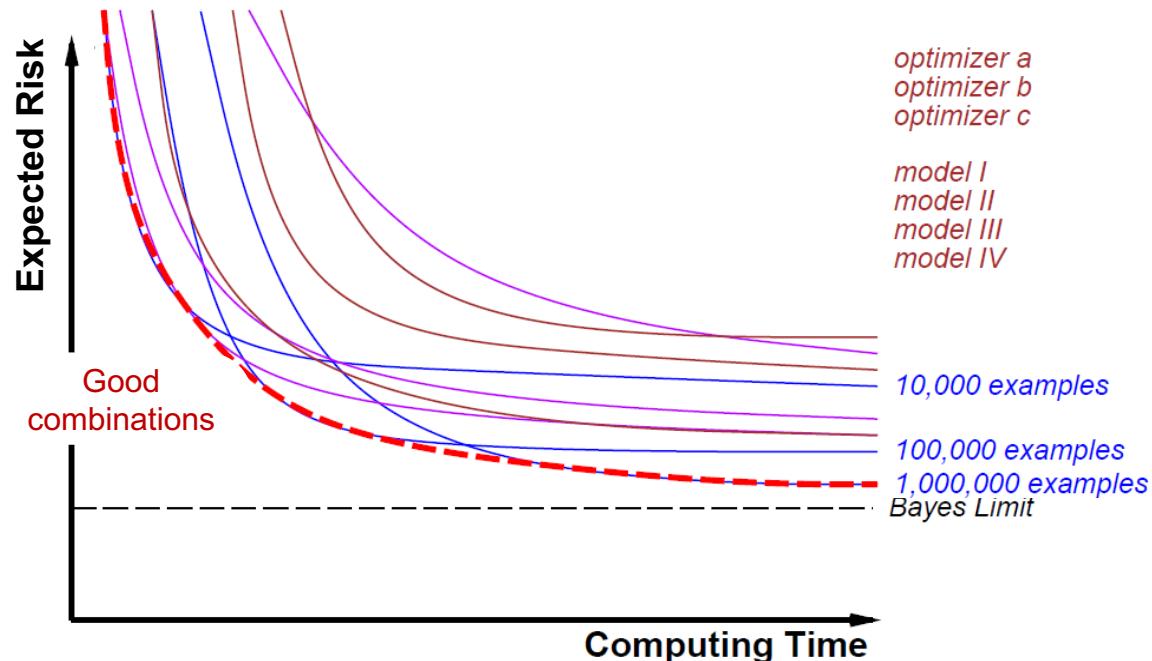
- When we vary the number of examples

# Expected risk versus training time



- When we vary the number of examples, the model, and the optimizer...

# Expected risk versus training time



- The optimal combination depends on the computing time budget

# Formalization

## The components of the expected risk

$$\mathbb{E}[R(\tilde{w}_n)] = \underbrace{R(w_*)}_{\mathcal{E}_{app}(\mathcal{H})} + \underbrace{\mathbb{E}[R(w_n) - R(w_*)]}_{\mathcal{E}_{est}(\mathcal{H}, n)} + \underbrace{\mathbb{E}[R(\tilde{w}_n) - R(w_n)]}_{\mathcal{E}_{opt}(\mathcal{H}, n, \epsilon)} \quad (4.29)$$

## Question

- Given a fixed model  $\mathcal{H}$  and a time budget  $\mathcal{T}_{\max}$ , choose  $n, \epsilon \dots$

$$\min_{n, \epsilon} \mathcal{E}(n, \epsilon) = \mathbb{E}[R(\tilde{w}_n) - R(w_*)] \text{ s.t. } \mathcal{T}(n, \epsilon) \leq \mathcal{T}_{\max}. \quad (4.30)$$

## Approach

- Statistics tell us  $\mathcal{E}_{est}(n)$  decreases with a rate in range  $1/\sqrt{n} \dots 1/n$ .
- For now, let's work with the fastest rate compatible with statistics

$$\mathcal{E}(n, \epsilon) \sim \frac{1}{n} + \epsilon \quad (4.32)$$

# Batch versus Stochastic

## Typical convergence rates

- Batch algorithm:  $\mathcal{T}(n, \epsilon) \sim n \log(1/\epsilon)$
- Stochastic algorithm:  $\mathcal{T}(n, \epsilon) \sim 1/n$

## Rate analysis

	Batch	Stochastic
$\mathcal{T}(n, \epsilon)$	$\sim n \log\left(\frac{1}{\epsilon}\right)$	$\frac{1}{\epsilon}$
$n^*$	$\sim \frac{\mathcal{T}_{\max}}{\log(\mathcal{T}_{\max})}$	$\mathcal{T}_{\max}$
$\mathcal{E}^*$	$\sim \frac{\log(\mathcal{T}_{\max})}{\mathcal{T}_{\max}} + \frac{1}{\mathcal{T}_{\max}}$	$\frac{1}{\mathcal{T}_{\max}}$

Processing more training examples beats optimizing more thoroughly.

This effect only grows if  $\mathcal{E}_{est}(n)$  decreases slower than  $1/n$ .

## 6- Comments

# Asymptotic performance of SG is fragile

## Diminishing stepsizes are tricky

- Theorem 4.7 (strongly convex function) suggests

$$\alpha_k = \frac{\beta}{\gamma + k}$$

SG converges very slowly if  $\beta < \frac{1}{c\mu}$

SG usually diverges when  $\alpha$  is above  $\frac{2\mu}{LM_G}$

## Constant stepsizes are often used in practice

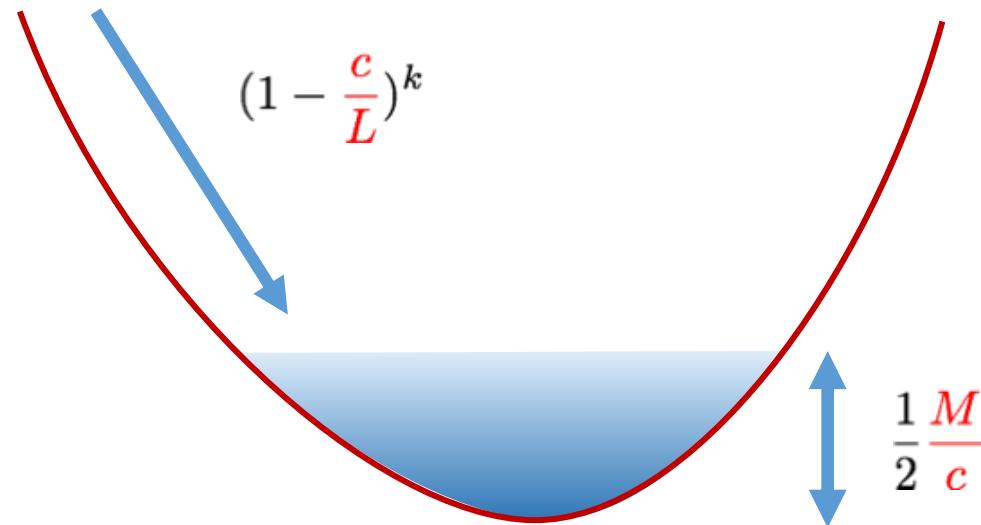
- Sometimes with a simple halving protocol.

**Spoiler** – Certain SG variants are more robust.

# Condition numbers

The ratios  $\frac{L}{c}$  and  $\frac{M}{c}$  appear in critical places

- Theorem 4.6. With  $\mu = 1, M_V = 0$ , the optimal stepsize is  $\bar{\alpha} = \frac{1}{L}$



# Distributed computing

## **SG is notoriously hard to parallelize**

- Because it updates the parameters  $w$  with high frequency
- Because it slows down with delayed updates.

## **SG still works with relaxed synchronization**

- Because this is just a little bit more noise.

