

# First-Order Optimization

Optimization Techniques (ENGG\*6140)

School of Engineering,  
University of Guelph, ON, Canada

Course Instructor: Benyamin Ghojogh  
Winter 2023

## Fundamental Theorem of Calculus

# Fundamental theorem of calculus

## Lemma (Fundamental theorem of calculus for multivariate functions)

Consider a differentiable function  $f(\cdot)$  with domain  $\mathcal{D}$ . For any  $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ , we have:

$$\begin{aligned} \star \quad f(\mathbf{y}) &= f(\mathbf{x}) + \underbrace{\nabla f(\mathbf{x})^\top}_{\text{gradient at } \mathbf{x}} \underbrace{(\mathbf{y} - \mathbf{x})}_{\text{direction}} + \int_0^1 \underbrace{\left( \underbrace{\nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))}_{\text{intermediate gradient}} - \nabla f(\mathbf{x}) \right)^\top}_{\text{error term}} \underbrace{(\mathbf{y} - \mathbf{x})}_{\text{direction}} dt \\ &= \underbrace{f(\mathbf{x})}_{\text{value at } \mathbf{x}} + \underbrace{\nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})}_{\text{linear approximation}} + \underbrace{o(\mathbf{y} - \mathbf{x})}_{\text{higher-order terms}}, \end{aligned} \quad (1)$$

where  $o(\cdot)$  is the small-o complexity.

## Lemma (Corollary of the fundamental theorem of calculus)

Consider a differentiable function  $f(\cdot)$ , with domain  $\mathcal{D}$ , whose gradient is L-smooth:

$$\boxed{|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})| \leq L \|\mathbf{x} - \mathbf{y}\|_2, \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{D}.} \quad (2)$$

For any  $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ , we have:

$$\boxed{f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|_2^2.} \quad (3)$$

~~$$|f(\mathbf{x}) - f(\mathbf{y})| \leq L \|\mathbf{x} - \mathbf{y}\|_2$$~~

# Fundamental theorem of calculus

$$\|a\|_2 \|y\|_2$$

Proof for the corollary of the fundamental theorem of calculus:

$$\begin{aligned}
 f(\mathbf{y}) &\stackrel{(1)}{=} f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \int_0^1 (\nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}))^\top (\mathbf{y} - \mathbf{x}) dt \\
 &\stackrel{(a)}{\leq} f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \int_0^1 \|\nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x})\|_2 \|\mathbf{y} - \mathbf{x}\|_2 dt \\
 &\stackrel{(b)}{\leq} f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \int_0^1 L t \|\mathbf{y} - \mathbf{x}\|_2^2 dt \\
 &= f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \underbrace{L \|\mathbf{y} - \mathbf{x}\|_2^2}_{=1} \int_0^1 t dt \\
 &= f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|_2^2,
 \end{aligned}$$

where (a) is because of the Cauchy-Schwarz inequality and (b) is because, according to Eq. (2),

$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|_2, \forall \mathbf{x}, \mathbf{y} \in \mathcal{D}$ , we have:

$$\|\nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x})\|_2 \leq L \|\mathbf{x} + t(\mathbf{y} - \mathbf{x}) - \mathbf{x}\|_2 = Lt \|\mathbf{y} - \mathbf{x}\|_2. \text{ Q.E.D.}$$



# Gradient descent: introduction

- Gradient descent is one of the fundamental first-order methods.
- It was first suggested by Cauchy in 1874 [1] and Hadamard in 1908 [2] and its convergence was later analyzed in [3].
- Unconstrained optimization:

$$\boxed{\underset{x}{\text{minimize}} \quad f(x)}. \quad (4)$$

- In numerical optimization for unconstrained optimization, we start with a random feasible initial point and iteratively update it by step  $\Delta x$ :

$$\star \quad \boxed{x^{(k+1)} := x^{(k)} + \Delta x.} \quad \left. \begin{array}{l} x^{(1)} \leftarrow x^{(0)} + \Delta x^{(0)} \\ x^{(2)} \leftarrow x^{(1)} + \Delta x^{(1)} \\ \vdots \end{array} \right\} (5)$$

- Continue until we converge to (or get sufficiently close to) the desired optimal point  $x^*$ .
- The step  $\Delta x$  is also denoted by  $p$  in the literature, i.e.,  $p := \Delta x$ .

