

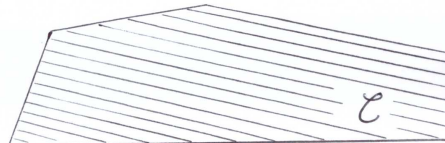
## SL = GRADIENT DESCENT ALGORITHMS

Gradient descent methods are a family of algorithms designed to optimize  $\min_{x \in \mathcal{C}} f(x)$ , where both  $f$  and  $\mathcal{C}$  are convex. We present in this chapter some popular gradient techniques, with perhaps the most important one: Stochastic Gradient Descent (SGD). We then adapt these algorithms to learning problems, where the objective is to minimize the empirical risk:  $\min_{f \in \mathcal{F}} \hat{R}_n(f)$  (or some regularized version of it). We first recall some important definitions and properties of functions, such as convexity, Lipschitzness, and subdifferentiability.

### I. PRELIMINARIES

#### I.1. Convexity:

Convexity of sets and functions were introduced in the chapter SL: SUPPORT VECTOR MACHINE. A set  $\mathcal{C}$  is convex if  $\forall x, y \in \mathcal{C}$ ,  $0 \leq \lambda \leq 1$ ,  $\lambda x + (1-\lambda)y \in \mathcal{C}$  [ $\mathcal{C}$  contains line segments between any two points]. A function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is convex if its domain is a convex set, and  $f(\lambda x + (1-\lambda)y) \leq \lambda f(x) + (1-\lambda)f(y)$ ,  $\forall x, y \in \text{dom } f$ ,  $\forall \lambda \in [0, 1]$ .



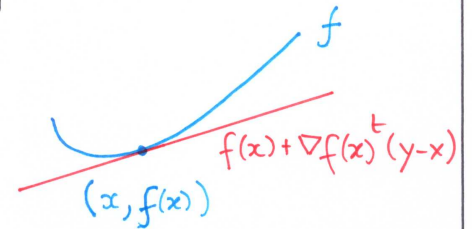
When  $f$  is differentiable / smooth, it is easy to check whether it is convex or not. (2)

→ First order condition ( $f$  differentiable)

" $f$  lies above its tangent line"

$$\forall x, y \in \text{dom } f,$$

$$f(y) \geq f(x) + \nabla f(x)^t (y-x).$$



→ Second order condition ( $f$  twice differentiable)

$$f \text{ is convex (strictly)} \Leftrightarrow \nabla^2 f(x) \succcurlyeq 0 \quad (\succ)$$

Ex: (i) Quadratic functions  $f(x) = \frac{1}{2} x^t P x + q^t x + r$

$$\nabla f(x) = P x + q$$

$$\nabla^2 f(x) = P \Rightarrow f \text{ is convex} \Leftrightarrow P \succcurlyeq 0$$

(ii) Least Squares objective  $f(\beta) = \|y - X\beta\|_2^2$

$$\nabla f(\beta) = -2 X^t (y - X\beta)$$

$$\nabla^2 f(\beta) = 2 X^t X \succcurlyeq 0 \quad (\text{convex } \forall X).$$

Theorem. Let  $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$  be convex functions  $i=1, \dots, n$ .

Then

(i)  $f(x) := \max_{1 \leq i \leq n} f_i(x)$  (pointwise maximum) is also

convex

(ii)  $f(x) := \sum_{i=1}^n \alpha_i f_i(x)$ ,  $\alpha_i \geq 0$  (weighted sum) is convex.

## I.2. Lipschitzness

(3)

Let  $\mathcal{C} \subset \mathbb{R}^d$  (not necessarily convex). A function  $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$  is L-Lipschitz on  $\mathcal{C}$  if  $\forall x, y \in \mathcal{C}$ ,

$$\|f(x) - f(y)\|_2 \leq L \|x - y\|_2$$

Euclidean norms in  $\mathbb{R}^k / \mathbb{R}^d$  respectively; but the notion of Lipschitzness can be generalized to arbitrary norms. We drop the subscript 2 in what follows.

• Consequences: An L-Lipschitz function cannot move too fast.

For example, if  $f: \mathbb{R} \rightarrow \mathbb{R}$  is differentiable  $\forall x \in \text{dom } f$ ,

$$\text{then } \forall h > 0, \quad \frac{|f(x+h) - f(x)|}{h} \leq L$$

$\downarrow h \rightarrow 0$   
 $f'(x)$

$\Rightarrow$  The derivative ( $\equiv$  rate of change of  $f$ ) is bounded by the Lipschitz constant (we may reach a similar conclusion using the mean value theorem:  $\exists u \in [x, y]$ , such that  $f(x) - f(y) = f'(u)(x - y)$ ).

Ex: (i)  $f(x) = |x|$  is 1-Lipschitz on  $\mathbb{R}$  since  $\forall x, y$ ,

$$|x| - |y| = |x - y + y| - |y| \leq |x - y| + |y| - |y| = |x - y|$$

$$\text{Likewise } |y| - |x| \leq |x - y| \Rightarrow ||x| - |y|| \leq |x - y|.$$

(ii)  $f(x) = \log(1 + e^x)$  is 1-Lipschitz on  $\mathbb{R}$  since

$$|f'(x)| = \left| \frac{e^x}{1 + e^x} \right| \leq 1$$

(iii)  $f(x)$  is not L-Lipschitz, for any L. Indeed, (4)

taking  $x=0$  and  $y=1+L$ , we have

$$f(y) - f(x) = (1+L)^2 > L(1+L) = L|y-x|$$

However,  $f(x)$  is L-Lipschitz on  $\mathcal{C} = \{x \mid |x| \leq \frac{L}{2}\}$  since  $\forall x, y \in \mathcal{C}$ ,

$$|f(x) - f(y)| = |x^2 - y^2| = |x - y| |x + y| \leq L |x - y| \leq 2x \frac{L}{2}$$

## I.3. Subgradients.

A convex and differentiable function  $f$  is such that the graph of  $f$  lies above its tangent line at any point of its domain. The slope of the tangent line being the gradient of  $f$  evaluated at that point. When dealing with non-differentiable functions, we introduce the notion of subdifferential  $\partial f(x)$  of  $f$  at  $x \in \text{dom } f$ :

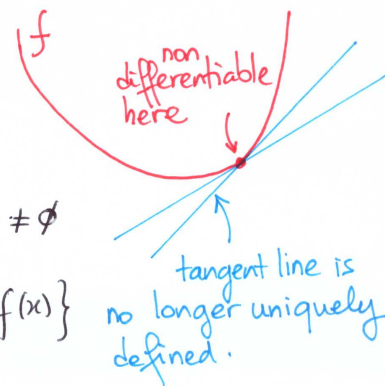
$$\partial f(x) = \{u \in \mathbb{R}^d \mid f(y) \geq f(x) + \langle u, y - x \rangle \quad \forall y \in \text{dom } f\}$$

$u \in \partial f(x)$  is called the SUBGRADIENT of  $f$  at  $x$ .

In fact, we have:

$$f \text{ is convex} \Leftrightarrow \forall x \in \text{dom } f \quad \partial f(x) \neq \emptyset$$

$$f \text{ is convex} \\ \& \text{ differentiable at } x \Leftrightarrow \partial f(x) = \{\nabla f(x)\}$$



Ex:  $\partial \|x\|_1 = \left\{ u \in \mathbb{R}^d \mid \begin{array}{ll} u_j = \text{sign } x_j & \text{if } x_j \neq 0 \\ u_j \in [-1, 1] & \text{if } x_j = 0 \end{array} \right\}$  (5)

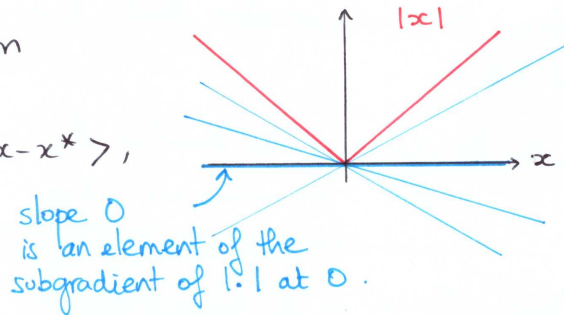
(see chapter SL: RIDGE REGRESSION & LASSO)

Theorem: The minimizers of a convex function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  are characterized by

$$x_* \in \underset{x \in \mathbb{R}^d}{\text{argmin}} f(x) \Leftrightarrow 0 \in \partial f(x_*)$$

Follows directly from  $\forall x \in \text{dom } f$

$f(x) \geq f(x_*) + \langle 0, x - x_* \rangle$ ,  
so that  
 $0 \in \partial f(x_*)$   
indeed ■



Theorem: Let  $\mathcal{C}$  be a convex open set  
 $f: \mathcal{C} \rightarrow \mathbb{R}$  a convex function.

