

# OPTIMIZATION FOR MACHINE LEARNING - STOCHASTIC GRADIENT

October 12, 2023

Program :      Today: Basics of SG  
Oct 19 : Theory of SG methods  
Oct 23 : Advanced SG methods

# INTRODUCTION TO STOCHASTIC GRADIENT

→ So far : minimize  $f(x)$   $x \in \mathbb{R}^d$   $f \in C^1 \Rightarrow$  can query a gradient  $\nabla f$  at every point

Gradient descent

$$x_{k+1} = x_k - \alpha_k \nabla f(x_k) \quad \alpha_k > 0$$

→ In ML,  $f$  generally has a specific structure that depends on some data

$\mathcal{X}$  : input space  $\mathcal{Y}$  : output space

Goal: Compute a mapping between  $\mathcal{X}$  and  $\mathcal{Y}$  such that for a given data distribution  $\mathcal{D}$  on  $(\mathcal{X}, \mathcal{Y})$ , the mapping maps  $x$  to  $y$  for any  $(x, y) \sim \mathcal{D}$   
 $\Rightarrow$  Actual task: find the best possible mapping

Quantify how good a mapping is at fitting the data

• Expected risk : For a mapping  $h: \mathcal{X} \rightarrow \mathcal{Y}$ ,

$$R(h) = \underbrace{P}_{\text{probability}}((x, y) \sim \mathcal{D}) [h(x) \neq y] = \underbrace{E}_{\text{expected value}} [\underbrace{11}_{\text{indicator function}}(h(x) \neq y)]$$

$$11(E) = \begin{cases} 1 & \text{if } E \text{ happens} \\ 0 & \text{otherwise} \end{cases}$$

$\Rightarrow$  A natural optimization problem

minimize  $R(h)$   $h \in \mathcal{H}$

where  $\mathcal{H}$  is some function space

↳ In general, one cannot design algorithms to solve this problem in its generic:

↳ Modifications

•  $\mathcal{H}$  can be hard to optimize over  $\Rightarrow$  parameterize

$$h \in \mathcal{H} \Rightarrow h \in \{h(\cdot; x) \mid x \in \mathbb{R}^d\}$$

•  $\mathbb{E}_{(a,y) \sim \mathcal{D}}[\cdot]$  not known in practice  $\Rightarrow$  approximate with samples

$$\mathbb{E}_{(a,y) \sim \mathcal{D}}[\cdot] \Rightarrow \frac{1}{m} \sum_{i=1}^m (\cdot)$$

$\uparrow$   
 $(a_i, y_i) \sim \mathcal{D}$

•  $\mathbb{1}(h(a) \neq y)$  "0-1 loss" is not easy to optimize  
 $\Rightarrow$  replace with a continuous loss (typically convex)

$$\mathbb{1}(h(a) \neq y) \Rightarrow l(h(a), y)$$

$l: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$

Ex)  $l(h, y) = \|h - y\|^2$   
 $l(h, y) = \|h - y\|_1$

With these changes, the problem becomes

minimize  $x \in \mathbb{R}^d$

Variables  $\rightarrow$

$$\frac{1}{m} \sum_{i=1}^m l(h(a_i; x), y_i)$$

$\uparrow$   
Average over data samples

$\uparrow$   
 $f_i(x)$ : loss at the data point  $(a_i, y_i)$   
"Empirical Risk Minimization"

We set

$$f(x) = \frac{1}{m} \sum_{i=1}^m f_i(x) = \frac{1}{m} \sum_{i=1}^m l(h(a_i; x), y_i)$$

↳ If every  $f_i$  is  $C^1$ , then  $f$  is  $C^1$  and

$$\forall x \in \mathbb{R}^d, \quad \nabla f(x) = \frac{1}{n} \sum_{i=1}^m \nabla f_i(x)$$

and so gradient descent (GD) performs the iteration

$$x_{k+1} = x_k - \alpha_k \nabla f(x_k) = x_k - \frac{\alpha_k}{n} \sum_{i=1}^m \nabla f_i(x_k)$$

⚠ In a "big data" setting, computing  $\nabla f(x_k)$  is **expensive** because it involves computing all  $\nabla f_i(x_k)$  for  $i=1 \dots m$ , i.e. going through the entire dataset.  
 **$m$  Samples**  $\{(a_1, y_1), \dots, (a_n, y_n)\}$

In addition, it may be possible to compute a better point than  $x_k$  without looking at all data points.