$$\left. \begin{array}{l} x^{(1)} \quad 2.3 \\ x^{(2)} \quad 1.5 \quad \dots \quad x^{(1000)} \quad 1.0001 \\ x^* = 1 \leftarrow x^{(1000)} \end{array} \right\}$$

## Gradient descent: step size

- Let the function  $f(\cdot)$  be differentiable and its gradient be  $L$ -smooth.
- Recall Eq. (3) from the corollary of the fundamental theorem of calculus:

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|_2^2. \quad \star$$

If we set  $\mathbf{x} = \mathbf{x}^{(k)}$  and  $\mathbf{y} = \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \Delta \mathbf{x}$  in this equation, we have:

$$f(\mathbf{x}^{(k)} + \Delta \mathbf{x}) \leq f(\mathbf{x}^{(k)}) + \nabla f(\mathbf{x}^{(k)})^\top \Delta \mathbf{x} + \frac{L}{2} \|\Delta \mathbf{x}\|_2^2$$

$$\Rightarrow \underbrace{f(\mathbf{x}^{(k)} + \Delta \mathbf{x}) - f(\mathbf{x}^{(k)})}_{\star} \leq \nabla f(\mathbf{x}^{(k)})^\top \Delta \mathbf{x} + \frac{L}{2} \|\Delta \mathbf{x}\|_2^2. \quad \star \quad \underline{(6)}$$

- Until reaching the minimum, we want to decrease the cost function  $f(\cdot)$  in every iteration; hence, we desire:

$$\boxed{f(\mathbf{x}^{(k)} + \Delta \mathbf{x}) - f(\mathbf{x}^{(k)})} < 0. \rightarrow f(\underbrace{\mathbf{x}^{(k)} + \Delta \mathbf{x}}_{\mathbf{x}^{(k+1)}}) < f(\mathbf{x}^{(k)}) \quad (7)$$

- According to Eq. (6), one way to achieve Eq. (7) is:

$$\star \quad \boxed{\nabla f(\mathbf{x}^{(k)})^\top \Delta \mathbf{x} + \frac{L}{2} \|\Delta \mathbf{x}\|_2^2 < 0.}$$

# Gradient descent: step size

- We have:

$$\star \underbrace{\nabla f(\mathbf{x}^{(k)})^\top \Delta \mathbf{x} + \frac{L}{2} \|\Delta \mathbf{x}\|_2^2}_{\text{minimize}} < 0.$$

- Hence, we should minimize  $\nabla f(\mathbf{x}^{(k)})^\top \Delta \mathbf{x} + \frac{L}{2} \|\Delta \mathbf{x}\|_2^2$  w.r.t.  $\Delta \mathbf{x}$ :

$$\underbrace{\text{minimize}_{\Delta \mathbf{x}} \nabla f(\mathbf{x}^{(k)})^\top \Delta \mathbf{x} + \frac{L}{2} \|\Delta \mathbf{x}\|_2^2}_{\star} \quad (8)$$

- This function is convex w.r.t.  $\Delta \mathbf{x}$  and we can optimize it by setting its derivative to zero:

$$\frac{\partial}{\partial \Delta \mathbf{x}} \left( \nabla f(\mathbf{x}^{(k)})^\top \Delta \mathbf{x} + \frac{L}{2} \|\Delta \mathbf{x}\|_2^2 \right) = \underbrace{\nabla f(\mathbf{x}^{(k)})}_{\Delta \mathbf{x}^\top \Delta \mathbf{x}} + \underbrace{L \Delta \mathbf{x}}_{\text{set } 0} = 0$$

$$\Rightarrow \Delta \mathbf{x} = -\frac{1}{L} \nabla f(\mathbf{x}^{(k)}).$$

(9)

- So, we have:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \frac{1}{L} \nabla f(\mathbf{x}^{(k)}).$$

## Gradient descent: step size

- Recall Eqs. (9) and (6):

$$\star \Delta \mathbf{x} = -\frac{1}{L} \nabla f(\mathbf{x}^{(k)}).$$

$$\star f(\mathbf{x}^{(k)} + \Delta \mathbf{x}) - f(\mathbf{x}^{(k)}) \leq \nabla f(\mathbf{x}^{(k)})^\top \Delta \mathbf{x} + \frac{L}{2} \|\Delta \mathbf{x}\|_2^2.$$