Then

$$f \text{ is } L\text{-Lipschitz on } \mathcal{C} \Leftrightarrow \forall x \in \mathcal{C}, \forall u \in \partial f(x) \quad \|u\| \leq L$$

↳ For differentiable  $L$ -Lipschitz functions, the gradient is bounded by the Lipschitz constant. This theorem extends the result to non-differentiable  $L$ -Lipschitz functions, stating that the norm of the subgradient is necessarily less than  $L$ .

proof =  $\square$  Suppose that  $\forall x, \forall u \in \partial f(x), \|u\| \leq L$ .  
Since  $u \in \partial f(x)$ , we have  $f(y) \geq f(x) + \langle u, y - x \rangle$ .

Thus  $f(x) - f(y) \leq \langle u, x - y \rangle$  (6)  
 $\leq \|u\| \|x - y\|$  (Cauchy-Schwartz)  
 $\leq L \|x - y\|$

Likewise, we can show that  $f(y) - f(x) \leq L \|x - y\|$ , which concludes the first part of the theorem.

$\Rightarrow$  Suppose that  $f$  is  $L$ -Lipschitz.

Take  $x \in \mathcal{C}$  and  $u \in \partial f(x)$ .

Since  $\mathcal{C}$  is open,  $\exists \varepsilon > 0$  s.t.  $y := x + \varepsilon \frac{u}{\|u\|} \in \mathcal{C}$ .

Thus  $\rightarrow \langle y - x, u \rangle = \varepsilon \langle \frac{u}{\|u\|}, u \rangle = \varepsilon \|u\|$   
 $\rightarrow \|x - y\| = \varepsilon$ .

Since  $u \in \partial f(x)$ , we have  $f(y) \geq f(x) + \langle u, y - x \rangle$ ,  
or  $f(y) - f(x) \geq \langle u, y - x \rangle = \varepsilon \|u\|$ . (\*)

Moreover,  $f(y) - f(x) \leq L \|y - x\| = L \varepsilon$  (\*\*)  
 $\uparrow$   
 $f$  is  $L$ -Lipschitz

Putting (\*) & (\*\*) together, we obtain  $\|u\| \leq L$  indeed ■

## II. GRADIENT DESCENT ALGORITHMS.

### II.1. Standard Approach

We consider the unconstrained optimization problem

$$\underset{x \in \mathbb{R}^d}{\text{minimize}} f(x)$$

where  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is a convex function. If in addition  $f$  is differentiable, a necessary & sufficient condition for

(7)

a point  $x^*$  to be optimal is that the gradient of  $f$  at  $x^*$  vanishes:  $\nabla f(x^*) = 0$  ( $\equiv$  the advantage of working with convex functions: obtain global properties from local properties: a local minimum is a global minimum).

The optimization task is usually performed iteratively: an algorithm computes a sequence of points  $x_0, x_1, \dots$  such that (hopefully):  $f(x_k) \rightarrow f(x^*)$  as  $k \rightarrow +\infty$ . Descent algorithms produce a sequence  $\{x_k\}_{k \geq 0}$  s.t.

$$x_{k+1} = x_k + \eta_k \Delta x_k \quad \& \quad f(x_{k+1}) < f(x_k)$$

$\nwarrow$  step size ( $> 0$ )       $\nearrow$  search direction

For a convex & differentiable function  $f$ ,

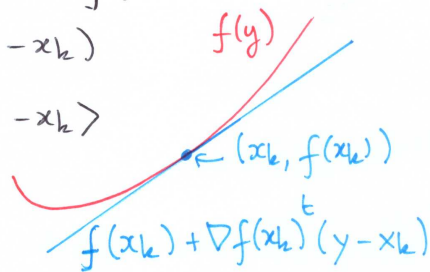
$$f(y) \geq f(x_k) + \nabla f(x_k)^t (y - x_k)$$

$$f(y) - f(x_k) \geq \langle \nabla f(x_k), y - x_k \rangle$$

To get  $f(y) < f(x_k)$ , the inner product must be negative, so that  $y - x_k$  must go in the direction of the negative gradient.

The search direction is then  $\Delta x_k = -\nabla f(x_k)$ , so that  $x_{k+1} = x_k - \eta_k \nabla f(x_k)$ .

If we loose differentiability, we loose the uniqueness of the tangent line, and  $y - x_k$  must be in the direction of one of the subgradients. We obtain the following procedure:



(8)

### Algorithm I : Gradient Descent Algorithm

1. Init:  $x_1 \in \mathbb{R}^d$  and a positive sequence  $\{\eta_j\}_{j \geq 1}$ .
2. For  $j = 1, \dots, k-1$  do
3.  $x_{j+1} = x_j - \eta_j g_j$ ,  $g_j \in \partial f(x_j)$
4. End For
5. Return either  $\bar{x} = \frac{1}{k} \sum_{j=1}^k x_j$  or  $x^0 \in \operatorname{argmin}_{x \in \{x_1, \dots, x_k\}} f(x)$

Goal: minimize  $f(x)$  ;  $f: \mathbb{R}^d \rightarrow \mathbb{R}$

Theorem: Let  $f =$  convex + Lipschitz function on  $\mathbb{R}^d$ , such that  $x^*$  exists.

Assume that  $\|x_1 - x^*\|_2 \leq R$ .

Then, if  $\eta_j = \eta = \frac{R}{L\sqrt{k}}$ ,  $\forall j \geq 1$ , we have

$$\rightarrow f\left(\frac{1}{k} \sum_{j=1}^k x_j\right) - f(x^*) \leq \frac{LR}{\sqrt{k}}$$

$$\rightarrow \min_{1 \leq j \leq k} f(x_j) - f(x^*) \leq \frac{LR}{\sqrt{k}}$$

Alternatively, we may consider  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ , and look for  $x \in \mathcal{C} \subset \mathbb{R}^d$  s.t.  $\mathcal{C} = \{x : \|x\| \leq R/2\}$  ( $\equiv \operatorname{diam}(\mathcal{C}) \leq R$ )

proof = First, note that

$$\min_{1 \leq j \leq k} f(x_j) - f(x^*) \leq \frac{1}{k} \sum_{j=1}^k f(x_j) - f(x^*)$$

$$f\left(\frac{1}{k} \sum_{j=1}^k x_j\right) - f(x^*)$$

Jensen.

⇒ We turn our attention to the term  $f(x_j) - f(x^*)$ . (9)

$$f(x_j) - f(x^*) \leq g_j^t(x_j - x^*) \quad (\text{def of subgradient})$$

$$= \frac{1}{2} (x_j - x_{j+1})^t (x_j - x^*)$$

Since  $\forall u, v \in \mathbb{R}^d$ ,  
 $2u^t v = \|u+v\|^2 - \|u-v\|^2$

(since  $2g_j = x_j - x_{j+1}$ )

$$= \frac{1}{2\eta} \left\{ \|x_j - x_{j+1}\|^2 + \|x_j - x^*\|^2 - \|x_{j+1} - x^*\|^2 \right\}$$

$$= \frac{\eta}{2} \|g_j\|^2 + \frac{1}{2\eta} \{d_j^2 - d_{j+1}^2\},$$

where  $d_j := \|x_j - x^*\|$ .

$f$  is  $L$ -Lipschitz  $\Rightarrow \|g_j\|^2 \leq L^2$  (Theorem p. 5),  
 so that

$$f(x_j) - f(x^*) \leq \frac{\eta L^2}{2} + \frac{1}{2\eta} (d_j^2 - d_{j+1}^2) \quad (1)$$

summing up from 1 to  $k$ : telescoping sum

$$\frac{1}{k} \sum_{j=1}^k f(x_j) - f(x^*) \leq \frac{\eta L^2}{2} + \frac{1}{2\eta k} (d_1^2 - d_{k+1}^2)$$

$$\leq \frac{\eta L^2}{2} + \frac{1}{2\eta k} d_1^2$$

$$\leq \frac{\eta L^2}{2} + \frac{R^2}{2\eta k}$$

Optimizing with respect to  $\eta$  ( $\eta = \frac{R}{L\sqrt{k}}$ ) gives the result. ■

Remark = The learning rate depends on the maximum number of iterations,  $k$ . We would prefer to obtain theoretical guarantees with a learning rate independent of  $k$ . To this end, note that (10)

$$\left( \sum_{j=1}^k \eta_j \right) \left( \min_{1 \leq j \leq k} f(x_j) - f(x^*) \right) \leq \sum_{j=1}^k \eta_j (f(x_j) - f(x^*))$$

