

# OPTIMIZATION FOR MACHINE LEARNING

October 2, 2023

Outline:            Lab on Acceleration (done)  
                          Lecture part: Stochastic gradient

⚠ Thursday: 10:15am - 1:30 pm (they 1 month break)

# STOCHASTIC GRADIENT METHODS

↳ So far we looked at generic problems

$$\underset{w \in \mathbb{R}^d}{\text{minimize}} \quad f(w) \quad f \in C^1$$

and their solving using gradient descent

↳ In ML, most optimization problems have a specific structure that can be exploited to design better algorithms

Typical form of an optimization problem in ML

$$\underset{w \in \mathbb{R}^d}{\text{minimize}} \quad f(w) := \frac{1}{n} \sum_{i=1}^n f_i(w)$$

↳ Finite sum structure

where  $n$  is the number of data points (or samples)  
and  $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$  is the loss function that depend on data point  $i$  in the dataset

Usually (especially in deep learning), but not always,  $n$  is very large

(ex) Linear least squares

$$f(w) = \frac{1}{2n} \|Xw - y\|^2 = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (x_i^T w - y_i)^2$$

Data set :  $\{(x_1, y_1), \dots, (x_n, y_n)\}$

$$f_i(w) = \frac{1}{2} (x_i^T w - y_i)^2$$

↳ Assuming that the  $f_i$ s are  $C^1$ , then  $f$  is  $C^1$  and we can apply gradient descent

$$w_{k+1} = w_k - \alpha_k \nabla f(w_k) \quad (\text{GD iteration})$$

$$= w_k - \alpha_k \times \frac{1}{m} \sum_{i=1}^m \nabla f_i(w_k)$$

↑ Gradient of  $f_i$  with respect to  $w$

Observations

- One iteration of GD involves computing  $\nabla f(w_k)$ , which itself involves computing  $\nabla f_1(w_k), \dots, \nabla f_m(w_k)$
- Each  $\nabla f_i(w_k)$  depends on data point  $i$  in the dataset
- ⇒ Computing  $\nabla f_i(w_k)$  involves one access to that data point
- Thus, one iteration of GD involves accessing the entire dataset
- ⇒ Expensive when  $m$  is large

↳ In what follows, we assume that  $m$  is too large for GD to be practical

## ① Stochastic gradient method

Recall: GD iteration  $w_{k+1} = w_k - \alpha_k \nabla f(w_k) \quad \alpha_k > 0$

$$= w_k - \frac{\alpha_k}{m} \sum_{i=1}^m \nabla f_i(w_k)$$

Stochastic Gradient (SG) iteration

$$w_{k+1} = w_k - \alpha_k \nabla f_{i_k}(w_k)$$

where  $i_k$  is drawn randomly in  $\{1, 2, \dots, m\}$   
 $\alpha_k > 0$

- Motivation:
- One iteration of SG only accesses 1 data point
  - ⇒ clearly cheaper than 1 iteration of GD ( $n$  times cheaper)
  - If the dataset is very large and there is an underlying data distribution, maybe seeing all data points at every iteration is not necessary for optimization

Does this work?

- SG does not always converge because it is a randomized algorithm (depends on a random draw) →  $i_0, i_1, \dots$
- But possible to show that SG works on average
- And in practice, the performance is really close to the average (even better)

Randomized algorithm: Running the method twice (on the same problem) gives different results

GD: deterministic, always converges (see previous lecture), expensive for large  $n$

SG: randomized, can fail to converge, cheap

## ② Basic SG theory

- Assumptions:
- $f \in C_{\mathbb{L}}^{1,1}$ , convex, every  $f_i$   $C^1$  convex
  - The random indices  $\{i_k\}_{k \in \mathbb{N}}$  are drawn at random so that:
    - (a)  $i_k$  is independent of  $i_0, \dots, i_{k-1}$
    - (b)  $\mathbb{E}_{i_k} [\nabla f_{i_k}(w_k)] = \nabla f(w_k)$
    - (c)  $\mathbb{E}_{i_k} [\|\nabla f_{i_k}(w_k)\|^2] \leq \|\nabla f(w_k)\|^2 + \sigma^2$   $\sigma \geq 0$

## Remarks

- (a) guarantees that every data point gets accessed infinitely often
- (b) on average, the stochastic gradient is the true gradient

$\mathbb{E}_{i_k}[\cdot]$ : expected value/mean of a quantity with respect to the random index  $i_k$

$$\mathbb{E}_{i_k}[\nabla f_{i_k}(w_k)] = \sum_{i=1}^n \underset{\substack{\uparrow \\ \text{probability} \\ \text{that } i_k \\ \text{takes the value } i}}{P(i_k=i)} \times \nabla f_i(w_k)$$

- (c) The variance of  $\|\nabla f_{i_k}(w_k)\|$  is bounded  
"the stochastic gradient  $\nabla f_{i_k}(w_k)$  doesn't vary too much from the true gradient  $\nabla f(w_k)$ "

Example: If at every iteration  $k$ ,  $i_k$  is drawn uniformly at random, i.e.  $P(i_k=i) = \frac{1}{n} \quad \forall i=1, \dots, n$

then (a), (b) and (c) are satisfied.

## L > Analysis:

Recall that the theory of GD is based on the fundamental inequality

$$f(w_{k+1}) \leq f(w_k) + \nabla f(w_k)^T (w_{k+1} - w_k) + \frac{L}{2} \|w_{k+1} - w_k\|^2$$

- For SG, the value of  $w_{k+1}$  is random (depends on the value of the random index  $i_k$ )  $\Rightarrow$  To remove the randomness, we take the

expectation  $E_{i_k}[\cdot]$  on both sides of the inequality

Theorem: Consider the  $k$ th iteration of SG  $w_{k+1} = w_k - \alpha_k \nabla f_{i_k}(w_k)$ .  
Then

$$E_{i_k} [f(w_{k+1})] \leq \underbrace{f(w_k)}_{\substack{\uparrow \\ \text{independent} \\ \text{of } i_k \text{ by (a)}}} - \alpha_k \underbrace{\nabla f(w_k)^T}_{\substack{\uparrow \\ = \nabla f(w_k) \\ \text{by (b)}}} E_{i_k} [\nabla f_{i_k}(w_k)] + \frac{L}{2} \alpha_k^2 E_{i_k} [\|\nabla f_{i_k}(w_k)\|^2] \leq \|\nabla f(w_k)\|^2 + \sigma^2 \text{ by (c)}$$

- ↳ With this inequality, we can find good choices for  $\alpha_k$  that guarantee decrease in  $f$  on average
- ↳ By combining these inequalities, we obtain convergence rates for SG

One example of CR rate (for convex  $f$ )

→ With appropriate choice of  $\alpha_k$ ,  $\forall k \geq 1$ ,

$$\underbrace{E_{i_0, \dots, i_{k-1}}}_{\substack{\text{Expected} \\ \text{value w.r.t. } i_0, \dots, i_{k-1}}} [f(w_k) - \min_{w \in \mathbb{R}^d} f(w)] \leq O\left(\frac{1}{\sqrt{k}}\right)$$

GD	$f(w_k) - \min_{w \in \mathbb{R}^d} f(w) \leq O\left(\frac{1}{k}\right)$	$mK$
SG	$E[f(w_k) - \min_{w \in \mathbb{R}^d} f(w)] \leq O\left(\frac{1}{\sqrt{k}}\right)$	$K$

guarantee after  $K$  iterations
number of data points accessed in  $K$  iterations

→ In terms of iterations, GD is best

→ But if we take the cost of every iteration into account, then SG is better!