- Using the first equation in the second equation gives:

$$\begin{aligned} f(\mathbf{x}^{(k)} + \Delta \mathbf{x}) - f(\mathbf{x}^{(k)}) &\leq \underbrace{-\frac{1}{L} \nabla f(\mathbf{x}^{(k)})^\top \nabla f(\mathbf{x}^{(k)})}_{\|\nabla f(\mathbf{x}^{(k)})\|_2^2} + \frac{L}{2} \left\| \frac{-1}{L} \nabla f(\mathbf{x}^{(k)}) \right\|_2^2 \\ &= -\frac{1}{L} \|\nabla f(\mathbf{x}^{(k)})\|_2^2 + \frac{L}{2} \left\| \frac{-1}{L} \nabla f(\mathbf{x}^{(k)}) \right\|_2^2 \\ &= -\frac{1}{L} \|\nabla f(\mathbf{x}^{(k)})\|_2^2 + \frac{L}{2} \left( \frac{1}{L} \right)^2 \|\nabla f(\mathbf{x}^{(k)})\|_2^2 \\ &= -\frac{1}{2L} \|\nabla f(\mathbf{x}^{(k)})\|_2^2 \leq 0, \quad \star \end{aligned}$$

$$\mathbf{x}^\top \mathbf{x} = \|\mathbf{x}\|_2^2$$

which satisfies Eq. (7). Eq. (9) means that it is better to move toward a scale of minus gradient for updating the solution. This inspires the name of algorithm which is *gradient descent*.

## Gradient descent: step size

- The problem is that often we either do not know the Lipschitz constant  $L$  or it is hard to compute. Hence, rather than Eq. (9),  $\Delta \mathbf{x} = \ominus \frac{1}{L} \nabla f(\mathbf{x}^{(k)})$ , we use:

step  $\rightarrow$   $\Delta \mathbf{x} = \ominus \eta \nabla f(\mathbf{x}^{(k)})$ , i.e.,  $\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} \ominus \eta \nabla f(\mathbf{x}^{(k)})$ , (10)

where  $\eta > 0$  is the step size, also called the learning rate in data science literature.

- If the optimization problem is maximization rather than minimization, the step should be  $\Delta \mathbf{x} = \eta \nabla f(\mathbf{x}^{(k)})$ , rather than Eq. (10). In that case, the name of method is gradient ascent.
- Recall Eq. (6):

$\star$   $f(\mathbf{x}^{(k)} + \Delta \mathbf{x}) - f(\mathbf{x}^{(k)}) \leq \nabla f(\mathbf{x}^{(k)})^\top \Delta \mathbf{x} + \frac{L}{2} \|\Delta \mathbf{x}\|_2^2$

$\nabla f(\mathbf{x}^{(k+1)})^\top \nabla f(\mathbf{x}^{(k)}) = \|\nabla f(\mathbf{x}^{(k)})\|_2^2$

Using Eq. (10) in this equation gives:

$$\begin{aligned} f(\mathbf{x}^{(k)} + \Delta \mathbf{x}) - f(\mathbf{x}^{(k)}) &\leq \ominus \eta \|\nabla f(\mathbf{x}^{(k)})\|_2^2 + \frac{L}{2} \eta^2 \|\nabla f(\mathbf{x}^{(k)})\|_2^2 \\ &= \eta \left( \frac{L}{2} \eta - 1 \right) \|\nabla f(\mathbf{x}^{(k)})\|_2^2 \end{aligned} \quad (11)$$

# Gradient descent: step size

- We found:

$$f(x^{(k)} + \Delta x) - f(x^{(k)}) \leq \eta \left( \frac{L}{2} \eta - 1 \right) \|\nabla f(x^{(k)})\|_2^2$$

*Handwritten notes:  $\eta > 0$ ,  $\frac{L}{2}\eta - 1$  must be  $< 0$ ,  $\|\nabla f(x)\|_2^2 > 0$*

- If  $x^{(k)}$  is not a stationary point, we have  $\|\nabla f(x^{(k)})\|_2^2 > 0$ .
- Recall Eq. (7):

$$f(x^{(k)} + \Delta x) - f(x^{(k)}) < 0$$

*Handwritten note:  $\geq 0$*

- Noticing  $\eta > 0$ , for satisfying this equation, we must set:

$$\frac{L}{2} \eta - 1 < 0 \implies \eta < \frac{2}{L} \quad 1 \leq 2 \quad (12)$$

$$f(x^{(k)} + \Delta x) - f(x^{(k)}) = 1 \implies f(x^{(k+1)}) = f(x^{(k)}) + 1$$