Thus, making use of relation (1) on page 9, we see that

$$\eta_j (f(x_j) - f(x^*)) \leq \frac{\eta_j^2 L^2}{2} + \frac{1}{2} (d_j^2 - d_{j+1}^2)$$

summing up:

$$\sum_{j=1}^k \eta_j (f(x_j) - f(x^*)) \leq \frac{L^2}{2} \sum_{j=1}^k \eta_j^2 + \frac{1}{2} \sum_{j=1}^k (d_j^2 - d_{j+1}^2)$$

$$\leq \frac{L^2}{2} \sum_j \eta_j^2 + \frac{R^2}{2}$$

dividing by  $\sum \eta_j$

$$\min_{1 \leq j \leq k} f(x_j) - f(x^*) \leq \frac{L^2}{2} \frac{\sum \eta_j^2}{\sum \eta_j} + \frac{R^2}{2} \frac{1}{\sum \eta_j}$$

For this term to go to 0 with  $k$ , we need  $\sum \eta_j \rightarrow +\infty$  &  $\frac{\sum \eta_j^2}{\sum \eta_j} \rightarrow 0$

One possible candidate is to take  $\eta_j = \frac{c}{\sqrt{j}}$ , for some  $c > 0$ . Indeed,

$$\rightarrow \sum_{j=1}^k \eta_j^2 \leq c^2 \sum_{j=1}^k \frac{1}{j} \leq c^2 \left( 1 + \int_1^k \frac{dx}{x} \right) \leq c^2 C \log k$$

$$\rightarrow \sum_{j=1}^k \eta_j \geq \sum_{j=1}^k \eta_k = \sum_{j=1}^k \frac{c}{\sqrt{k}} = c\sqrt{k}. \quad (11)$$

We obtain:

$$\min_{1 \leq j \leq k} f(x_j) - f(x^*) \leq \frac{L^2}{2} \frac{c_1 \log k}{\sqrt{k}} + \frac{R^2}{2} \frac{1}{c\sqrt{k}}$$

Optimizing with respect to  $c$ :  $c = \frac{R}{L\sqrt{c_1 \log k}}$ ,

we get a rate  $O\left(LR\sqrt{\frac{\log k}{k}}\right)$ .  
 We get an extra  $\sqrt{\log k}$  factor.

However, noticing that

$$\begin{aligned} \min_{1 \leq j \leq k} f(x_j) - f(x^*) &\leq \min_{\frac{k}{2} \leq j \leq k} f(x_j) - f(x^*) \\ &\leq \frac{1}{\sum_{j=k/2}^k \eta_j} \left( \sum_{j=k/2}^k f(x_j) - f(x^*) \right) \\ &\leq c_2 \frac{LR}{\sqrt{k}}, \end{aligned}$$

since with  $\eta_j = \frac{c}{\sqrt{j}}$ , we still have  $\sum_{j=k/2}^k \eta_j \geq c\sqrt{k}$ ,

$$\text{but } \sum_{j=k/2}^k \eta_j^2 \leq c^2 \sum_{j=k/2}^k \frac{1}{j} \leq c' \int_{k/2}^k \frac{1}{x} dx \leq c'',$$

for  $c', c'' > 0$  constants. However, we need to ensure that  $\|x_{k/2} - x^*\| \leq R$  which is not always true (unless we are restricted to  $\mathcal{C} = \{x \mid \|x\| \leq R/2\}$ ).

## II.2. Projected Gradient Descent.

When dealing with constrained optimization problem of the form  $\min_{x \in \mathcal{C}} f(x)$ , where  $\begin{cases} f, \mathcal{C} = \text{convex} \\ f: \mathbb{R}^d \rightarrow \mathbb{R} \end{cases}$ , we need to make sure that at each iteration  $x_j \in \mathcal{C}$ . The standard form of gradient descent does not ensure this. The projected gradient descent algorithm will incorporate this condition by adding an extra step at each iteration, and projecting the current estimate onto the convex set  $\mathcal{C}$ . We have the following result.

Theorem. Let  $\mathcal{C}$  be a closed convex set  $\subset \mathbb{R}^d$ .

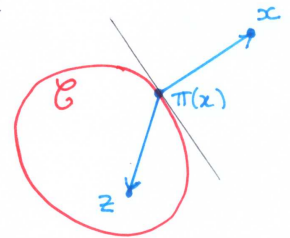
$\forall x \in \mathbb{R}^d$ , let  $\pi(x) \in \mathcal{C}$  be such that

$$\|x - \pi(x)\| = \min_{z \in \mathcal{C}} \|x - z\|.$$

Then  $\pi(x)$  is unique, and

$$\langle \pi(x) - x, \pi(x) - z \rangle \leq 0$$

$$\forall z \in \mathcal{C}$$



proof. Uniqueness.

$$\text{If } \pi_1, \pi_2 \in \mathcal{C} \text{ both satisfy } \begin{cases} \langle \pi_1 - x, \pi_1 - z \rangle \leq 0 \\ \langle \pi_2 - x, \pi_2 - z \rangle \leq 0 \end{cases} \quad \forall z,$$

Then taking  $z = \pi_2$  in the first equality  
 &  $z = \pi_1$  in the second,

$$\begin{cases} \langle \pi_1 - x, \pi_1 - \pi_2 \rangle \leq 0 \\ \langle x - \pi_2, \pi_1 - \pi_2 \rangle \leq 0 \end{cases} \rightarrow \text{adding them} \quad \|\pi_1 - \pi_2\|^2 \leq 0 \Rightarrow \pi_1 = \pi_2.$$

• Proof that  $\langle \pi - x, \pi - z \rangle \leq 0$ . (13)

By definition,  $\|x - \pi\|^2 \leq \|x - z\|^2 \quad \forall z \in \mathcal{C}$ .

Fix  $u \in \mathcal{C}$  and put  $z := (1 - \varepsilon)\pi + \varepsilon u$ ;  $\varepsilon \in (0, 1]$ .

$\uparrow$   
 $\in \mathcal{C}$  since  $\mathcal{C}$  is convex, and both  $u, \pi$  belong to  $\mathcal{C}$ .

$$\Rightarrow \|x - \pi\|^2 \leq \|x - z\|^2 = \|x - \pi - \varepsilon(u - \pi)\|^2$$

$$\Downarrow \quad \|x - \pi\|^2 - 2\varepsilon \langle x - \pi, u - \pi \rangle + \varepsilon^2 \|u - \pi\|^2$$

$$\langle x - \pi, u - \pi \rangle \leq \varepsilon \|u - \pi\|^2$$

Valid  $\forall \varepsilon \in (0, 1]$ . Taking the limit  $\varepsilon \rightarrow 0$  yields the desired result. ■

### Algorithm 2: Projected Gradient Descent Algorithm

1. Init:  $x_1 \in \mathcal{C}$  and a positive sequence  $\{\gamma_j\}_{j \geq 1}$ .
2. For  $j = 1, \dots, k-1$  do
3.  $z_{j+1} = x_j - \gamma_j g_j$ ,  $g_j \in \partial f(x_j)$
4.  $x_{j+1} = \pi(z_{j+1})$
5. End For
6. Return either  $\bar{x} = \frac{1}{k} \sum_{j=1}^k x_j$  or  $x^0 \in \operatorname{argmin}_{x \in \{x_1, \dots, x_k\}} f(x)$

↖ Note that the projection is in terms of the Euclidean norm  $\|\cdot\|_2$ . Goal: minimize  $f(x)$ ;  $f: \mathbb{R}^d \rightarrow \mathbb{R}$   
 $x \in \mathcal{C}$  ← close & convex

Theorem. Let  $\mathcal{C} =$  closed, convex, non-empty subset of  $\mathbb{R}^d$ , such that  $\operatorname{diam}(\mathcal{C}) \leq R$ . (14)

Let  $f$  be a convex,  $L$ -Lipschitz function on  $\mathcal{C}$ , such that  $x^* \in \operatorname{argmin}_{x \in \mathcal{C}} f(x)$  exists.   
 ↗ but  $f$  is defined  $\forall x \in \mathbb{R}^d$

• Then; with  $\gamma_j = \frac{R}{L\sqrt{k}}$ , we have

$$f(\bar{x}) - f(x^*) \leq \frac{LR}{\sqrt{k}} \quad \& \quad f(x^0) - f(x^*) \leq \frac{LR}{\sqrt{k}}$$

• In addition, if  $\gamma_j = \frac{R}{L\sqrt{j}}$ , there exist  $c > 0$  s.t.

$$f(\bar{x}) - f(x^*) \leq \frac{cLR}{\sqrt{k}} \quad \& \quad f(x^0) - f(x^*) \leq \frac{cLR}{\sqrt{k}}$$

proof: We proceed as before; and evaluate the difference

$$f(x_j) - f(x^*) \leq g_j^t (x_j - x^*)$$

$$= \frac{1}{2} (x_j - z_{j+1})^t (x_j - x^*)$$

(with  $\eta = \gamma_j$ )

$$= \frac{1}{2\eta} \left\{ \|x_j - z_{j+1}\|^2 + \|x_j - x^*\|^2 - \|z_{j+1} - x^*\|^2 \right\}.$$