## **Communication overhead give room for new opportunities**

- There is ample time to compute things while communication takes place.
  - Opportunity for optimization algorithms with higher per-iteration costs
- ☒ SG may not be the best answer for distributed training.

# Smoothness versus Convexity

## Analyses of SG that only rely on convexity

- Bounding  $\|w_k - w^*\|^2$  instead of  $F(w_k) - F^*$  and assuming  $\mathbb{E}_{\xi_k}[g(w_k, \xi_k)] = \hat{g}(w_k) \in \partial F(w_k)$  gives a result similar to Lemma 4.4.

$$\begin{aligned} & \mathbb{E}_{\xi_k}[\|w_{k+1} - w_*\|_2^2] - \|w_k - w_*\|_2^2 \\ &= -2\alpha_k \hat{g}(w_k)^T (w_k - w_*) + \alpha_k^2 \mathbb{E}_{\xi_k}[\|g(w_k, \xi_k)\|_2^2], \end{aligned} \quad (\text{A.2})$$

Expected decrease

Noise

- Ways to bound the expected decrease

General convexity :  $\hat{g}(w_k)^T (w_k - w_*) \geq F(w_k) - F(w_*) \geq 0$

Strong convexity :  $\hat{g}(w_k)^T (w_k - w_*) \geq c\|w_k - w_*\|^2 \geq 0$

- Proof does not easily support second order methods.

# Reparametrization and neural nets training algorithms

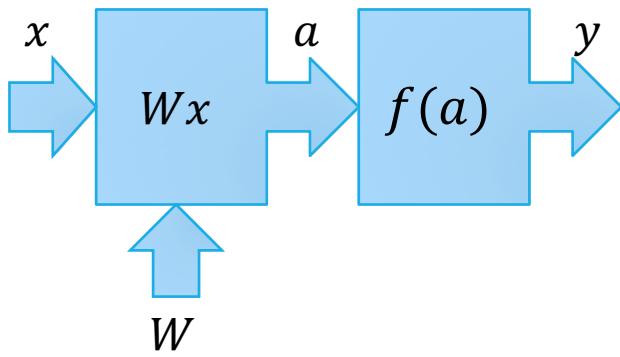
---

LÉON BOTTOU

FACEBOOK AI RESEARCH, NEW YORK

# Rescaling weights

---



## Propagation

- $y = f(Wx)$

## Back-propagation

- $g_a = f'(a) g_y$

- $g_x = g_a W$

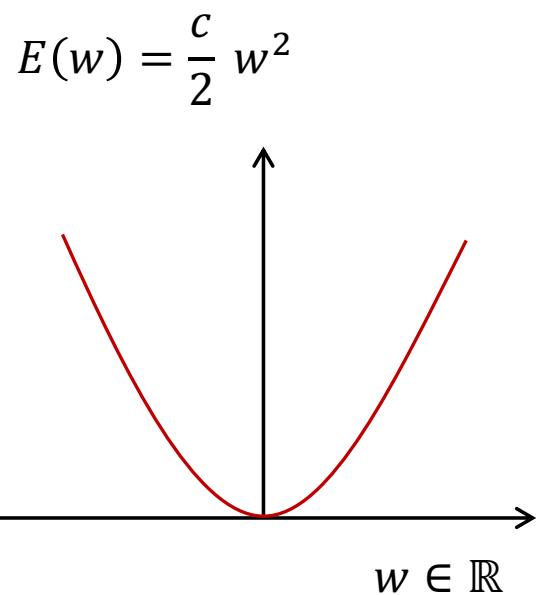
- $\Delta W = -\eta x g_a$

Consider the change  $f_{new}(a) = f(2a)$  and  $W_{new} = W/2$ .

- This leaves  $y(x)$  unchanged.
- What can you say about  $\Delta W_{new}$  ?

# Parabola

---



## Gradient descent

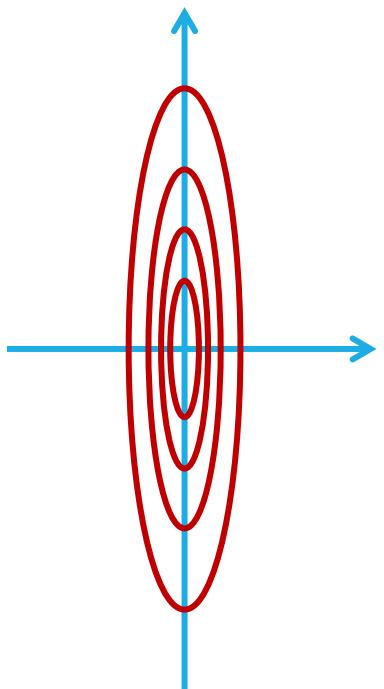
$$\blacksquare w_{t+1} = w_t - \eta \frac{dE}{dw}(w_t)$$

## Questions

- How does  $\eta$  affect the convergence?
- What's the best value of  $\eta$  ?

# More dimensions

---



Two dimensions.

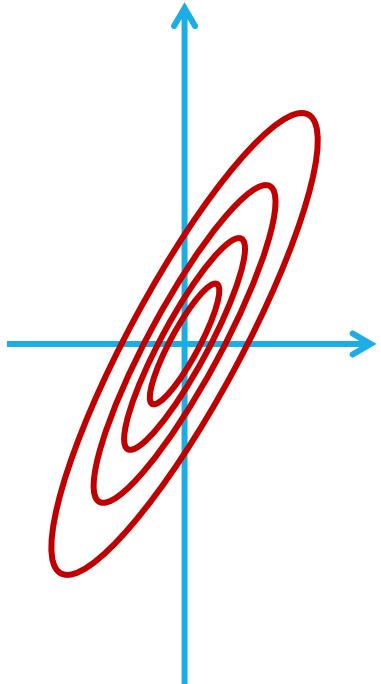
Two different curvatures.

Same questions

- How does  $\eta$  affect the convergence?
- What's the best value of  $\eta$  ?

# More dimensions

---



## Gradient descent

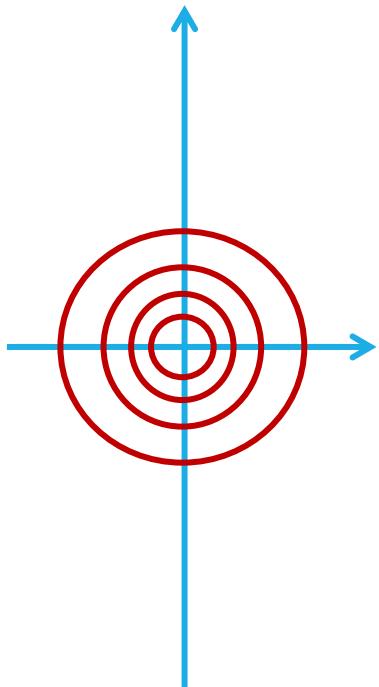
- $E(w) = \frac{1}{2} w^T H w$
- $w_{t+1} = w_t - \eta \frac{dE}{dw}(w_t)$

## Questions

- How does  $\eta$  affect the convergence?
- What's the best value of  $\eta$  ?

# Second order rescaling

---



## Rescale $w$

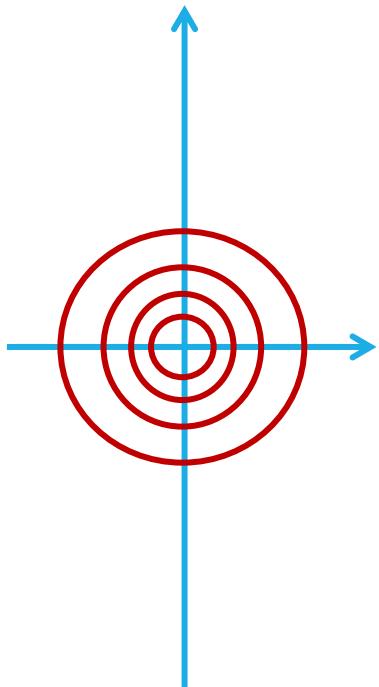
- $w_{new} \leftarrow H^{\frac{1}{2}} w$
- $E = \frac{1}{2} w^T H w = \frac{1}{2} w_{new}^T w_{new}$

## Questions

- Write gradient descent in  $w_{new}$  space.
- Write the equivalent  $w$  update.