*Handwritten notes:  $x^{(k+1)}$ ,  $x^{(k)}$ ,  $\star$*

## Gradient descent: step size

- Recall Eq. (11):

$$f(\mathbf{x}^{(k)} + \Delta \mathbf{x}) - f(\mathbf{x}^{(k)}) \leq -\eta \|\nabla f(\mathbf{x}^{(k)})\|_2^2 + \frac{L}{2} \eta^2 \|\nabla f(\mathbf{x}^{(k)})\|_2^2 \quad *$$

*Handwritten annotations:*  $\triangleleft \alpha$  (pointing to the first term), upper bound (pointing to the second term)

- On the other hand, we can minimize this equation by setting its derivative w.r.t.  $\eta$  to zero:

$$\begin{aligned} \frac{\partial}{\partial \eta} \left( -\eta \|\nabla f(\mathbf{x}^{(k)})\|_2^2 + \frac{L}{2} \eta^2 \|\nabla f(\mathbf{x}^{(k)})\|_2^2 \right) &= -\|\nabla f(\mathbf{x}^{(k)})\|_2^2 + L\eta \|\nabla f(\mathbf{x}^{(k)})\|_2^2 \\ &= (-1 + L\eta) \|\nabla f(\mathbf{x}^{(k)})\|_2^2 \stackrel{\text{set}}{=} 0 \implies \eta = \frac{1}{L}. \end{aligned} \quad *$$

If we set:

$$\eta < \frac{1}{L}, \quad (13)$$

then Eq. (11) becomes:

$$\begin{aligned} f(\mathbf{x}^{(k)} + \Delta \mathbf{x}) - f(\mathbf{x}^{(k)}) &\leq -\frac{1}{L} \|\nabla f(\mathbf{x}^{(k)})\|_2^2 + \frac{1}{2L} \|\nabla f(\mathbf{x}^{(k)})\|_2^2 = -\frac{1}{2L} \|\nabla f(\mathbf{x}^{(k)})\|_2^2 < 0 \\ \implies f(\mathbf{x}^{(k+1)}) &\leq f(\mathbf{x}^{(k)}) - \frac{1}{2L} \|\nabla f(\mathbf{x}^{(k)})\|_2^2. \end{aligned} \quad (14)$$

# Gradient descent: step size

- Recall Eqs. (13) and (14):

$$\eta < \frac{1}{L}, \quad (13)$$
$$f(\mathbf{x}^{(k+1)}) \leq f(\mathbf{x}^{(k)}) - \frac{1}{2L} \|\nabla f(\mathbf{x}^{(k)})\|_2^2. \quad (14)$$

- Eq. (13) means that there should be an upper-bound, dependent on the Lipschitz constant, on the step size. Hence,  $L$  is still required.
- Eq. (14) shows that every iteration of gradient descent decreases the cost function:

$$f(\mathbf{x}^{(k+1)}) \leq f(\mathbf{x}^{(k)}), \quad (15)$$

and the amount of this decrease depends on the norm of gradient at that iteration.

- So, the series of solutions converges to the optimal solution while the function value decreases iteratively until the local minimum:



$$\{\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots\} \rightarrow \mathbf{x}^*,$$
$$f(\mathbf{x}^{(0)}) \geq f(\mathbf{x}^{(1)}) \geq f(\mathbf{x}^{(2)}) \geq \dots \geq f(\mathbf{x}^*).$$

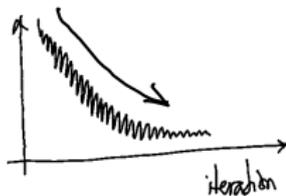
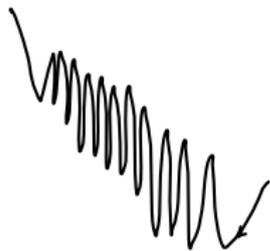
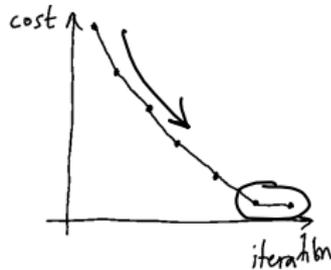
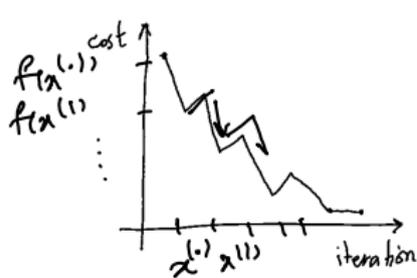


- If the optimization problem is a convex problem, the solution is the global solution; otherwise, the solution is local.