We turn our attention to the term  $\|z_{j+1} - x^*\|^2$ :

$$\|z_{j+1} - x^*\|^2 = \|z_{j+1} - x_{j+1}\|^2 + \|x_{j+1} - x^*\|^2 + 2 \langle z_{j+1} - x_{j+1}, x_{j+1} - x^* \rangle$$

$$\geq \|x_{j+1} - x^*\|^2$$

↖  $\pm x_{j+1}$   
 & expanding

$$\underbrace{\|z_{j+1} - x_{j+1}\|^2 + \|x_{j+1} - x^*\|^2 + 2 \langle z_{j+1} - x_{j+1}, x_{j+1} - x^* \rangle}_{\geq 0}$$

Moreover,  $\|x_j - z_{j+1}\|^2 = \eta^2 \|g_j\|^2 \leq \eta^2 L^2$ . We obtain:

$$\frac{1}{k} \sum_{j=1}^k f(x_j) - f(x^*) \leq \frac{1}{k} \sum_{j=1}^k \frac{1}{2\eta} (\eta^2 L^2 + \|x_j - x^*\|^2 - \|x_{j+1} - x^*\|^2) \quad (15)$$

$$\leq \frac{\eta L^2}{2} + \frac{1}{2\eta k} \underbrace{\|x_1 - x^*\|^2}_{\leq R^2}$$

Minimizing the RHS yields  $\eta = \frac{R}{L\sqrt{k}}$ , and the desired upper bound.

The case  $\eta_j = \frac{R}{L\sqrt{j}}$  is treated as in the proof of the performance of the standard gradient descent, since we have  $\|x_{k/2} - x^*\|^2 \leq R^2$ .

### I.3. Stochastic Gradient Descent.

We turn our attention to the optimization of random functions  $x \mapsto l(x, z)$ ,

where  $x$  = optimization parameter  
 $z$  = RV with distribution  $\mathbb{P}_z$ .

We assume that  $x \mapsto l(x, z)$  is convex,  $\mathbb{P}_z$ -a.s.

(so that in a practical setting, for each observation  $z_i$  we make,  $x \mapsto l(x, z_i)$  is convex).

In particular,  $x \mapsto \mathbb{E}\{l(x, z)\}$  is also convex in  $x$ .

Our goal is to solve the optimization problem

$$\min_{x \in \mathcal{C}} \mathbb{E}\{l(x, z)\} \quad \leftarrow \text{put } f(x) := \mathbb{E}\{l(x, z)\}$$

↑ convex set & close

Remark: In the learning context,  $z$  correspond to pairs  $(X, Y)$ ,  $l$  = loss function, and  $\mathbb{E}\{l(x, z)\}$  represents the risk of a predictor parametrized by the variable  $x$ .

Ex: square loss + linear classifier:

$$\mathbb{E}\{l(x, z)\} = \mathbb{E}(Y - \beta_0 - \beta^t x)^2$$

$$= R(f); \text{ for } f \in \mathcal{F} = \{x \mapsto \beta_0 + \beta^t x\}$$

→ Usually, the distribution  $\mathbb{P}_z$  is unknown, and we cannot compute explicitly the objective function  $\mathbb{E}\{l(x, z)\}$ . Instead, we are given the observations  $z_1, z_2, \dots$  and the associated values of  $l(x, z_i)$ . This means that a gradient descent method cannot be directly applied here, as the direction of the steepest descent cannot be calculated. Stochastic Gradient Descent circumvents this problem by taking a step at each iteration in a random direction, as long as the expected value of this direction corresponds to the negative gradient.

→ Consider two distinct scenarios =

(i) At each iteration of the algorithm, consider a new pair  $z_i = (X_i, Y_i)$ . Each pair is used at most once. We do ONE PASS on the data.

(ii) External Randomization: the SGD algorithm can be applied in the context of ERM, where the objective function  $f(x)$  is of the form  $\frac{1}{n} \sum_{i=1}^n f_i(x)$ .



Once observations  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  are collected, the functions  $f_i(x) = l(y_i, g(x_i))$ ,  $g \in \mathcal{F}$ , are deterministic. We artificially introduce randomness as follows: let  $I$  be a RV uniformly distributed on  $\{1, \dots, n\}$ . We then have the

representation  $f(x) = \mathbb{E}_I \{f_I(x)\}$ , which falls into the context of stochastic (convex) optimization, with  $Z=I$ , and  $l(x, I) = f_I(x)$ .

↓ In the notation of ERM,  $g_n \in \operatorname{argmin}_{g \in \mathcal{F}} \hat{R}_n(g)$ ,  
 where  $\hat{R}_n(g) = \frac{1}{n} \sum_{i=1}^n l(y_i, g(x_i))$   
 $= \mathbb{E}_I \{l(y_I, g(x_I))\}$

When performing  $k$  passes of the algorithm, we may take  $k$  as large as possible, regardless of the value of  $n$ , by doing MULTIPLE PASSES on the data (the performance of the minimizer is however limited by the number of observations).

Summary:

- In (i), observations are received on-line, each  $(X_i, Y_i) \sim P_{X,Y}$ , and used only one.  
Objective: minimize the true risk  $\mathbb{E}\{l(Y, f(X))\}$
- In (ii), we have received  $n$  observations. We pick randomly  $k$  observations from  $\{(X_1, Y_1), \dots, (X_n, Y_n)\}$ , with replacement.  
Objective: minimize the empirical risk  $\frac{1}{n} \sum_{i=1}^n l(Y_i, f(X_i))$ .

- In the SGD algorithm, each time an observation  $Z = z_j$  is made, a subgradient  $\tilde{g}_j$  of  $l(x, z_j)$  at the current estimate  $x_j$  is taken. For any  $x$ , we have   
 ↗ And used to update the parameters

$$l(x, z_j) - l(x_j, z_j) \geq \langle \tilde{g}_j, x - x_j \rangle$$

↓ Taking  $\mathbb{E}(\cdot)$  with respect to the distribution of  $Z$ , conditionally on  $x_j$ , we obtain

$$\mathbb{E}\{l(x, Z)\} - \mathbb{E}\{l(x_j, Z) | x_j\}$$

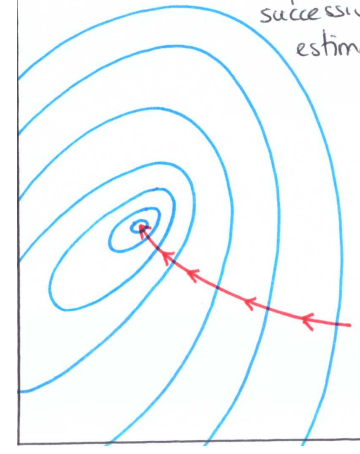
Always hold in scenario (ii) →  $\geq \langle \mathbb{E}(\tilde{g}_j | x_j), x - x_j \rangle$ ,

so that  $\mathbb{E}(\tilde{g}_j | x_j) =: g_j$  is a subgradient of  $\mathbb{E}\{l(x_j, Z) | x_j\}$

↳ "on average",  $\tilde{g}_j$  points towards the right direction

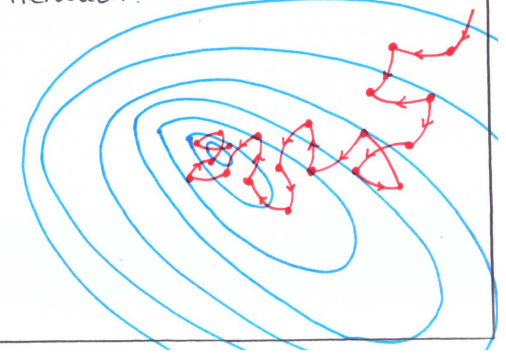
↳ Gradient Descent

Updates point toward the direction of the largest decrease ⇒ smooth paths of successive estimates



↳ Stochastic Gradient Descent

Updates point "on average" towards the direction of the largest decrease. ⇒ Rough path: "Not a descent" algorithm as such; no guarantee to decrease the objective at each iteration.



### Algorithm 3: Stochastic Gradient Descent (SGD)

19

1. Init:  $x_1 \in \mathcal{C}$
2. positive sequence  $\{\eta_j\}_{j \geq 1}$
3. independent RVs  $z_1, \dots, z_k$  with distribution  $\mathbb{P}_z$
4. For  $j=1, \dots, k-1$  do
  5.  $y_{j+1} = x_j - \eta_j \tilde{g}_j$ ,  $\tilde{g}_j \in \partial \ell(x_j, z_j)$  ← a random sequence
  6.  $x_{j+1} = \pi(y_{j+1})$
7. End For
8. Return  $\bar{x}_k = \frac{1}{k} \sum_{j=1}^k x_j$

Goal:  $\min_{x \in \mathcal{C}} \mathbb{E}\{\ell(x, z)\}$ ;  
 $\ell(\cdot, z): \mathbb{R}^d \rightarrow \mathbb{R}$

Typically, the objective function  $\mathbb{E}\{\ell(x, z)\}$  is unknown, and  $x^*$  cannot be computed.

