



# MORE ON STOCHASTIC GRADIENT

Setup:

- minimize  $f(w) = \frac{1}{m} \sum_{i=1}^m f_i(w)$   $w \in \mathbb{R}^d$
- $i=1..m$ ,  $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$  that depends on the  $i^{\text{th}}$  element in a dataset of  $m$  elements

- Every  $f_i$  is  $C^1$ , and evaluating  $\nabla f_i$  (or  $f_i$ ) requires to access the corresponding data point
- $m \gg 1$  so that evaluating  $\nabla f = \frac{1}{m} \sum_{i=1}^m \nabla f_i$  is expensive

GD:  $w_{k+1} = w_k - \alpha_k \nabla f(w_k)$  : expensive to use in that setting

SG  
(Stochastic Gradient)  $w_{k+1} = w_k - \alpha_k \nabla f_{i_k}(w_k)$  where  $i_k \in \{1, \dots, m\}$  drawn at random

## ① More on convergence rates

↳ SG has worse convergence rates than GD

Ex) If  $f_i$  is convex  $\forall i$ , then

• After  $K$  iterations of GD,  $f(w_k) - \min_{w \in \mathbb{R}^d} f(w) \leq O\left(\frac{1}{K}\right)$

• After  $K$  iterations of SG,  $\mathbb{E}\left[f(w_k) - \min_{w \in \mathbb{R}^d} f(w)\right] \leq O\left(\frac{1}{\sqrt{K}}\right)$

⚠ The number of iterations is not a good metric for comparing GD and SG

⇒ We care about the number of accesses to data points

⇒ 1 iteration of GD = 1  $\nabla f$  =  $m$  accesses to data points

1 iteration of SG =  $\nabla f_{i_k}$  = 1 access to 1 data point

Def: An epoch is a cost unit corresponding to  $m$  accesses to data points in a dataset of size  $m$ .

↳ 1 iteration of GD costs 1 epoch

↳ 1 iteration of SG costs  $\frac{1}{m}$  epoch

⇒ 1 epoch is the cost of  $m$  iterations of SG

↳ If we provide convergence rates in terms of epochs and not in terms of iterations, then SG will have better guarantees for  $m \gg 1$

Ex) If  $f$  is convex  $\forall i$ , then

• After  $N_E$  epochs, GD will have done  $N_E$  iterations and so

$$f(w_{N_E}) - \min_{w \in \mathcal{M}} f(w) \leq O\left(\frac{1}{N_E}\right) \quad (K=N_E \text{ above})$$

• After  $N_E$  epochs, SG will have done  $m \times N_E$  iterations and so

$$\mathbb{E}\left[f(w_{mN_E}) - \min_{w \in \mathcal{M}} f(w)\right] \leq O\left(\frac{1}{\sqrt{m} \sqrt{N_E}}\right) \quad (K=mN_E \text{ above})$$

For  $m \gg 1$ ,  $\frac{1}{\sqrt{m} \sqrt{N_E}} \ll \frac{1}{N_E}$  : Better rate for SG!

## ② More on randomness

↳ SG is by design a randomized method

⇒ Several runs of SG can give very different results "Variance of SG"

⇒ But on average the method works (see convergence rates)

↳ Variants of SG have been developed to guarantee that the behavior of any run is closer to the average than for "Vanilla" SG ("Vanilla" = basic)

## Batch stochastic gradient

cardinality (number of elements) in  $S_k$

$$w_{k+1} = w_k - \alpha_k \left( \frac{1}{|S_k|} \sum_{i \in S_k} \nabla f_i(w_k) \right)$$

average of all  $\nabla f_i$ 's corresponding to indices in  $S_k$

with replacement:  
 $S_k$  can contain multiple copies of the same index

where  $S_k$  is a set of indices drawn randomly in  $\{1, \dots, m\}$  (with or without replacement)

- SG is a special case of batch SG for which  $|S_k| = 1$
- GD is a special case of batch SG  $|S_k| = m$  and indices drawn without replacement  $\Rightarrow S_k = \{1, \dots, m\}$
- Case of interest: mini-batch SG  $1 < |S_k| < m \quad \forall k$

$\hookrightarrow$  In general, using a batch of more than 1 data point reduces the variance of the algorithm (compared to vanilla SG) and it can even improve the performance for  $|S_k| > 1$

$\hookrightarrow$  But every iteration of batch SG is more expensive than an iteration of vanilla SG, so finding the right value for the batch size ( $|S_k|$ ) is not straightforward.

NB: There are other ways to reduce the variance by changing the stochastic gradient iteration

## (3) Advanced SG methods

Generic iteration defined coordinate wise

$$\forall j=1..d, \quad [w_{k+1}]_j = [w_k]_j - \alpha_k \frac{[m_k]_j}{[v_k]_j}$$

$\alpha_k > 0$

$j^{\text{th}}$  coordinate of  $m_k$

$\mathbb{R}^d \ni m_k$ : direction (depends on a stochastic gradient, or a batch of stochastic gradients)

$\mathbb{R}^d \ni v_k$ : vector of normalization for the various coordinates

$$\text{SG} : m_k = \nabla f_{i_k}(w_k), \quad v_k = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

### a) SG with momentum

$$v_k = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \text{but } m_k = \beta_1 m_{k-1} + (1-\beta_1) \nabla f_{i_k}(w_k)$$

$\beta_1 \in (0,1)$

→ Combines previous step with current stochastic gradient ( $\approx$  heavy ball)

→ Was used for the groundbreaking results on neural network performance in 2012

### b) AdaGrad ( $\approx 2014$ )

$$m_k = \nabla f_{i_k}(w_k) \quad [v_k]_j = \sqrt{\sum_{l=0}^k [\nabla f_{i_l}(w_l)]_j^2}$$

→ Normalize every coordinate according to the value of all previous stochastic gradients

→ Useful for problems with sparse gradients (lots of zero coefficients), used in recommender systems

### c) Adam (2015)

! "Combines the ideas of AdaGrad and SG with momentum"  
change  $v_h$                       change  $m_h$

→ Default training algorithm in many deep learning tasks  
especially MLP

$$m_k = (1 - \beta_1) \sum_{l=0}^k \beta_1^{k-l} \nabla \ell(w_l)$$

$$[v_h]_j = \sqrt{\frac{(1 - \beta_2) \sum_{l=0}^h \beta_2^{h-l} [\nabla \ell(w_l)]_j^2}{1 - \beta_2^{h+1}}}$$