# Gradient descent: cost versus iterations

$$\{x^{(0)}, x^{(1)}, x^{(2)}, \dots\} \rightarrow x^*,$$

$$\underbrace{f(x^{(0)})} \geq \underbrace{f(x^{(1)})} \geq \underbrace{f(x^{(2)})} \geq \dots \geq \underbrace{f(x^*)}.$$



## Line-Search

# Line-search

$$\eta < \frac{1}{L} \rightarrow \eta < 1$$

- We saw the step size of gradient descent requires knowledge of the Lipschitz constant for the smoothness of gradient. However, we may not know the exact Lipschitz constant. Hence, we can find the suitable step size  $\eta$  by a search which is named the line-search.
- In line-search of every optimization iteration, we start with  $\eta = 1$  and if it does not satisfy Eq. (7) with step  $\Delta \mathbf{x} = -\eta \nabla f(\mathbf{x}^{(k)})$ :

$$\boxed{f(\mathbf{x}^{(k)} + \Delta \mathbf{x}) < f(\mathbf{x}^{(k)})} \Rightarrow \boxed{f(\mathbf{x}^{(k)} - \eta \nabla f(\mathbf{x}^{(k)})) < f(\mathbf{x}^{(k)})}, \quad (16)$$

we halve it,  $\eta \leftarrow \eta/2$ .

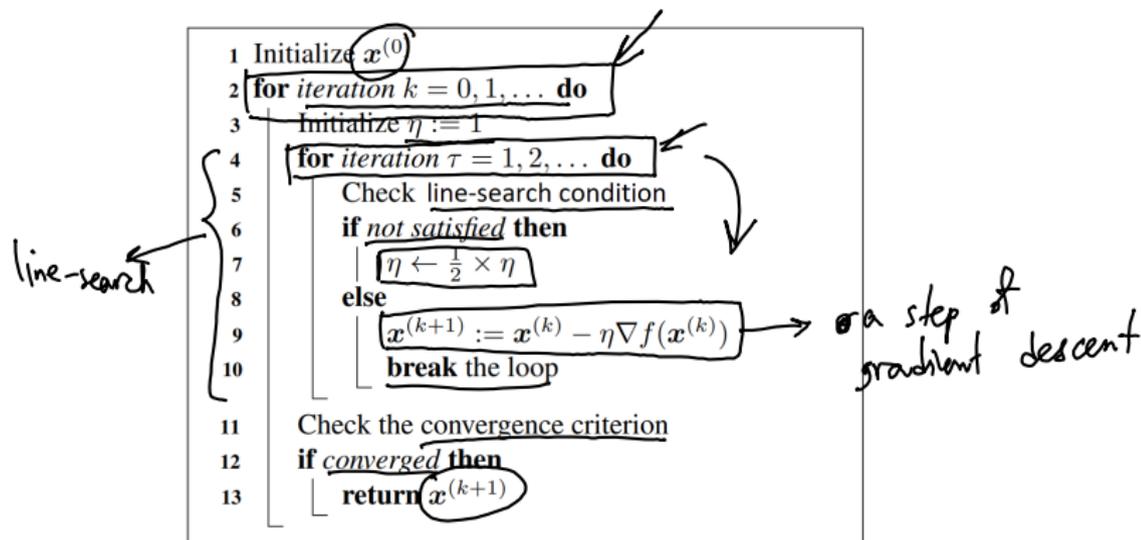
- This halving step size is repeated until this equation is satisfied, i.e., until we have a decrease in the objective function. Note that this decrease will happen when the step size becomes small enough to satisfy Eq. (13):

$$\boxed{\eta < \frac{1}{L}}$$

The algorithm of gradient descent with line-search is shown in Algorithm ?? . As this algorithm shows, line-search has its own internal iterations inside every iteration of gradient descent.

# Gradient descent with line-search

The algorithm of gradient descent with line-search:



As this algorithm shows, line-search has its own internal iterations inside every iteration of gradient descent.