Theorem: Let  $\mathcal{C} =$  closed convex subset of  $\mathbb{R}^d$ ,  $\text{diam } \mathcal{C} \leq R$ .

- $f(x) = \mathbb{E}\{\ell(x, z)\}$  attains its minimum on  $\mathcal{C}$  at  $x^*$ ;  $f: \mathbb{R}^d \rightarrow \mathbb{R}$
- $\ell(x, z)$  is convex on  $x$ ,  $\mathbb{P}_z$ -a.s.
- $\mathbb{E}\|\tilde{g}\|^2 \leq L^2 \quad \forall \tilde{g} \in \partial \ell(x, z), \forall x$ .

Then, with  $\eta_j = \eta = \frac{R}{L\sqrt{k}}$ ,  $\mathbb{E}\{f(\bar{x}_k)\} - f(x^*) \leq \frac{LR}{\sqrt{k}}$

proof = We have  $f(x_j) - f(x^*) \leq g_j^t(x_j - x^*)$

by definition of subgradient of  $f$

$$= \mathbb{E}\{g_j^t(x_j - x^*) \mid x_j\}$$

since  $\mathbb{E}\{\tilde{g}_j \mid x_j\} = g_j$

$$= \frac{1}{\eta} \mathbb{E}\{(y_{j+1} - x_j)^t (x_j - x^*) \mid x_j\}$$

$\eta_j = \eta$

$$\Rightarrow f(x_j) - f(x^*) = \frac{1}{\eta} \left\{ \mathbb{E}[\|x_j - y_{j+1}\|^2 + \|x_j - x^*\|^2 - \|y_{j+1} - x^*\|^2 \mid x_j] \right\} \quad (20)$$

$$\leq \frac{1}{\eta} \left\{ \mathbb{E}(\eta^2 \|\tilde{g}_j\|^2 \mid x_j) + \mathbb{E}(\|x_j - x^*\|^2 \mid x_j) - \underbrace{\mathbb{E}(\|x_{j+1} - x^*\|^2 \mid x_j)}_{\text{argue as on page 14 for this term}} \right\}$$

Taking  $\mathbb{E}(\dots)$ , summing over  $j$ , using Jensen's inequality, and taking  $\eta = \frac{R}{L\sqrt{k}}$  yields the upper bound. ▀

\* Remark: On-line algorithms that compute the minimum (or root) of a function  $x \mapsto \mathbb{E}\{\ell(x, z)\}$  date back to Robbins & Monro in the 50's. Such algorithms are known as Stochastic Approximation Algorithms in the statistics literature. Under mild assumptions on  $\ell$  and on the learning rate  $\eta_j$  ( $\eta_j \rightarrow 0$ ,  $\sum \eta_j = +\infty$ ,  $\sum \eta_j^2 < \infty$ ), almost sure convergence and convergence in mean square towards the minimum (or root) can be established, using tools from martingale theory. The two most famous procedures are due to Robbins & Monro (1951) and Kiefer & Wolfowitz (1952). Some basic convergence results of the Robbins & Monro algorithm are presented in the Appendix of this chapter.

### III - APPLICATION TO LEARNING PROBLEMS.

(21)

In the context of empirical risk minimization, the goal is to minimize the empirical risk  $\hat{R}_n(f_\beta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_\beta(x_i))$  with respect to a vector of parameters  $\beta$ . The key difference between GD & SGD algorithms applied in this context is that at each iteration, either we update according to:

$$\beta_{j+1} \leftarrow \beta_j - \eta_j \left\{ \frac{1}{n} \sum_{i=1}^n \nabla_{\beta} \ell(y_i, f_{\beta}(x_i)) \right\} \quad (\text{GD})$$

or  
 gradient, or subgradient  
 slow: one update requires a pass through  $n$  observations.  
 (smooth trajectories of the estimates  $\{\beta_j\}$ )

$$\beta_{j+1} \leftarrow \beta_j - \eta_j \left\{ \nabla_{\beta} \ell(y_i, f_{\beta}(x_i)) \right\} \quad (\text{SGD})$$

much faster: one update requires only one observation.  
 ( $\Rightarrow$  rough trajectories)  
 (May not settle at the minimum, even if we reach it: wiggles around).  
 On the other hand, the noisy updates can allow the algorithm to avoid local minima.

Theoretical guarantees for particular choices of  $\{\eta_j\}$

MINI-BATCH SGD is a compromise between GD (all training data used at once) and the SGD (only one observation is used at a time) that splits the learning sample into small batches, used to update the parameters at each iteration:

$$\beta_{j+1} \leftarrow \beta_j - \eta_j \left\{ \frac{1}{B} \sum_{i \in \mathcal{B}} \nabla_{\beta} \ell(y_i, f_{\beta}(x_i)) \right\}$$

sum over all observations belonging to the mini-batch, of size  $B$ .

Updates are more frequent than for GD; local minima may be avoided; smoothes up the trajectories of the SGD updates.

(popular in SGD in deep neural nets)

• Example 1 = We implement soft-SVM using SGD. Recall that the optimization problem is

(22)

$$\min_{\beta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (1 - y_i \langle \beta, x_i \rangle)_+ + \frac{\lambda}{2} \|\beta\|^2$$

Consider for simplicity the homogeneous case; by setting the bias term  $\beta_0$  to 0. Note that a bias term can be introduced by adding one more feature equal to 1 into the input vectors  $x_i$  (in which case we regularize the bias term as well).

hinge loss  
 $\varphi(z) = \max(0, 1+z)$

The hinge loss is not differentiable at  $z = -1$ , and we need to consider subgradients. Note that the function  $\mathcal{L}(\beta) := (1 - y \langle \beta, x \rangle)_+$  is  $\|x\|$ -Lipschitz, and that its subgradients are all of the form  $\alpha x$ , with  $|\alpha| \leq 1$ . Specifically, the vector

$$v = \begin{cases} 0 & \text{if } y \langle \beta, x \rangle \geq 1 \\ -yx & \text{if } y \langle \beta, x \rangle < 1 \end{cases} \text{ belongs to } \partial \mathcal{L}(\beta).$$

More generally, if  $h(\beta) = \max_{1 \leq i \leq m} h_i(\beta)$  where  $h_i = \text{convex} + \text{different}$ , then  $\nabla h_j(\beta) \in \partial h(\beta)$ , where  $j = \text{argmax}_{1 \leq i \leq m} h_i(\beta)$ .

Indeed,  $h_j$  convex  $\Rightarrow \forall u \quad h_j(u) \geq \underbrace{h_j(\beta)}_{=h(\beta)} + \langle \nabla h_j(\beta), u - \beta \rangle$

$$\Rightarrow h(u) \geq h_j(u) \geq h(\beta) + \langle \nabla h_j(\beta), u - \beta \rangle \Rightarrow \nabla h_j(\beta) \in \partial h(\beta).$$

SGD: Current estimate:  $\beta_j$ .

. Select at random an observation  $(x_i, y_i)$  from  $\mathcal{Z}_n$ .

Then  $(v_j + \lambda \beta_j)$ ; where  $v_j = \begin{cases} -y_i x_i & \text{if } y_i \langle \beta_j, x_i \rangle < 1 \\ 0 & \text{otherwise} \end{cases}$

is a subgradient of  $(1 - y_i \langle \beta_j, x_i \rangle)_+ + \frac{\lambda}{2} \|\beta_j\|^2$  at  $\beta_j$ .

$$\Rightarrow \beta_{j+1} \leftarrow \beta_j - \eta_j (v_j + \lambda \beta_j) \text{ (with no projection step).}$$

In the case of soft-SVM, and more generally, when a quadratic regularizer is added to the convex empirical risk, we can achieve faster rates of convergence than the ones stated on page 19, since the objective function belongs to the family of  $\lambda$ -STRONGLY CONVEX functions.