# Second order rescaling

---



Rescale  $w$

- $w_{new} \leftarrow H^{\frac{1}{2}} w$
- $E = \frac{1}{2} w^T H w = \frac{1}{2} w_{new}^T w_{new}$

Gradient descent in  $w_{new}$  space

- $\Delta w_{new} = -\eta \frac{dE}{dw_{new}} = -\eta H^{-\frac{1}{2}} \frac{dE}{dw}$
- $\Delta w = -\eta H^{-1} \frac{dE}{dw}$

# Practical issues

---

- Objective function is not quadratic.
  - Local quadratic approximation is reasonable.
  - Hessian  $H$  changes with  $w$ .
  - When objective is non-convex,  $H$  can have negative eigenvalues
- Estimate the Hessian  $H$  on the fly.
- The Hessian is often too large to store or invert.

# Standard solutions for batch optimization

---

Idea: estimate a compact approximation of  $H^{-1}$

using the observed gradients  $g(w_t), g(w_{t-1}), \dots, g(w_{t-k}), \dots$   
then use approximate line search along direction  $\hat{H}^{-1}g(w_t)$

Idea: use exact line search and ensure conjugate search directions.

- Let  $d_{t-1}, d_{t-2}, \dots, d_{t-k}$  be the last  $k$  search directions.  
We want to choose  $d_t = g(w_t) + \sum \lambda_i d_{t-i}$  such that  $d_t^T H d_{t-i} = 0$

Very good algorithm have been developed.

- Conjugate gradient ( $k = 1$ , exact line search)
- LBFGS ( $k > 1$ , approximate line search)

# Why line searches?

---

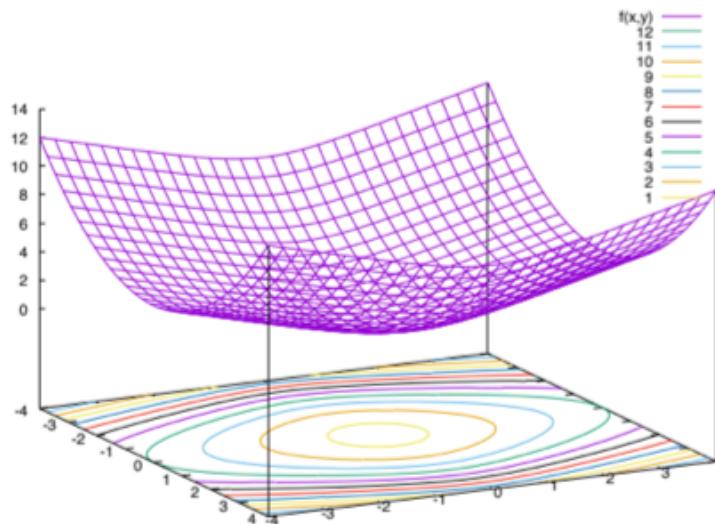


Figure 1: The function  $F(w_1, w_2) = \frac{1}{2}w_1^2 + \log(e^{w_2} + e^{-w_2})$ .

$$\nabla F(w_1, w_2) = \begin{bmatrix} w_1 \\ \tanh(w_2) \end{bmatrix} \quad \nabla^2 F(w_1, w_2) = \begin{bmatrix} 1 & 0 \\ 0 & \cosh(w_2)^{-2} \end{bmatrix}$$

# Why line searches?

---

$$\nabla F(w_1, w_2) = \begin{bmatrix} w_1 \\ \tanh(w_2) \end{bmatrix} \quad \nabla^2 F(w_1, w_2) = \begin{bmatrix} 1 & 0 \\ 0 & \cosh(w_2)^{-2} \end{bmatrix}$$

Following Bottou et al. [2016, §6.5], assume we are optimizing this function with starting point  $(3, 3)$ . The first update moves the current point along direction  $-\nabla F \approx [-3, -1]$  which unfortunately points slightly away from the optimum  $(0, 0)$ . Rescaling with the inverse Hessian yields a substantially worse direction  $-(\nabla^2 F)^{-1} \nabla F \approx [-3, -101]$ . The large second coefficient calls for a small stepsize. Using stepsize  $\gamma \approx 0.03$  moves the current point to  $(2.9, 0)$ . Although the new gradient  $\nabla F \approx [-2.9, 0]$  points directly towards the optimum, the small stepsize that was necessary for the previous update is now ten times too small to effectively leverage this good situation.

## Standard solutions for $\hat{h}_t$

- Idea: estimate a cost function using the training set. We do not know how to perform line search without computing the cost function over all the examples.
  - Meanwhile, trying to do without line search tends to perform poorly.
- $g(w_t) + \sum \lambda_i d_{t-i}$  such that  $d_t^T H d_{t-i} = 0$
- Conjugate search directions.
- Conjugate gradient (k = 1, exact line search)
- LBFGS ( $k > 1$ , approximate line search)

# About diagonal rescaling applied to neural nets

---

JEAN LAFOND, NICOLAS VASILACHE, LÉON BOTTOU  
FACEBOOK AI RESEARCH, NEW YORK

# Motivation

---

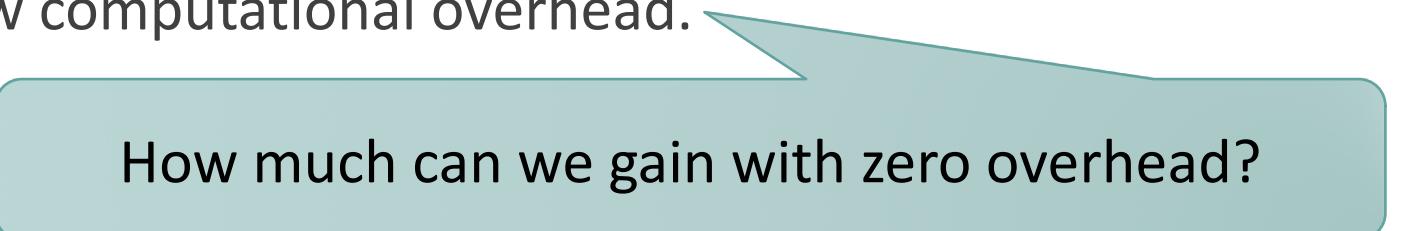
- Lots of stochastic gradient (SG) variants for training neural nets.
  - with **heuristic justifications**,
  - and **hard-to-interpret experiments** because baselines are different.
- How would they compare with a “*well tuned stochastic gradient*” ?
- What is a “*well tuned stochastic gradient*” ?
  - ❑ Initialization tricks.
  - ❑ Different stepsizes in different layers/units.

Diagonal rescaling:  $w_{t+1} = w_t - \gamma_t D_t g(w_t, \xi_t)$  with  $D_t$  diagonal

# Alternate Motivation

---

- The convergence rate of SG algorithms already match the lower bounds for optimization algorithms using noisy first order oracles [Agarwal et al, 2011].  
(well, when both are known...)
- Improved algorithms typically reduce #iterations by a multiplicative constant ...  
... with a computational overhead that is also a multiplicative constant.
- Diagonal rescaling has very low computational overhead.



How much can we gain with zero overhead?