# Backtracking line-search

- A more sophisticated line-search method is the Armijo line-search [4], also called the backtracking line-search. Rather than Eq. (16), it checks if the cost function is *sufficiently* decreased:

$$f(\mathbf{x}^{(k)} + \mathbf{p}) \leq f(\mathbf{x}^{(k)}) + c \mathbf{p}^\top \nabla f(\mathbf{x}^{(k)}), \quad (17)$$

*Handwritten annotations:* A box is drawn around the right-hand side of the inequality. An arrow labeled 'b' points to the term  $c \mathbf{p}^\top \nabla f(\mathbf{x}^{(k)})$ . Another arrow labeled 'b-3' points to the same term. A star is drawn to the left of the inequality, and another star is to the right. A small arrow labeled  $\Delta \mathbf{x}$  points to  $\mathbf{p}$ .

where  $c \in (0, 0.5]$  is the parameter of Armijo line-search and  $\mathbf{p} = \Delta \mathbf{x}$  is the search direction for update.

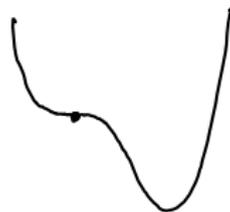
- The value of  $c$  should be small, e.g.,  $c = 10^{-4}$  [5].
- This condition is called the Armijo condition or the Armijo-Goldstein condition.
- In gradient descent, the search direction is  $\mathbf{p} = \Delta \mathbf{x} = -\eta \nabla f(\mathbf{x}^{(k)})$  according to Eq. (10). Hence, for gradient descent, it checks:

$$f(\mathbf{x}^{(k)} - \eta \nabla f(\mathbf{x}^{(k)})) \leq f(\mathbf{x}^{(k)}) - \eta c \|\nabla f(\mathbf{x}^{(k)})\|_2^2. \quad (18)$$

*Handwritten annotation:* A box is drawn around the entire inequality.

- Another more sophisticated line-search is Wolfe conditions [6]. We will learn it later in the course.

$$\|\nabla f(x^{(k)})\|_2 < \varepsilon$$

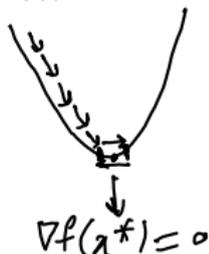


Convergence criterion

# Convergence criteria

- For all numerical optimization methods including gradient descent, there exist several methods for convergence criterion to stop updating the solution and terminate optimization.
- Some of them are:
  - Small norm of gradient:

$$\|\nabla f(\mathbf{x}^{(k+1)})\|_2 \leq \epsilon, \quad *$$



where  $\epsilon$  is a small positive number.

- ★ The reason for this criterion is the first-order optimality condition (recall that at the local optimum, we have  $\|\nabla f(\mathbf{x}^*)\|_2 = 0$ ).
  - ★ If the function is not convex, this criterion has the risk of stopping at a saddle point.
- Small change of cost function:

$$|f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^{(k)})| \leq \epsilon.$$

- Small change of gradient of function:

$$\|\nabla f(\mathbf{x}^{(k+1)}) - \nabla f(\mathbf{x}^{(k)})\| \leq \epsilon.$$

- Reaching maximum desired number of iterations, denoted by

$$k < \max_k \quad *$$

## Momentum



# Gradient descent with momentum

- Gradient descent and other first-order methods can have a momentum term. **Momentum**, proposed in [7], makes the change of solution  $\Delta \mathbf{x}$  a little similar to the previous change of solution.
- Hence, the change adds a history of previous change to Eq. (10):

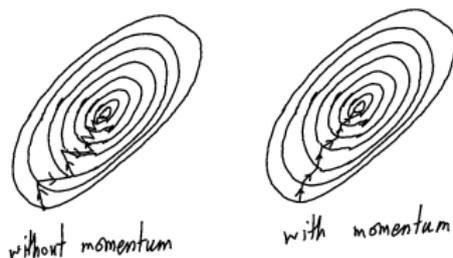
$$(\Delta \mathbf{x})^{(k)} := \alpha(\Delta \mathbf{x})^{(k-1)} - \eta^{(k)} \nabla f(\mathbf{x}^{(k)}), \quad (19)$$

where  $\alpha > 0$  is the momentum parameter which weights the importance of history compared to the descent direction.