## ① STOCHASTIC GRADIENT ALGORITHM

↳ Aka stochastic gradient descent (SGD)

Setup: minimize  $f(x) = \frac{1}{n} \sum_{i=1}^m f_i(x)$        $f_i: \mathbb{R}^d \rightarrow \mathbb{R}_{C^1}$   
 $x \in \mathbb{R}^d$   
"Finite-sum problem"

Key assumption: Every value of  $f_i$  or its gradient  $\nabla f_i$  is computed using the  $i$ th sample of a dataset with  $n$  elements

# Stochastic Gradient iteration

(Robbins & Monro, 1954)

$$\forall k \geq 0, \quad x_{k+1} = x_k - \alpha_k \nabla f_{i_k}(x_k) \quad \alpha_k > 0$$

$\hookrightarrow$  "stochastic gradient"

where  $i_k$  is drawn randomly  
in  $\{1, \dots, m\}$

$\hookrightarrow$  One iteration of SG (Stochastic Gradient) accesses  
1 data point only  
 $\Rightarrow$   $m$  times cheaper than one iteration  
of GD (which accesses  $m$  data points)

More general: Batch stochastic gradient

Iteration:  $x_{k+1} = x_k - \alpha_k \left( \frac{1}{|S_k|} \sum_{i \in S_k} \nabla f_i(x_k) \right)$   $\alpha_k > 0$

$\xrightarrow{\text{cardinality of } S_k}$

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

$\frac{1}{|S_k|} \sum_{i \in S_k} \nabla f_i(x_k)$ : "Batch stochastic gradient"

- $|S_k| = 1 \Rightarrow$  SG
- $|S_k| = m$  and draws without replacement  $\Rightarrow$  GD aka batch GD

$$\frac{1}{|S_k|} \sum_{i \in S_k} \nabla f_i(x_k) = \frac{1}{m} \sum_{i=1}^m \nabla f_i(x_k) = \nabla f(x_k)$$

For fixed batch size  $|S_k| = m_b \forall k$ , we distinguish two  
regimes outside of  $m_b = 1$  and  $m_b = m$ :

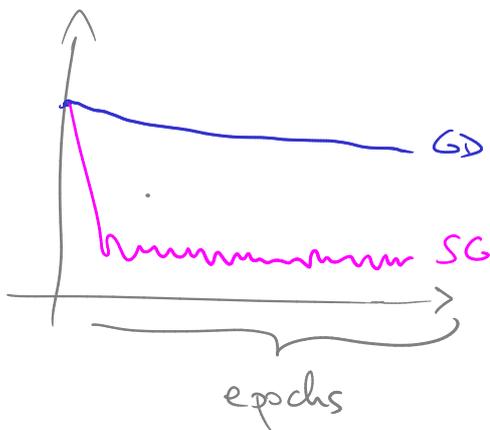
- $m_b \gg 1$  : large batch regime, per-iteration cost is (much) larger than SG  
 (much) larger than SG  
 $\Rightarrow m_b$  accesses to data points per iteration

$\rightarrow$  In practice, can be similar to GD in terms of speed of convergence