# Summary

---

1. A pointless remark?
2. Zero overhead reparametrization.
3. Natural gradient with a twist
4. Problems

# 1- A pointless remark?

---

# Diagonal Hessian ...

---

- If the Hessian of  $f(x, y)$  is diagonal for all  $(x, y)$  then  $f(x, y) = f_1(x) + f_2(y)$
- Proof:

$$f(x, y) = f(0,0) + \int_0^x \frac{\partial f}{\partial x}(t, 0) dt + \int_0^y \frac{\partial f}{\partial y}(x, r) dr$$

$$\frac{\partial f}{\partial y}(x, r) = \frac{\partial f}{\partial y}(0, r) + \int_0^x \frac{\partial^2 f}{\partial x \partial y}(t, r) dt = \frac{\partial f}{\partial y}(0, r)$$

# Diagonal approximation

---

- Conducting the optimization as if the Hessian of  $f(x, y)$  were diagonal  
is (locally) like optimizing a separated function  $f(x, y) = f_1(x) + f_2(y)$ .
- In general the best way to optimize a separated function  
consists of solving two separate optimization problems.  
(e.g. Newton optimization on  $f(x, y) = \frac{1}{2}x^2 + \log(e^y + e^{-y})$  does not work well.)
- Rephrasing for iterative optimization algorithms:  
update  $x$  as if  $y$  were fixed, update  $y$  as if  $x$  were fixed, ...

Isn't that what a gradient is about?

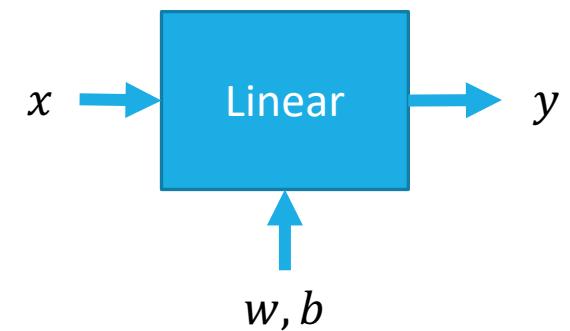
# 2- Zero overhead reparametrization

---

# A linear layer

---

$$\forall j \in \{1 \dots m\} \quad y_j = b_j + \sum_{i=1}^n x_i w_{ij}$$

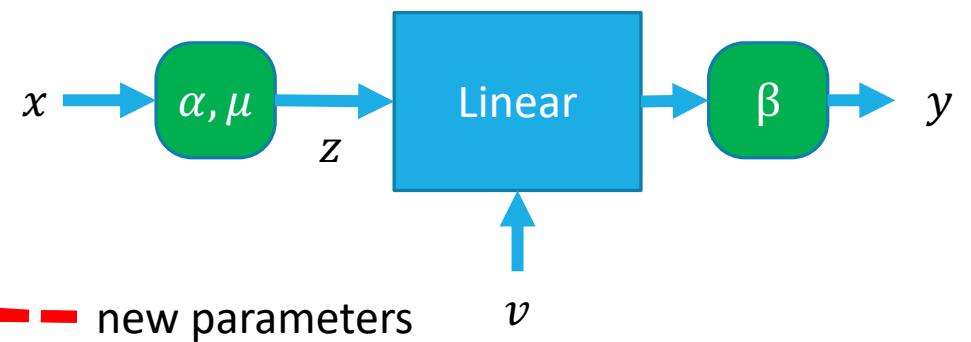


Using the notation  $g_j = \frac{\partial E}{\partial y_j}$  and the convention  $x_0 = 1$ ,

$$\forall (i, j) \in \{0 \dots n\} \times \{1 \dots m\} \quad \frac{\partial E}{\partial w_{ij}} = x_i g_j$$

# Reparametrization

$$y_j = \beta_j \left( v_{0j} + \sum_{i=1}^n \alpha_i (x_i - \mu_i) v_{ij} \right)$$



Compact notation

$$z_0 = 1 \text{ and } z_i = \alpha_i(x_i - \mu_i)$$

$$y_j = \beta_j \sum_{i=0}^n v_{ij} z_i$$

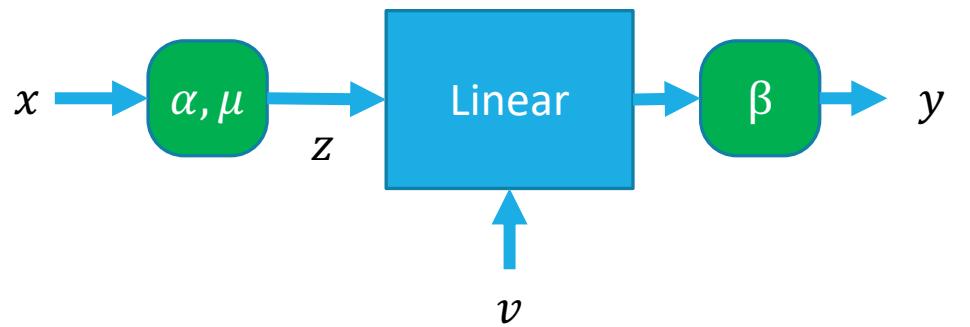
Parameter change

$$w_{ij} = \alpha_i \beta_j v_{ij} \quad (\text{for } i > 0)$$

$$b_j = \beta_j v_{0j} - \sum_{i=1}^n \mu_i w_{ij} = \beta_j v_{0j} - \sum_{i=1}^n \mu_i \alpha_i \beta_j v_{ij}$$

# SG on the new parameters

---



$$\delta v_{ij} = \left\langle \frac{\partial E}{\partial v_{ij}} \right\rangle = \langle \beta_j g_j z_i \rangle$$

Angle brackets = average over minibatch

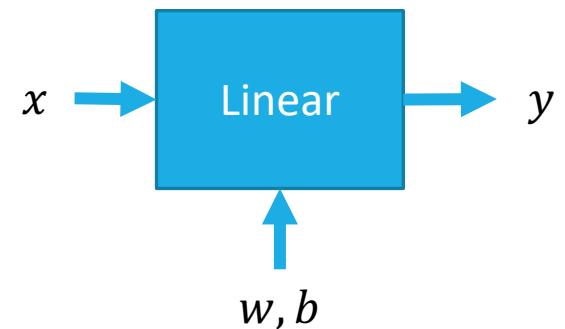
SG update is proportional to this quantity.

# ... expressed on the old parameters

---

$$\delta w_{ij} = \alpha_i \beta_j \delta v_{ij} = \langle \beta_j^2 g_j \alpha_i^2 (x_i - \mu_i) \rangle$$

$$\delta b_j = \beta_j \delta v_{0j} - \sum_{i=0}^n \mu_i \delta w_{ij} = \langle \beta_j^2 g_j \rangle - \sum_{i=1}^n \mu_i \delta w_{ij}$$



## Remarks

- This is not really diagonal but quasi-diagonal.
- Computational **overhead scales like  $n + m$** , not  $nm$ .
- We can choose  $\alpha_i$ ,  $\mu_i$ ,  $\beta_j$  freely **as long as we change them slowly**.
- How to choose them?

# 3- Natural gradient with a twist

---

# Classic natural gradient

---

1. Construct a Riemannian geometry on the optimization domain

$$d(w, w + \delta w) \approx \frac{1}{2} \delta w^T G(w) \delta w$$

2. Write steepest descent

$$w_{t+1} = w_t + \operatorname{argmin}_{\delta w} \left\{ \left\langle \frac{\partial E}{\partial w} \right\rangle^T \delta w \quad \text{s.t. } \frac{1}{2} \delta w^T G(w_t) \delta w \leq \eta_t^2 \right\}$$

3. Use a Lagrange coefficient  $1/\gamma_t$

$$w_{t+1} = w_t + \operatorname{argmin}_{\delta w} \left\{ \left\langle \frac{\partial E}{\partial w} \right\rangle^T \delta w + \frac{1}{2\gamma_t} \delta w^T G(w_t) \delta w \right\}$$

4. Solve.

$$w_{t+1} = w_t - \gamma_t G(w_t)^{-1} \left\langle \frac{\partial E}{\partial w} \right\rangle \quad [ \text{looks like (Gauss-)Newton} ]$$

# Another natural gradient...

---

- What's the difference between choosing  $\eta_t$  and choosing  $\gamma_t$  ?

$$w_{t+1} = w_t + \operatorname{argmin}_{\delta w} \left\{ \left\langle \frac{\partial E}{\partial w} \right\rangle^T \delta w \quad \text{s.t.} \quad \frac{1}{2} \delta w^T G(w_t) \delta w \leq \eta_t^2 \right\}$$

$$w_{t+1} = w_t + \operatorname{argmin}_{\delta w} \left\{ \left\langle \frac{\partial E}{\partial w} \right\rangle^T \delta w + \frac{1}{2\gamma_t} \delta w^T G(w_t) \delta w \right\}$$

- Keep solving...

$$w_{t+1} = w_t - \eta_t \frac{1}{\sqrt{\left\langle \frac{\partial E}{\partial w} \right\rangle^T G(w_t)^{-1} \left\langle \frac{\partial E}{\partial w} \right\rangle}}$$

 Same directions but different lengths.