- We use this  $(\Delta \mathbf{x})^{(k)}$  in Eq. (5) for updating the solution:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} + (\Delta \mathbf{x})^{(k)}.$$

- Because of faithfulness to the track of previous updates, momentum reduces the amount of oscillation of updates in gradient descent optimization.



## Steepest Descent

# Steepest Descent

- **Steepest descent** is similar to gradient descent but there is a difference between them.
- In steepest descent, we move toward the negative gradient as much as possible to reach the smallest function value which can be achieved at every iteration.
- Hence, the step size at iteration  $k$  of steepest descent is calculated as [8]:

$$\eta^{(k)} := \arg \min_{\eta} f(\mathbf{x}^{(k)} - \eta \nabla f(\mathbf{x}^{(k)})), \quad (20)$$

and then, the solution is updated using Eq. (10) as in gradient descent:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \eta \nabla f(\mathbf{x}^{(k)}).$$

- Another interpretation of steepest descent [9, Chapter 9.4]: The first-order Taylor expansion of function is  $f(\mathbf{x} + \mathbf{v}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^\top \mathbf{v}$ . Hence, the step size in the normalized steepest descent, at iteration  $k$ , is obtained as:

$$\Delta \mathbf{x} = \arg \min_{\mathbf{v}} \{ \nabla f(\mathbf{x}^{(k)})^\top \mathbf{v} \mid \|\mathbf{v}\|_2 \leq 1 \}, \quad (21)$$

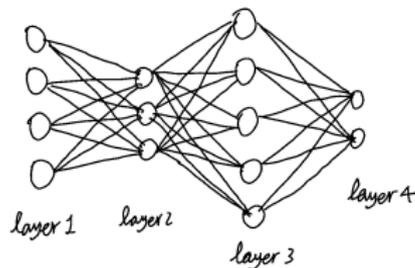
which is used in Eq. (5) for updating the solution:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} + \Delta \mathbf{x}.$$

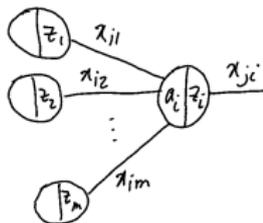
## Backpropagation

# Neural network

- Neural network:



- Every neuron in neural network:



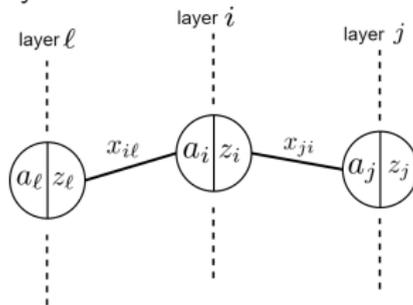
- Let  $x_{ji}$  denote the weight connecting neuron  $i$  to neuron  $j$ . Let  $a_i$  and  $z_i$  be the output of neuron  $i$  before and after applying its activation function  $\sigma_i(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$ , respectively.

$$a_i = \sum_{\ell=1}^m x_{i\ell} z_\ell,$$

$$z_i := \sigma_i(a_i).$$

# Backpropagation

- Consider three neurons in three layers of a network:



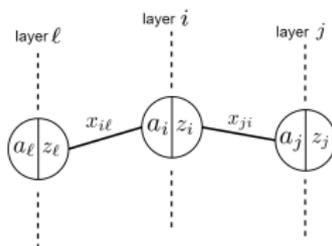
- We have  $a_i = \sum_\ell x_{i\ell} z_\ell$  which sums over the neurons in layer  $\ell$ . By chain rule, the gradient of error  $e$  w.r.t. to the weight between neurons  $\ell$  and  $i$  is:

$$\frac{\partial e}{\partial x_{i\ell}} = \frac{\partial e}{\partial a_i} \times \frac{\partial a_i}{\partial x_{i\ell}} \stackrel{(a)}{=} \delta_i \times z_\ell, \quad (22)$$

where (a) is because  $a_i = \sum_\ell x_{i\ell} z_\ell$  and we define:

$$\delta_i := \frac{\partial e}{\partial a_i}.$$

# Backpropagation



- If layer  $i$  is the last layer,  $\delta_i$  can be computed by derivative of error (loss function) w.r.t. the output.
- However, if  $i$  is one of the hidden layers,  $\delta_i$  is computed by chain rule as:

$$\delta_i = \frac{\partial e}{\partial a_i} = \sum_j \left( \frac{\partial e}{\partial a_j} \times \frac{\partial a_j}{\partial a_i} \right) = \sum_j \left( \delta_j \times \frac{\partial a_j}{\partial a_i} \right). \quad (23)$$