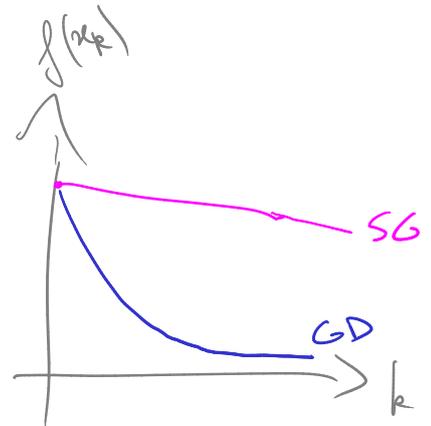
- $1 < m_b \ll m$  : mini-batch regime
  - per-iteration cost larger than SG but much smaller than GD
  - practical performance can be better than that of SG

### ③ Comparing SG, GD and batch SG

Typical plot



Not the plot you want



$\hookrightarrow$  Iteration count is not the right metric to compare algorithms for solving our problem of interest

$\hookrightarrow$  A better metric: Number of epochs

Def. An epoch is the cost of  $m$  accesses to a point in a dataset of size  $m$ .

1 iteration of GD = 1  $\nabla f(\cdot) = \nabla f_{i_1}, \dots, \nabla f_{i_m} = m$  accesses to data points

1 iteration of GD costs 1 epoch

1 iteration of SG costs  $\frac{1}{n}$  epoch

1 iteration of batch SG with fixed batch size  $m_b$  costs  $\frac{m_b}{n}$  epoch(s)

## ④ Towards theory for SG

We consider again

$$\text{minimize}_{x \in \mathbb{R}^d} f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$$

Assumptions:

•  $f_i$  is  $C^1 \quad \forall i \in \{1, \dots, n\}$

•  $f$  is  $C^{1,1}_L$ , i.e.  $C^1$  with  $L$ -Lipschitz continuous gradient

$\nabla f$  is  $L$ -Lipschitz continuous with  $L > 0$

$$\Leftrightarrow \forall (x, w) \in (\mathbb{R}^d)^2, \quad \|\nabla f(x) - \nabla f(w)\| \leq L \|x - w\|$$

Proposition: Under these assumptions, for any  $(x, w) \in (\mathbb{R}^d)^2$ ,

$$\boxed{f(w) \leq f(x) + \nabla f(x)^T (w - x) + \frac{L}{2} \|w - x\|^2} \quad (1)$$

Apply (1) with  
↳ For GD,  $\forall x = x_k$  and  $w = x_k - \alpha_k \nabla f(x_k)$ ,

$$f(x_k - \alpha_k \nabla f(x_k)) \leq f(x_k) - \alpha_k \|\nabla f(x_k)\|^2 + \frac{L}{2} \alpha_k^2 \|\nabla f(x_k)\|^2$$

• Allows to identify good values for  $\alpha_k$  that guarantee  $f(x_k - \alpha_k \nabla f(x_k)) < f(x_k)$

• One of the main ingredients to show asymptotic

or non-asymptotic convergence properties of GD

↳ For SG, we apply (1) with  $x = x_k$  and  $w = x_k - \alpha_k \nabla f_{i_k}(x_k)$

$$\begin{aligned} f(x_k - \alpha_k \nabla f_{i_k}(x_k)) &\leq f(x_k) + \nabla f(x_k)^\top (x_k - \alpha_k \nabla f_{i_k}(x_k) - x_k) \\ &\quad + \frac{L}{2} \|x_k - \alpha_k \nabla f_{i_k}(x_k) - x_k\|^2 \\ &= f(x_k) - \alpha_k \nabla f(x_k)^\top \nabla f_{i_k}(x_k) \\ &\quad + \frac{L}{2} \|\alpha_k \nabla f_{i_k}(x_k)\|^2 \\ &= f(x_k) - \alpha_k \nabla f(x_k)^\top \nabla f_{i_k}(x_k) + \frac{L\alpha_k^2}{2} \|\nabla f_{i_k}(x_k)\|^2 \end{aligned}$$

$$(2) \quad f(x_k - \alpha_k \nabla f_{i_k}(x_k)) \leq f(x_k) - \alpha_k \nabla f(x_k)^\top \nabla f_{i_k}(x_k) + \frac{L\alpha_k^2}{2} \|\nabla f_{i_k}(x_k)\|^2$$