# Singularities

---

$$w_{t+1} = w_t - \eta_t \frac{1}{\sqrt{\left\langle \frac{\partial E}{\partial w} \right\rangle^T G(w_t)^{-1} \left\langle \frac{\partial E}{\partial w} \right\rangle}} G(w_t)^{-1} \left\langle \frac{\partial E}{\partial w} \right\rangle$$

Too easily close to zero.

Two ways to improve this...

- Smooth  $\left\langle \frac{\partial E}{\partial w} \right\rangle^T G(w_t)^{-1} \left\langle \frac{\partial E}{\partial w} \right\rangle \approx \mathbb{E} \left[ \frac{\partial E}{\partial w} \right]^T G(w_t)^{-1} \mathbb{E} \left[ \frac{\partial E}{\partial w} \right] \leq \mathbb{E} \left[ \frac{\partial E}{\partial w}^T G(w_t)^{-1} \frac{\partial E}{\partial w} \right]$
- Add regularization term  $\mu > 0$

# The pointless remark matters

---

If we assume that  $G(w_t)$  is block-diagonal

$$w_{t+1} = w_t - \eta_t \frac{1}{\sqrt{\mu + \mathbb{E}\left[\frac{\partial E^T}{\partial w} G(w_t)^{-1} \frac{\partial E}{\partial w}\right]}} G(w_t)^{-1} \left\langle \frac{\partial E}{\partial w} \right\rangle$$

Should we compute this scaling coefficient globally or separately for each block ?

Computing them separately **changes the global direction** of the update...

# Sanity check : $G(w) = I$

---

- Compute scaling coefficients separately for each weight  $w_{ij}$ .
- Estimate  $\mathbb{E} \left[ \frac{\partial E}{\partial w_{ij}}^T \frac{\partial E}{\partial w_{ij}} \right]$  with a running average  $R_{ij} \leftarrow (1 - \lambda)R_{ij} + \lambda \left\langle \left( \frac{\partial E}{\partial w_{ij}} \right)^2 \right\rangle$

❑ RMSPROP (Tieleman & Hinton 2011)

$$w_{ij} \leftarrow w_{ij} - \frac{\gamma}{\sqrt{\mu + R_{ij}}} \left\langle \frac{\partial E}{\partial w_{ij}} \right\rangle$$

# Sanity check : 1990's style network

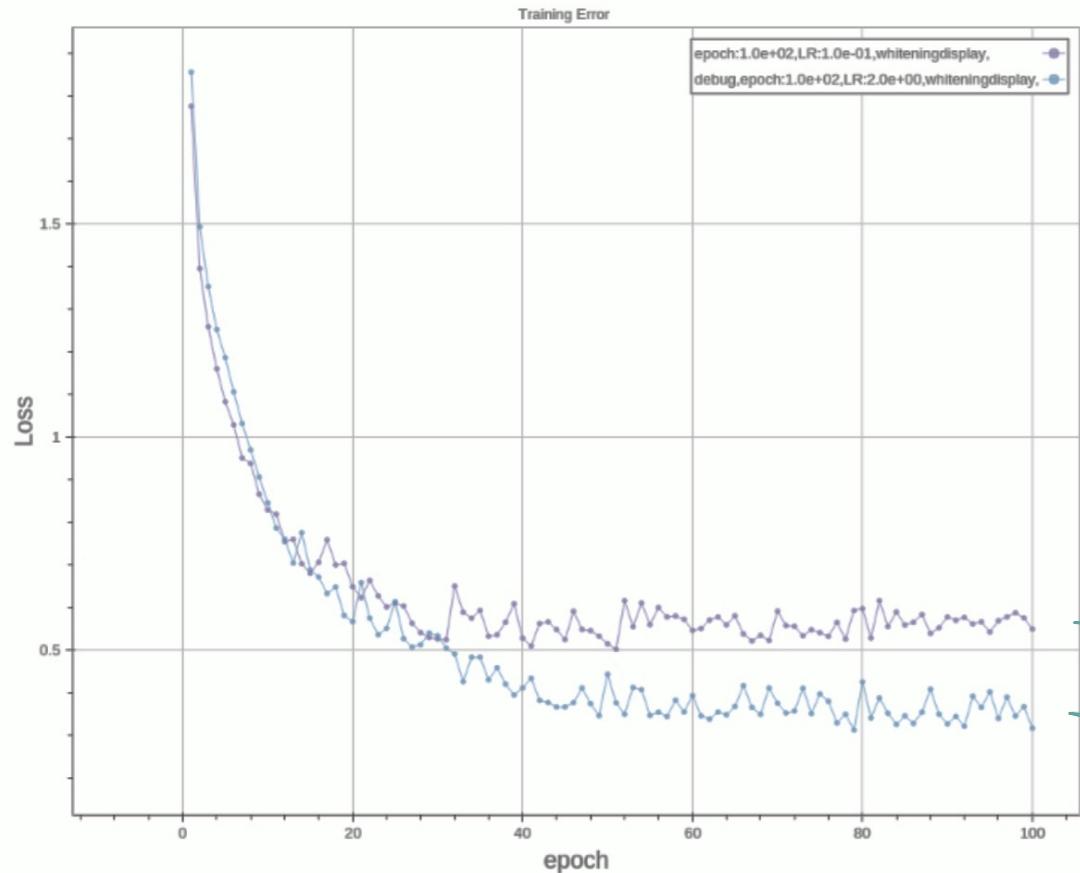
---

- Hyperbolic tangent activation function  $x_i^{(k+1)} = \tanh(x_i^{(k)}) \approx \pm 1$
- Gauss-Newton / Empirical Fisher Information :  $G(w) = \mathbb{E} \left[ \left( \frac{\partial E}{\partial w} \right) \left( \frac{\partial E}{\partial w} \right)^T \right]$
- Block diagonal approximation :  $G_j = \mathbb{E} \left[ \left( \frac{\partial E}{\partial w_{ij}} \right) \left( \frac{\partial E}{\partial w_{i'j}} \right) \right] = \mathbb{E}[g_j^2 x_i x_{i'}] \approx \mathbb{E}[g_j^2] I$
- Compute separate scaling ratio per block (= per neuron)

$$\sqrt{\mathbb{E} \left[ \frac{\partial E^T}{\partial w} G(w_t)^{-1} \frac{\partial E}{\partial w} \right]} = \sqrt{\mathbb{E} \left[ \frac{\sum_{ii'} g_j^2 x_i x_{i'}}{\mathbb{E}[g_j^2]} \right]} \approx \sqrt{fanin}$$

Well known heuristic.  
Divide stepsizes by  
 $\sqrt{share \times fanin}$

# Divide stepsizes by $\sqrt{share \times fanin}$



- Data : CIFAR 10  
60000 32x32 color images, 10 classes
- Network : CNN  
C6(5x5)-P(2x2)-C16(5x5)-F84-F16-F10
- ReLU :  $x_i^{(k+1)} = \max(y_i^{(k)}, 0)$