A function  $f$  is said to be  $\lambda$ -strongly convex if its curvature remains bounded away from 0. Specifically, if  $f(x) - \frac{\lambda}{2} \|x\|^2$  is convex  $\forall x$ . There are equivalent definitions of  $\lambda$ -strongly convex functions. Another characterization is

$$\forall x, y \quad \forall u \in \partial f(x), \quad \langle u, x-y \rangle \geq f(x) - f(y) + \frac{\lambda}{2} \|x-y\|^2.$$

If  $v \in \partial(f(x) - \frac{\lambda}{2} \|x\|^2)$ , then  $v = u + \lambda x$  for some  $u \in \partial f(x)$ . Then

$$f(y) - \frac{\lambda}{2} \|y\|^2 \geq f(x) - \frac{\lambda}{2} \|x\|^2 + \langle u + \lambda x, y-x \rangle$$

$$\Leftrightarrow x \mapsto f(x) - \frac{\lambda}{2} \|x\|^2 \text{ is convex}$$

$$\langle u, x-y \rangle \geq f(x) - f(y) + \frac{\lambda}{2} (\|y\|^2 - \|x\|^2 - 2\langle x, y-x \rangle) = \|x-y\|^2$$

**Algorithm 4: SGD for  $\lambda$ -strongly convex functions.**

1. Init:  $x_1 \in \mathcal{C}$  / indpt RVs  $z_1, \dots, z_k \sim \mathbb{P}_z$ .
2. For  $j=1, \dots, k-1$  do
3.  $y_{j+1} = x_j - \eta_j \tilde{g}_j$ ,  $\eta_j = \frac{1}{\lambda_j}$ ,  $\tilde{g}_j \in \partial \ell(x_j, z_j)$
4.  $x_{j+1} = \pi(y_{j+1})$
5. End For
6. Return  $\bar{x}_k = \frac{1}{k} \sum_{j=1}^k x_j$

**Goal:**  $\min_{x \in \mathcal{C}} \mathbb{E} \{ \ell(x, z) \}$ , where  $\mathbb{E} \{ \ell(x, z) \}$  is  $\lambda$ -str. conv. in  $x$ .

Theorem: Let

- $\mathcal{C}$  = closed convex subset of  $\mathbb{R}^d$
- $f(x) = \mathbb{E} \{ \ell(x, z) \}$  attains its minimum on  $\mathcal{C}$  at  $x^*$
- $f$  is  $\lambda$ -strongly convex in  $x$ .
- $\mathbb{E} \|\tilde{g}\|^2 \leq L^2 \quad \forall \tilde{g} \in \partial \ell(x, z), \forall x$ .

Then, with  $\eta_j := \frac{1}{\lambda_j}$ ,  $\mathbb{E} \{ f(\bar{x}_k) \} - f(x^*) \leq \frac{L^2}{2\lambda k} (1 + \log k)$

The diameter of  $\mathcal{C}$  does not appear in the bound.

The "log k" can be removed, see Rakhlin, Shamir, Sridharan (2012).  $\Rightarrow$  Faster convergence in  $\frac{1}{k}$  instead of  $\frac{1}{\sqrt{k}}$ .

proof = Put  $g_j := \mathbb{E}(\tilde{g}_j | x_j) \in \partial f(x_j)$ . Since  $f$  is  $\lambda$ -strongly convex, we can write

(0)  $\langle g_j, x_j - x^* \rangle \geq f(x_j) - f(x^*) + \frac{\lambda}{2} \|x_j - x^*\|^2$

On the other hand, we have

(1)  $\langle g_j, x_j - x^* \rangle \leq \frac{1}{2\eta_j} \{ \mathbb{E}(\|x_j - x^*\|^2 - \|x_{j+1} - x^*\|^2) \} + \frac{\eta_j}{2} L^2$

Indeed, since  $x_{j+1} = \pi(y_{j+1})$ , we have that

$$\|y_{j+1} - x^*\|^2 \geq \|x_{j+1} - x^*\|^2.$$

$$\Rightarrow \|x_j - x^*\|^2 - \|x_{j+1} - x^*\|^2 \geq \|x_j - x^*\|^2 - \|y_{j+1} - x^*\|^2$$

$$\underbrace{\|x_j - \eta_j \tilde{g}_j - x^*\|^2}_{\|x_j - x^*\|^2 - 2\langle x_j - x^*, \eta_j \tilde{g}_j \rangle + \eta_j^2 \|\tilde{g}_j\|^2} \geq \|x_j - x^*\|^2 - \|y_{j+1} - x^*\|^2$$

Rearranging the terms, taking expectations, and using  $\mathbb{E} \|\tilde{g}\|^2 \leq L^2$  gives (1).

Putting (0) & (1) together, and summing over  $j$  yields:

$$\sum_{j=1}^k \mathbb{E} \{ f(x_j) \} - f(x^*) \leq \mathbb{E} \left\{ \sum_{j=1}^k \frac{\|x_j - x^*\|^2 - \|x_{j+1} - x^*\|^2}{2\eta_j} - \frac{\lambda}{2} \|x_j - x^*\|^2 \right\} + \frac{L^2}{2} \sum_{j=1}^k \eta_j \leq \frac{1}{\lambda} (1 + \log k)$$

With  $\eta_j = \frac{1}{\lambda_j}$ , we get

$$\frac{\lambda}{2} \sum_{j=1}^k \{ (j-1) \|x_j - x^*\|^2 - j \|x_{j+1} - x^*\|^2 \} = -\frac{\lambda}{2} k \|x_{k+1} - x^*\|^2 \leq 0$$

$$\Rightarrow \sum_{j=1}^k \mathbb{E} \{ f(x_j) \} - f(x^*) \leq \frac{L^2}{2\lambda} (1 + \log k).$$

→ Back to soft-SVM. Recall that we wish to optimize

$$\frac{1}{n} \sum_{i=1}^n (1 - y_i \langle \beta, x_i \rangle)_+ + \frac{\lambda}{2} \|\beta\|^2.$$

Given the current estimate  $\beta_j$ , and with  $\mathcal{C} = \mathbb{R}^d$ , we obtain

$$\beta_{j+1} \leftarrow \beta_j - \frac{1}{\lambda_j} (v_j + \lambda \beta_j), \text{ where}$$

$$v_j = \begin{cases} 0 & \text{if } y_i \langle \beta_j, x_i \rangle \geq 1 \\ -y_i x_i & \text{if } \text{---} < 1 \end{cases}$$

associated with the randomly chosen observation  $(x_i, y_i)$ .

The RHS is equal to

$$\left(1 - \frac{1}{j}\right) \beta_j - \frac{1}{\lambda_j} v_j = \frac{j-1}{j} \left\{ \frac{j-2}{j-1} \beta_{j-1} - \frac{1}{\lambda_{j-1}} v_{j-1} \right\} - \frac{1}{\lambda_j} v_j$$

$$\frac{j-1}{j} = \dots = -\frac{1}{\lambda_j} \sum_{i=1}^j v_i.$$

Starting with  $\beta_1 \equiv 0$ .

We obtain the following simple algorithm:

Goal:  $\min_{\beta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (1 - y_i \langle \beta, x_i \rangle)_+ + \frac{\lambda}{2} \|\beta\|^2.$

Init:  $S_1 = 0$

For  $j=1, \dots, k-1$  do

•  $\beta_j = \frac{1}{\lambda_j} S_j$

• Select  $(x_i, y_i)$  at random from  $\mathcal{D}_n$ .

• If  $y_i \langle \beta_j, x_i \rangle < 1$

Put  $S_{j+1} \leftarrow S_j + y_i x_i$

Else

Put  $S_{j+1} \leftarrow S_j$

End For

Return  $\bar{\beta} = \frac{1}{k} \sum_{j=1}^k \beta_j.$

"Pegasos" Algorithm

See Shalev-Shwartz, Singer, Srebro, Cote (2011)

SGD for Soft-SVM.

• Example 2 = K-means.

The K-means objective is the minimization of the cost

$$C_0(M, \pi) = \frac{1}{n} \sum_{k=1}^K \sum_{i=1}^n \pi_{ik} \|x_i - m_k\|_2^2,$$

where  $\pi_{ik} = \begin{cases} 1 & \text{if } x_i \in k\text{-th cluster} \\ 0 & \text{otherwise} \end{cases}$

$M = \{m_1, \dots, m_k\}$

= set of K representatives, summarizing the data  $\mathcal{D}_n = \{x_1, \dots, x_n\}$ ,  $x_i \in \mathbb{R}^d$ .

For more details, we refer to the chapter UL: CLUSTERING, page 4. We make the K-means algorithm presented p.7 there

online, by considering a SGD:

(27)

At each iteration, pick randomly an observation  $x_i$  and update the nearest centroid using the gradient of  $\sum_{i=1}^K \pi_{ik} \|x_i - m_k^{(j)}\|_2^2$  with respect to the current estimate  $m_k^{(j)}$ :

$$\nabla_{m_k^{(j)}} (\dots) = -2\pi_{ik}(x_i - m_k^{(j)})$$

↑ All zero, except for the closest centroid of  $x_i$ .

$$\text{Update: } m_k^{(j+1)} \leftarrow m_k^{(j)} + \eta_j (x_i - m_k^{(j)})$$

↑ learning rate SGD for K-MEANS.

See MacQueen (1967).

- SGD is a general algorithm that may be applied to a wide variety of optimization problems. For instance, a stochastic version of the coordinate descent algorithm was successfully applied to the lasso criterion (see Shalev-Shwartz & Tewari (2011)). In the next chapter, we discuss two important classes of model where SGD are particularly useful: the perceptron & neural networks.