inequality between random variables that depend (in particular) on  $i_k$

Assumption: The random indices  $\{i_k\}_k$  are independent

Under that assumption, we take the expectation with respect to  $i_k$  on both sides of (2)

$$\begin{aligned} \mathbb{E}_{i_k} [f(x_k - \alpha_k \nabla f_{i_k}(x_k))] &\leq \mathbb{E}_{i_k} [f(x_k) - \alpha_k \nabla f(x_k)^\top \nabla f_{i_k}(x_k) + \frac{L}{2} \alpha_k^2 \|\nabla f_{i_k}(x_k)\|^2] \\ &= \mathbb{E}_{i_k} [f(x_k)] - \alpha_k \mathbb{E}_{i_k} [\nabla f(x_k)^\top \nabla f_{i_k}(x_k)] \\ &\quad + \frac{L}{2} \alpha_k^2 \mathbb{E}_{i_k} [\|\nabla f_{i_k}(x_k)\|^2] \end{aligned}$$

(Here we assume that  $\alpha_k$  is independent of  $i_k$ , e.g.  $\alpha_k = \frac{1}{L}$ )

independent of  $i_k$  by assumption

$x_k$  is completely determined by the values of  $i_0, i_1, \dots, i_{k-1}$

$$\Rightarrow \mathbb{E}_{i_k} [f(x_k)] = f(x_k)$$

$$\Rightarrow \mathbb{E}_{i_k} [\nabla f(x_k)^T u] = \nabla f(x_k)^T \mathbb{E}_{i_k} [u] \quad \forall u \in \mathbb{R}^d$$

Overall, we obtain

$$(3) \quad \mathbb{E}_{i_k} [f(x_k - \alpha_k \nabla f_{i_k}(x_k))] \leq f(x_k) - \alpha_k \nabla f(x_k)^T \mathbb{E}_{i_k} [\nabla f_{i_k}(x_k)] + \frac{L\alpha_k^2}{2} \mathbb{E}_{i_k} [\|\nabla f_{i_k}(x_k)\|^2]$$

- Counterpart of the GD inequality
- used to prove  $\mathbb{E}_{i_k} [f(x_{k+1})] \leq f(x_k)$  under the appropriate assumptions  
 $\Rightarrow$  Descent (decrease in  $f$ )  
or average

Assumptions on  $\nabla f_{i_k}(x_k)$

At every iteration  $k$ ,

$$a) \quad \mathbb{E}_{i_k} [\nabla f_{i_k}(x_k)] = \nabla f(x_k)$$

"On average, the stochastic gradient behaves like the true gradient"

$$b) \quad \mathbb{E}_{i_k} [\|\nabla f_{i_k}(x_k)\|^2] \leq \|\nabla f(x_k)\|^2 + \sigma^2 \quad \begin{matrix} \sigma < \infty \\ \text{with } \sigma \geq 0 \end{matrix}$$

$$\stackrel{\text{under a1}}{\Leftarrow} \mathbb{E}_{i_k} [\|\nabla f_{i_k}(x_k)\|^2 - \|\mathbb{E}_{i_k} [\nabla f_{i_k}(x_k)]\|^2] \leq \sigma^2 \rightarrow \text{Noise/Variance level}$$

"The norm of the stochastic gradient does not vary too much from the true gradient norm"

Using Assumptions (a) and (b) in (3), we obtain

$$\mathbb{E}_h \left[ f(x_k - \alpha_k \nabla f_k(x_k)) \right] \leq \underbrace{f(x_k) - \alpha_k \|\nabla f(x_k)\|^2}_{\substack{\downarrow \\ \text{what we had} \\ \text{for gradient descent}}} + \underbrace{\frac{L\alpha_k^2}{2} \|\nabla f(x_k)\|^2}_{\substack{\downarrow \\ \text{accounts for} \\ \text{the uncertainty in} \\ \text{choosing it}}} + \frac{L\alpha_k^2}{2} \sigma^2$$