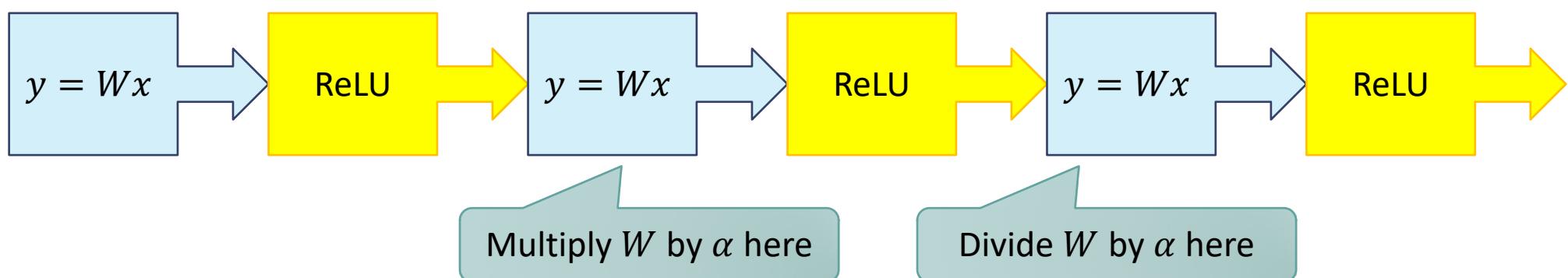
Stochastic gradient  
Common stepsize

Per-layer stepsize  
 $\propto 1/\sqrt{share * fanin}$

# Criticism

---

- With tanh activation,  $\mathbb{E}[g_j^2 x_i x_{i'}] \approx \mathbb{E}[g_j^2] I$  is a dubious approximation.
  - With ReLU activation, this is not credible.
- 
- ReLU networks  $x_i^{(k+1)} = \max(y_i^{(k)}, 0)$  have hyperbolic geometry in parameter space.



# Whitening reparametrization

---

- Make  $\mathbb{E}[g_j^2 z_i z_{i'}] \approx \mathbb{E}[g_j^2] I$  with:

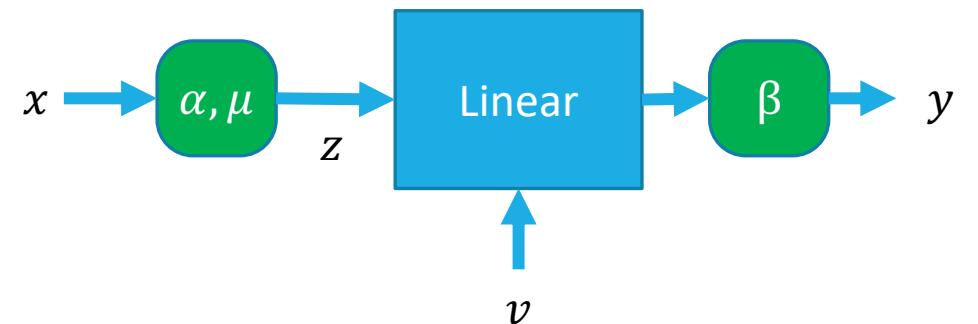
$$\mu_i = \mathbb{E}[x_i] \quad \alpha_i^2 = \frac{1}{\text{var}[x_i]}$$

estimated with running averages...

- Then take Fisher into account with:

$$\beta_j^2 = \frac{1}{\sqrt{\text{share} \times \text{fanin}}}$$

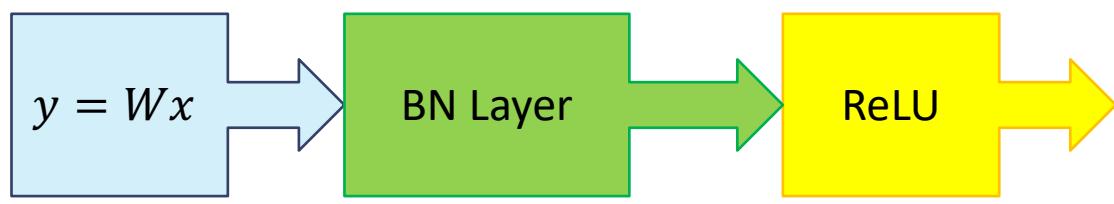
*(there are cleaner ways to get the same result)*



# Alternative : batch normalization

---

These days, everybody uses batch normalization (Ioffe & Szegedy, 2015)

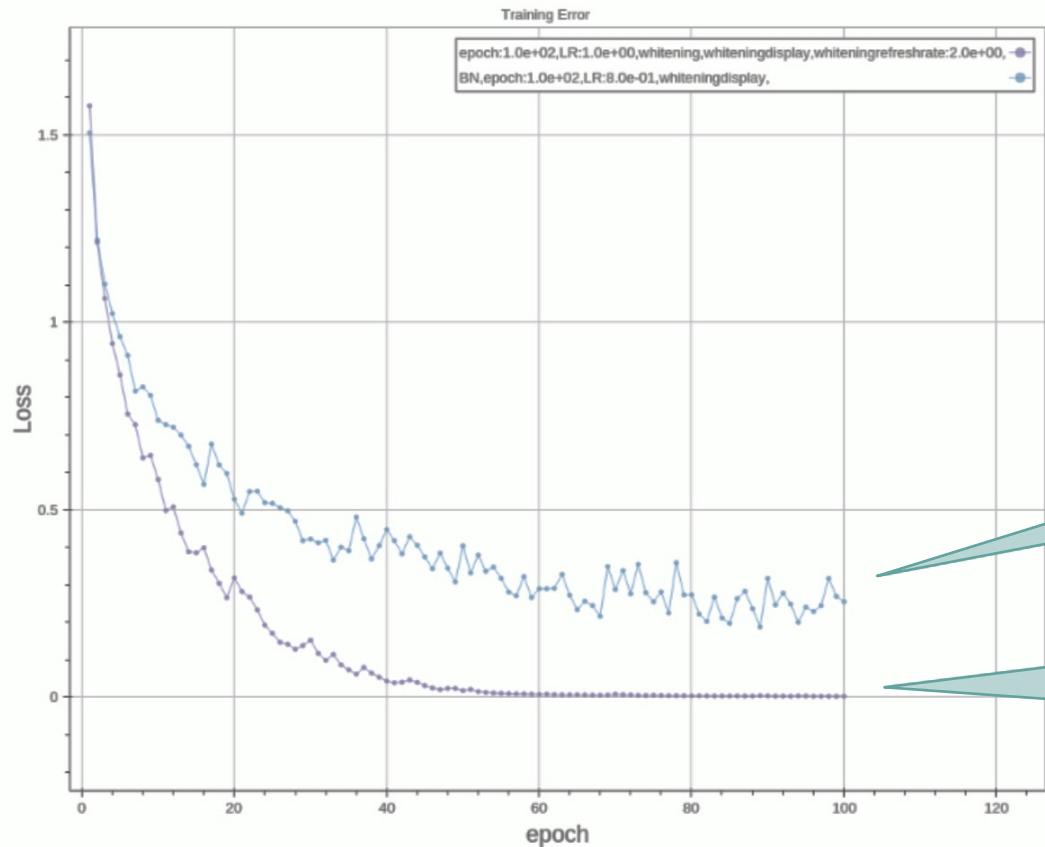


- $y_i = (y_i - \mu_i)/\sigma_i$  with  $\mu_i$  and  $\sigma_i$  computed on current minibatch
- Gradient backpropagation accounts for the computation of  $\mu_i, \sigma_i$

Batch normalization has strange side effects :

- The network output for a given pattern depends on the other minibatch examples.
- When the minibatch examples are randomly picked, this can be viewed as additional noise...