## References:

- S. Shalev-Shwartz & S. Ben-David (2014). Understanding Machine Learning: From Theory to Algorithms. Cambridge.
- S. Shalev-Shwartz, Y. Singer, N. Srebro, A. Cotter (2011). Pegasos: Primal Estimated Sub-Gradient Solver for SVM. Mathematical Programming, vol 127, no 1, p. 3-30.
- S. Shalev-Shwartz & A. Tewari (2011). Stochastic Methods For  $l_1$ -Regularized Loss Minimization. Journal of Machine Learning Research. Vol 12, p. 1865-1892.
- P. Rigollet. Lecture notes on "Mathematics of Machine Learning".

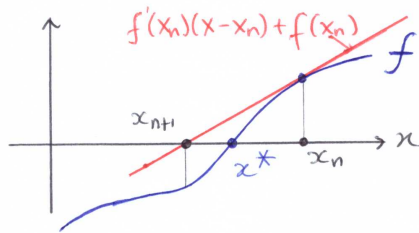
APPENDIX = STOCHASTIC APPROXIMATION ALGORITHMS (a)

I - THE ROBBINS - MONRO ALGORITHM

I.1. Introductory remarks.

• Goal: Find the root of a real-valued function  $f(\cdot)$ .  
 If the function  $f$  were known,  $f$  assumed differentiable,  
 $f: \mathbb{R} \rightarrow \mathbb{R}$ , Newton's procedure could be used: the  
 sequence  $x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$  converges to the root  $x^*$   
 of  $f$ .  $f' =$  derivative of  $f$ .

↳ Idea: approximate  $f$  by its tangent line: given a current  
 approximation  $x_n$ , the tangent line at  $x_n$  is  
 $f'(x_n)(x - x_n) + f(x_n)$ .



the  $x$ -intercept of this  
 line is used as the  
 next approximation:

$$f'(x_n)(x_{n+1} - x_n) + f(x_n) = 0$$

If  $f$  is not differentiable, a less efficient algorithm  
 would be  $x_{n+1} = x_n - \gamma f(x_n)$  (\*)  
 $\gamma$  step size;  $\gamma > 0$

Will converge provided the initial estimate  $x_0$  is not too  
 far away from the root.

→ Assume now that  $f$  is not directly observable, but that only noisy  
 measurements of  $f(x)$  are available  $\forall x$ .

We observe:  $f(x) + U =$  noisy measurement (b)

noise term, with zero mean  $\mathbb{E}U = 0$

One can use the procedure (\*), with  $f(x_n)$  replaced by  
 a good estimate of its value, obtained by averaging many  
 observations:  $f(x_n) \approx \frac{1}{m} \sum_{i=1}^m (f(x_n) + U_i)$

call this term  $F(x_i, U_i)$ ,  
 so that  $\frac{1}{n} \sum_{i=1}^n F(x_i, U_i) \xrightarrow{a.s.} f(x)$ ,  
 under the assumption that  $\mathbb{E}U = 0$ .

$$\Rightarrow x_{n+1} = x_n - \frac{\gamma}{m} \sum_{i=1}^m F(x_i, U_i)$$

↳ Robbins & Monro argued that taking an excessive  
 number of observations at each iteration is inefficient,  
 since  $x_n$  is used only as an intermediate step in  
 the approximation of the root. They suggested to  
 consider instead:

$$x_{n+1} = x_n - \gamma_n F(x_n, U_{n+1})$$

The sequence  $\gamma_n > 0$  converges to  
 zero as  $n \rightarrow \infty$ , and is such that  
 $\sum \gamma_n = \infty$ . The decreasing step  
 size actually provides an averaging  
 of the observations, required to  
 prove that the resulting sequence  
 $x_n$  converges to the root  $x^*$  of  $f$ .

Example:

Let  $U_1, \dots, U_{n+1}$  = sequence  
 of iid RVs with mean  $\mathbb{E}U$ .

$$\begin{aligned} \bar{X}_{n+1} &:= \frac{1}{n+1} \sum_{i=1}^{n+1} U_i \\ &= \frac{n}{n+1} \frac{1}{n} \sum_{i=1}^n U_i + \frac{1}{n+1} U_{n+1} \end{aligned}$$

$$\bar{X}_{n+1} = \left(\frac{n}{n+1}\right) \bar{X}_n + \frac{1}{n+1} U_{n+1} \quad (c)$$

$$= \bar{X}_n - \frac{1}{n+1} (\bar{X}_n - U_{n+1}) = \bar{X}_n - \gamma_n F(\bar{X}_n, U_{n+1}),$$

with  $\gamma_n = \frac{1}{n+1}$  (note that  $\sum \gamma_n = \infty$ ; while  $\sum \gamma_n^2 < \infty$ )

$$\begin{aligned} \bullet F(x, u) &= x - u \\ &= (x - \mathbb{E}U) + (\mathbb{E}U - u) \quad \forall x, \text{ you do not observe } x - \mathbb{E}U, \\ &= f(x) + \underbrace{(\mathbb{E}U - u)}_{\text{noise term}} \rightarrow \text{but only } x - u. \end{aligned}$$

$$\bullet f(x) = x - \mathbb{E}U$$

And indeed,  $\bar{X}_{n+1} \xrightarrow{a.s.} \mathbb{E}U$  (SLLN); corresponding to the root of  $f$ .

$\Rightarrow$  We are interested in the convergence of the sequence of estimates  $X_{n+1} = X_n - \gamma_n F(X_n, U_{n+1})$  (sequence of RVs). Put  $f(x) = \mathbb{E}[F(x, U)]$ . We show that under some conditions on  $F$  and  $f$ , the sequence converges (in some sense) to the unique root  $x^*$  of  $f$ .

### I.2. The Robbins-Monro algorithm.

We first establish convergence under restrictive assumptions on  $f$  and  $F$ . The conditions are later weakened.

Theorem. Let  $\bullet F: \mathbb{R}^2 \rightarrow \mathbb{R}$  measurable

$\bullet U$ : a RV with distribution  $G$ .

Consider the sequence of RVs  $X_0, X_1, \dots$ , satisfying

$$\rightarrow X_0 \in \mathcal{L}^2$$

$$\rightarrow X_{n+1} = X_n - \gamma_n F(X_n, U_{n+1}), \text{ where } \{U_n\} = \text{i.i.d with distribution } G.$$

Under the assumptions that

$$[H1] \exists K > 0 \quad \forall x \in \mathbb{R} \quad |F(x, u)| \leq K \\ \forall u \in \mathbb{R}$$

$$[H2] f = \text{continuous, strictly increasing. Moreover, there exists a unique } x^* \in \mathbb{R} \text{ such that } f(x^*) = 0$$

$$[H3] \forall n \quad \gamma_n > 0; \quad \sum \gamma_n = \infty; \quad \sum \gamma_n^2 < \infty$$

Then  $X_n \xrightarrow{a.s.} x^*$

proof = Step 1: Put  $\mathcal{F}_n = \sigma(X_0, X_1, \dots, X_n)$

$$S_n = \sum_{k=0}^{n-1} \gamma_k^2 \mathbb{E}\{F^2(X_k, U_{k+1}) | \mathcal{F}_k\}$$

We show that the sequence  $M_n = (X_n - x^*)^2 - S_n$  is an  $\mathcal{F}_n$ -supermartingale.

$$\begin{aligned} \mathbb{E}(M_{n+1} | \mathcal{F}_n) &= \mathbb{E}\{(X_{n+1} - x^*)^2 | \mathcal{F}_n\} - \mathbb{E}\{S_{n+1} | \mathcal{F}_n\} \\ &= \mathbb{E}\{(X_n - \gamma_n F(X_n, U_{n+1}) - x^*)^2 | \mathcal{F}_n\} \\ &\quad - \mathbb{E}\{S_{n+1} | \mathcal{F}_n\} \end{aligned}$$

$$\begin{aligned} &= (X_n - x^*)^2 - 2(X_n - x^*) \gamma_n \mathbb{E}\{F(X_n, U_{n+1}) | \mathcal{F}_n\} \\ &\quad + \gamma_n^2 \mathbb{E}\{F^2(X_n, U_{n+1}) | \mathcal{F}_n\} \end{aligned}$$

On the other hand,

$$\begin{aligned} \mathbb{E}\{S_{n+1} | \mathcal{F}_n\} &= \mathbb{E}\left\{\sum_{k=0}^n \gamma_k^2 \mathbb{E}\{F^2(X_k, U_{k+1}) | \mathcal{F}_k\} \mid \mathcal{F}_n\right\} \\ &= S_n + \gamma_n^2 \mathbb{E}\{F^2(X_n, U_{n+1}) | \mathcal{F}_n\} \end{aligned}$$

Thus

$$\begin{aligned} \mathbb{E}\{M_{n+1} | \mathcal{F}_n\} &= (X_n - x^*)^2 - 2(X_n - x^*) \gamma_n \mathbb{E}\{F(X_n, U_{n+1}) | \mathcal{F}_n\} \\ &\quad - S_n \\ &= M_n - 2(X_n - x^*) \gamma_n \mathbb{E}\{F(X_n, U_{n+1}) | \mathcal{F}_n\} \end{aligned}$$



In addition,  $(X_n - x^*) \mathbb{E}\{F(X_n, U_{n+1}) | \mathcal{F}_n\}$  (e)

$$= (X_n - x^*) \mathbb{E}\{F(X_n, U_{n+1}) - F(x^*, U_{n+1}) | \mathcal{F}_n\}$$

$$= (X_n - x^*) (f(X_n) - f(x^*))$$

↑  
same sign, since  $f$  is strictly positive.