# Batch normalization vs whitening

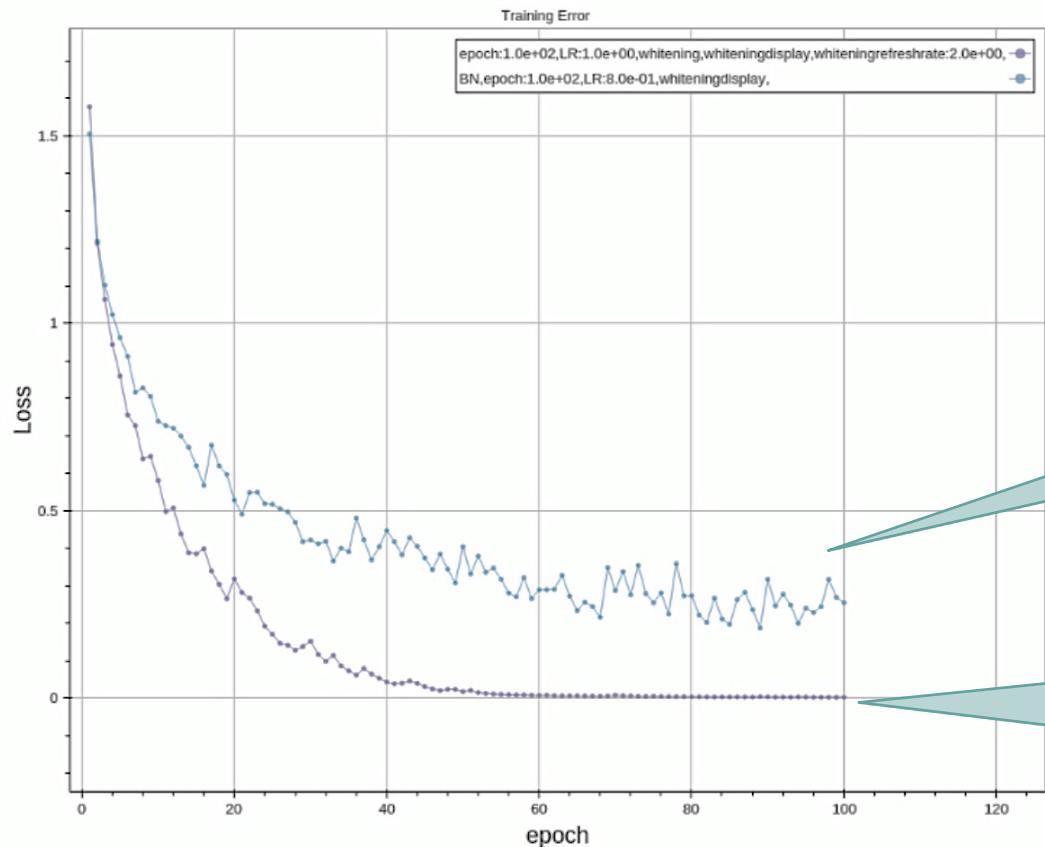


- Data : CIFAR 10  
60000 32x32 color images, 10 classes
- Network : CNN  
C6(5x5)-P(2x2)-C16(5x5)-F84-F16-F10

SG with common stepsize  
and batch normalization

Whitening reparametrization  
 $\mu_i = E(x_i)$   $\alpha_i^2 = 1/\text{Var}(x_i)$   
 $\beta_j^2 = 1/\sqrt{\text{share} * \text{fanin}}$

# Batch normalization vs whitening



Network structure was copied from Soumith Chintala's Torch tutorial.

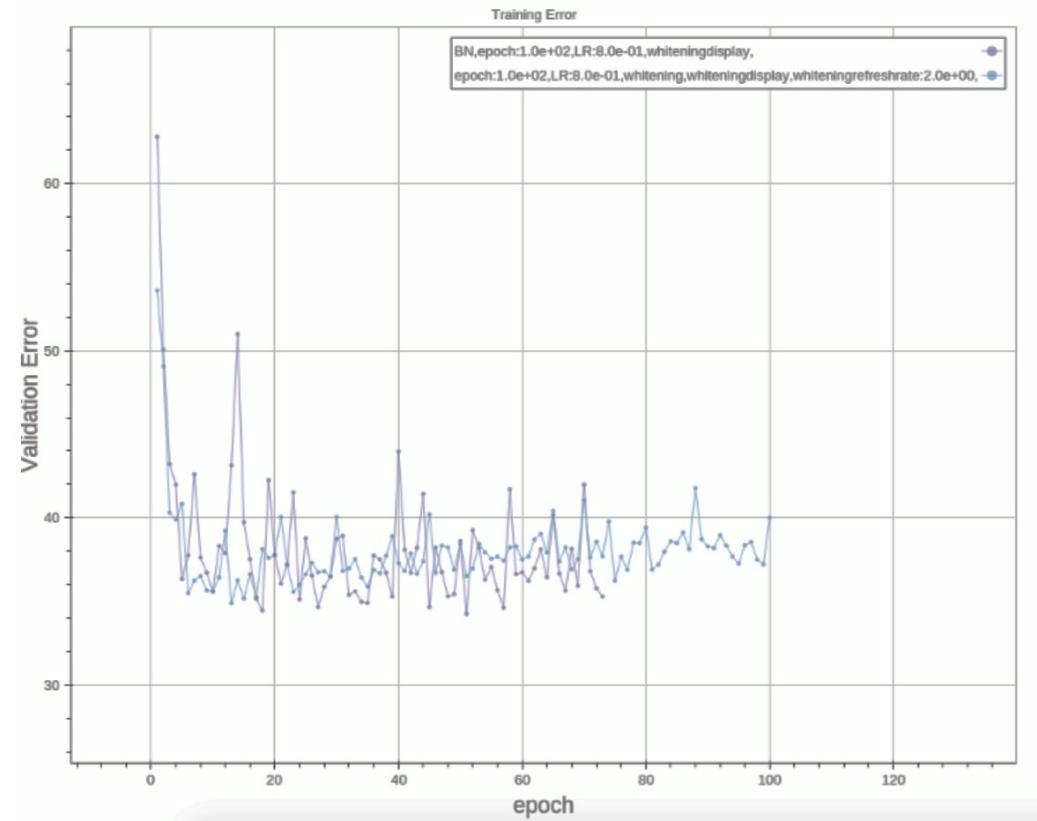
Batch normalization seem stuck.  
(normal SG was stuck even higher)

Whitening quickly reaches zero training cost!  
(overparametrized network)

# Test error sanity check

---

- Whitening optimizes better  
→ but also overfits!
- We could reduce the network size  
and train even faster...  
→ but how to compare?
- Sanity check:  
Use the same network size  
with L2 regularization  
→ No longer overfits.



# 4- Problems

---

# Scaling up to ImageNet

---

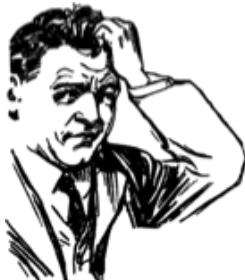
- Data: ImageNet – 1 million 224x224 color images, 1000 classes
- Network: Alexnet (Krizhevsky et al, 2012)

First try 10% of ImageNet (100000 randomly picked examples, 1000 classes)

❑ Same result as CIFAR10: Whitening quickly reaches near zero training cost. BN stuck higher.

Then try 100% of ImageNet (1M examples)

❑ Whitening now trains slower than batch-normalization. Slower than plain SG too!