Thus  $\mathbb{E}(M_{n+1} | \mathcal{F}_n) \leq M_n$  i.e.  $M_n = \overline{SMG}$ .

• Step 2: The sequence  $\{M_n\}$  converges almost surely and in mean to some  $M_\infty \in \mathcal{L}^1$ .

$$M_n = (X_n - x^*)^2 - S_n = M_n^+ - M_n^-$$

where  $M_n^+ = \max(0, M_n)$

$$M_n^- = -\min(0, M_n)$$

Since  $(X_n - x^*)^2$  and  $S_n$  are positive, we conclude that  $M_n^+ \leq (X_n - x^*)^2$  and  $M_n^- \leq S_n$

Thus  $\mathbb{E}M_n^- \leq \mathbb{E}S_n \leq K \sum_{k=0}^{\infty} \gamma_k^2 < +\infty$  under [H3]

Doob's first martingale convergence theorem implies that  $M_n \xrightarrow{a.s.} M_\infty \in \mathcal{L}^1$

Moreover,  $M_n = \overline{SMG} \Rightarrow \mathbb{E}M_{n+1} \leq \mathbb{E}M_n$ , and

$$\mathbb{E}M_n \geq -K \sum_{k=0}^{\infty} \gamma_k^2$$

$\Rightarrow$  sequence  $\{\mathbb{E}M_n\}$  is non-increasing & bounded below  $\Rightarrow$  it converges to  $\mathbb{E}M_\infty$ .

• Step 3:  $\sum_{k=0}^{\infty} \gamma_k \mathbb{E}\{f(X_k)(X_k - x^*)\} < +\infty$

Indeed,

$$\mathbb{E}\{f(X_n)(X_n - x^*)\} = \mathbb{E}\left[\mathbb{E}\{(X_n - x^*)F(X_n, U_{n+1}) | \mathcal{F}_n\}\right]$$

↑  
definition of  $f$

and  $\gamma_n \mathbb{E}\{(X_n - x^*)F(X_n, U_{n+1}) | \mathcal{F}_{n+1}\}$  (f)

$$= \frac{1}{2} (M_n - \mathbb{E}\{M_{n+1} | \mathcal{F}_n\})$$

$$\Rightarrow \gamma_n \mathbb{E}\{f(X_n)(X_n - x^*)\} = \frac{1}{2} (\mathbb{E}M_n - \mathbb{E}M_{n+1})$$

$$\Rightarrow \sum_{k=0}^n \gamma_k \mathbb{E}\{f(X_k)(X_k - x^*)\} = \frac{1}{2} (\mathbb{E}M_0 - \mathbb{E}M_{n+1})$$

$$\rightarrow \frac{1}{2} (\mathbb{E}M_0 - \mathbb{E}M_\infty)$$

from Step 2.

• Step 4: Conclusion.

First, note that since  $M_n \xrightarrow{a.s.} M_\infty$  and

$$0 \leq S_n = \sum_{k=0}^{n-1} \gamma_k^2 \mathbb{E}\{F^2(X_k, U_{k+1}) | \mathcal{F}_k\} \leq S_{n+1} \leq K \underbrace{\sum_{k=0}^{\infty} \gamma_k^2}_{< +\infty}$$

$\Rightarrow S_n$  converges.

$\Rightarrow (X_n - x^*)^2$  converges almost surely to some RV  $L$ .

It remains to show that  $L \equiv 0$ .

Suppose  $\exists \omega \in \Omega$  such that  $L(\omega) > 0$ .

$\Rightarrow$  there exists  $m_1 > 0$  and  $n_0$  such that  $\forall n \geq n_0$ ,

$$|X_n - x^*|(\omega) \geq m_1$$

Since  $f$  is continuous, and  $x^*$  such that  $f(x^*) = 0$ , there exists  $m_2 > 0$  such that  $\forall n \geq n_0$ ,

$f(X_n)(X_n - x^*) \geq m_2$ , which implies that

$$\sum_{k=0}^{\infty} \gamma_k \mathbb{E}f(X_k)(X_k - x^*) \geq m_2 \sum_{k=0}^{\infty} \gamma_k = +\infty$$

↑  
[H3]

which contradicts what we established in Step 3.

We conclude that  $L \equiv 0$  and that  $X_n \xrightarrow{a.s.} x^*$  ▣

Remarks: (i) Assumption [H1] can be replaced with the weaker assumption that (g)

$$[H1'] \exists K > 0 \forall x \in \mathbb{R}^d \mathbb{E}\{\|F(x, U)\|^2\} \leq K(1 + \|x\|^2)$$

for  $F: \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}$

Almost sure convergence still holds, replacing [H1] with [H1'].

The idea is to show that the sequence

$$M_n = B_{n-1} \|X_n - x^*\|^2 - \sum_{k=0}^{n-1} B_k \delta_k (K \gamma_k - 2f(X_k)(X_k - x^*)),$$

with  $B_n := \prod_{i=1}^{n-1} \frac{1}{1 + K \gamma_i^2}$ , is a supermartingale, bounded in  $\mathcal{L}^1$ , and then proceed as before.

(ii) Assumption [H2] can be weakened either to

$$[H2'] \exists x^* \in \mathbb{R}^d \text{ \& } c > 0 \text{ s.t. } \forall x \in \mathbb{R}^d f(x) \cdot (x - x^*) \geq c \|x - x^*\|^2$$

or

$$[H2''] \exists x^* \in \mathbb{R}^d \text{ s.t. } f(x^*) = 0 \text{ and } \forall x \neq x^*, f(x) \cdot (x - x^*) > 0.$$

Then it is possible to show that under  
 [H1'] + [H2'] + [H3],  $X_n \xrightarrow{\mathcal{L}^2} x^*$  and  
 [H1'] + [H2''] + [H3],  $X_n \xrightarrow{\text{a.s.}} x^*$ .

Note that [H2''] is weaker than [H2'], and that these two assumptions imply that  $x^*$  is the unique root of the equation  $f(x) = 0$ .

(iii) Assumption  $\sum_{k \geq 0} \delta_k^2 < \infty$  in [H3] can also be weakened considerably.

### I.3. Applications.

(h)

(a) Quantile estimation.

Take  $F(x, u) = \mathbb{1}(u \leq x) - \alpha$ , for  $\alpha \in (0, 1)$ .

$G$  continuous (so that the quantile is well defined).

$$\text{Then } f(x) = \mathbb{E}[F(x, U)] = \mathbb{P}(U \leq x) - \alpha = G(x) - \alpha.$$

The root of this function is precisely  $q_\alpha$ :  $\alpha$ -quantile of  $G$ ;  $G(q_\alpha) = \alpha$ .

$\Rightarrow$  Use the recursion  $X_{n+1} = X_n - \gamma_n (\mathbb{1}(U_{n+1} \leq X_n) - \alpha)$ .

(b) Maximum likelihood estimation.

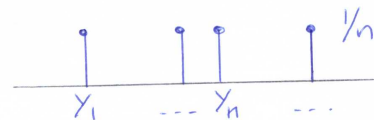
Let  $Y_1, \dots, Y_n$  = iid observation with distribution  $P_\theta$ , and density  $g_\theta$  (twice differentiable)

The MLE is typically solution of an equation

$$\frac{1}{n} \sum_{i=1}^n F(\theta, Y_i) = 0,$$

where  $F(\theta, y) := g'_\theta(y) / g_\theta(y)$ .

Take  $* \mathbb{P}_n = \frac{1}{n} \sum_{i=1}^n \delta_{Y_i}$  = empirical probability of obs  $Y_1, \dots, Y_n$



$$* f(\theta) = \mathbb{E}_{\mathbb{P}_n} F(\theta, U) \text{ , where } U \sim \mathbb{P}_n.$$

$\Rightarrow$  The sequence  $\theta_{m+1} = \theta_m - \gamma_m F(\theta_m, U_{m+1}) \xrightarrow[m \rightarrow \infty]{\text{a.s.}} \text{MLE } \hat{\theta}_{\text{MLE}}$   
 drawn uniformly in  $\{Y_1, \dots, Y_n\}$  at each iteration