*"This is very strange! With SG, we do not usually expect training error convergence to depend on the training set size."*

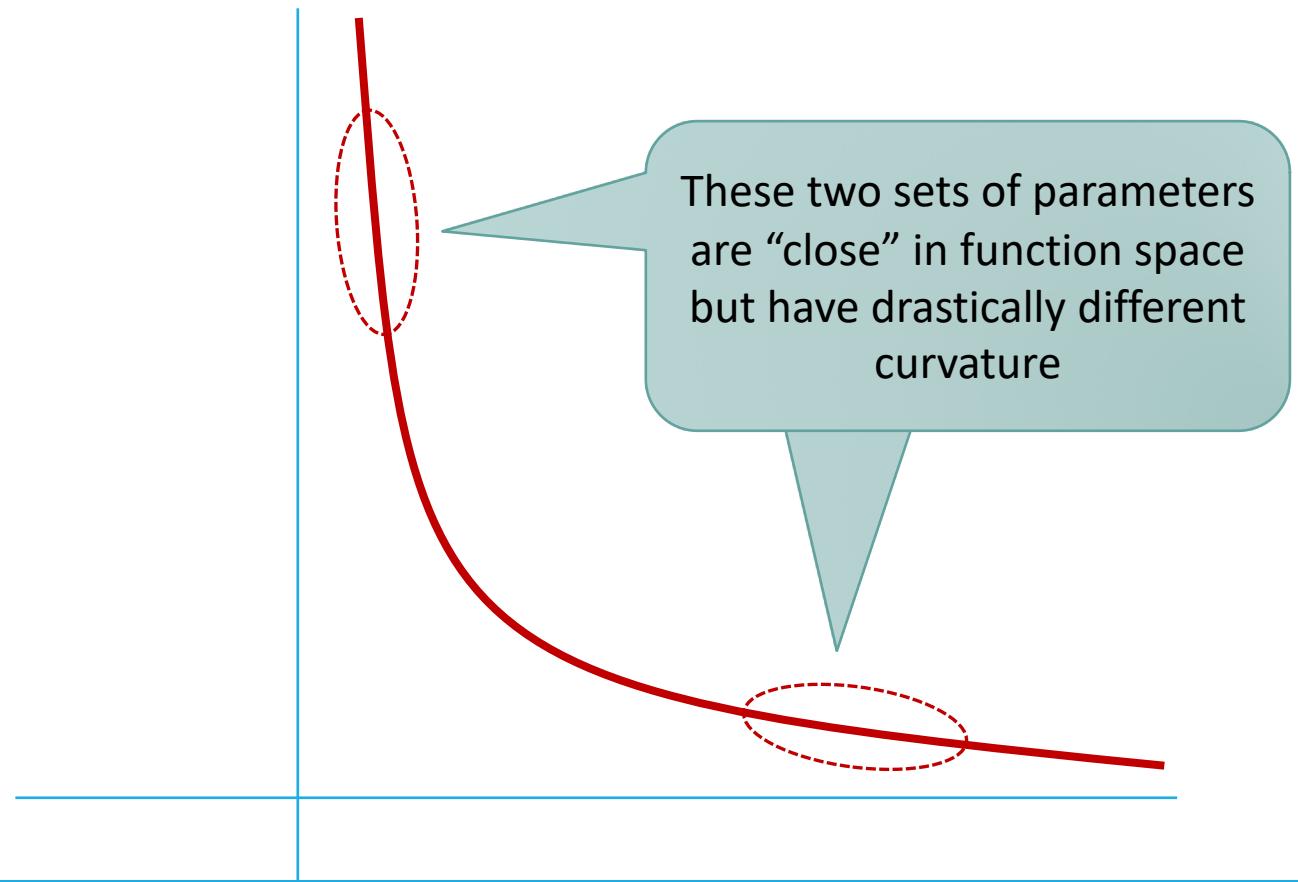
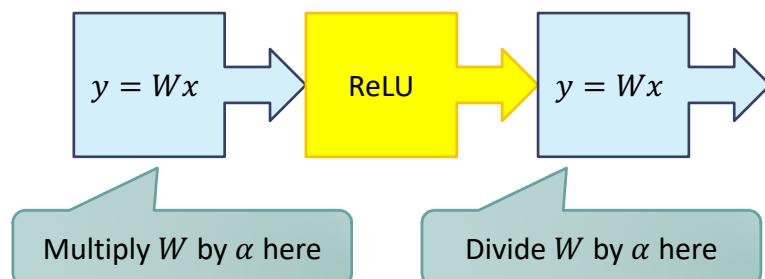
# Investigation

---

- Good old investigation method : look at the numbers during training.  
This is hard when the network has 60M weights.
- Key observation : during the first epochs,  $E(x_i)$  and  $Var(x_i)$  can change **very quickly** !
  - ❑ The running estimates of  $E(x_i)$  and  $Var(x_i)$  are behind.  
They are sometimes so wrong that we make a big unwarranted update  
that takes the  $y_j$  in places where the ReLU has no gradient... (plateau)
  - ❑ Speeding up the running estimates becomes unstable.
  - ❑ BN does not have this problem because it gets  $E(x_i)$  and  $Var(x_i)$  on the current minibatch.

# Why do statistics change so quickly?

The parameters of a ReLU networks  
have hyperbolic geometry

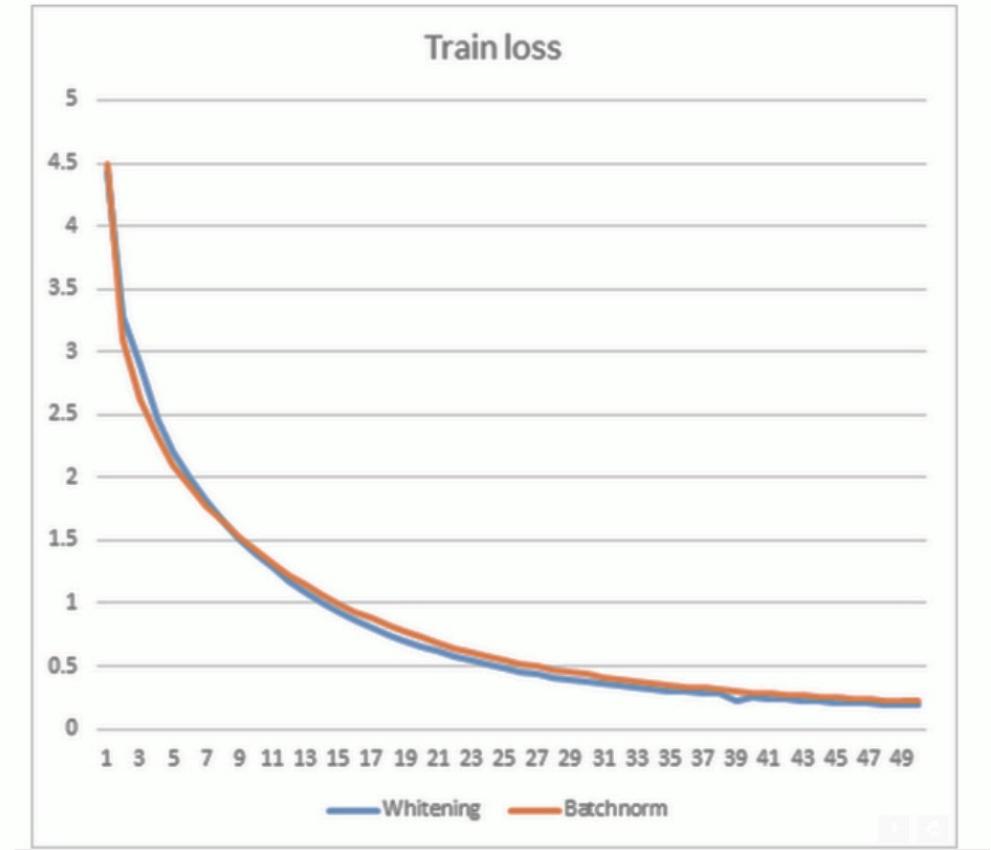


# Mitigation

---

- Statistics **seem** more stable after 2 epochs.
- Running 2 BN epochs before whitening fixes the problem →
- Same speed as BN in iterations, but 20-25% faster in compute time (whitening has less overhead!)
- Looking for a better fix.

*Work in progress...*



# Conclusion

---

# Conclusions

---

- We were looking for something *robust, efficient, and principled*.
- This is still work in progress but...
- ❑ We now understand why RMSPROP works.
- ❑ We now understand why Batch-Normalization works.