- The term  $\partial a_j / \partial a_i$  is calculated by chain rule as:

$$\frac{\partial a_j}{\partial a_i} = \frac{\partial a_j}{\partial z_i} \times \frac{\partial z_i(a)}{\partial a_i} = x_{ji} \sigma'(a_i), \quad (24)$$

where  $(a)$  is because  $a_j = \sum_i x_{ji} z_i$  and  $z_i = \sigma(a_i)$  and  $\sigma'(\cdot)$  denotes the derivative of activation function. Putting Eq. (24) in Eq. (23) gives:

$$\delta_i = \sigma'(a_i) \sum_j (\delta_j x_{ji}).$$

# Backpropagation

- We found:

$$\delta_i = \sigma'(a_i) \sum_j (\delta_j x_{ji}).$$

- Putting this equation in Eq. (22),  $\frac{\partial e}{\partial x_{i\ell}} = \delta_i \times z_\ell$ , gives:

$$\frac{\partial e}{\partial x_{i\ell}} = z_\ell \sigma'(a_i) \sum_j (\delta_j x_{ji}). \quad (25)$$

- **Backpropagation** uses the gradient in Eq. (25) for updating the weight  $x_{i\ell}, \forall i, \ell$  by gradient descent:

$$x_{i\ell}^{(k+1)} := x_{i\ell}^{(k)} - \eta^{(k)} \frac{\partial e}{\partial x_{i\ell}}.$$

- This tunes the weights from last layer to the first layer for every iteration of optimization.
- Therefore, **backpropagation**, proposed in 1986 [7], is actually gradient descent with chain rule in derivatives because of having layers of parameters. It is the most well-known optimization method used for training neural networks.

## Accelerated gradient method

# Accelerated Gradient method

- It was shown in the literature that gradient descent is not optimal in convergence rate and can be improved.
- It was at that time that Nesterov proposed **Accelerated Gradient Method (AGM)** [10], in 1983, to make the convergence rate of gradient descent optimal [11, Chapter 2.2].
- AGM is also called the **Nesterov's accelerated gradient method** or **Fast Gradient Method (FGM)**. A series of Nesterov's papers improved AGM [10, 12, 13, 14].
- Consider a sequence  $\{\gamma^{(k)}\}$  which satisfies:

$$\prod_{i=0}^k (1 - \gamma^{(i)}) \geq (\gamma^{(k)})^2, \quad \forall k \geq 0, \quad \gamma^{(k)} \in [0, 1]. \quad (26)$$

- An example sequence, satisfying this condition, is

$$\gamma^{(0)} = \gamma^{(1)} = \gamma^{(2)} = \gamma^{(3)} = 0, \gamma^{(k)} = 2/k, \forall k \geq 4.$$

- The AGM updates the solution iteratively as [10]:

$$\mathbf{x}^{(k+1)} := \mathbf{y}^{(k)} - \eta^{(k)} \nabla f(\mathbf{y}^{(k)}), \quad (27)$$

$$\mathbf{y}^{(k+1)} := (1 - \gamma^{(k)}) \mathbf{x}^{(k+1)} + \gamma^{(k)} \mathbf{x}^{(k)}, \quad (28)$$

until convergence.

# Comparison of convergences rates

- Consider a convex and differentiable function  $f(\cdot)$ , with domain  $\mathcal{D}$ , whose gradient is  $L$ -smooth (see Eq. (2)). Let  $f^*$  be the minimum of cost function and  $\mathbf{x}^*$  be the minimizer. Starting from the initial point  $\mathbf{x}^{(0)}$ , after  $t$  iterations of the optimization algorithm, we will have the following.
- The convergence rate of gradient descent:

$$f(\mathbf{x}^{(t+1)}) - f^* \leq \frac{2L\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2}{t+1} = \mathcal{O}\left(\frac{1}{t}\right). \quad (29)$$

- The convergence rate of accelerated gradient method:

$$f(\mathbf{x}^{(t+1)}) - f^* \leq \frac{2L\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2}{(t+1)^2} = \mathcal{O}\left(\frac{1}{t^2}\right). \quad (30)$$

## Stochastic gradient methods

# Stochastic gradient descent

- Assume we have a dataset of  $n$  data points,  $\{\mathbf{a}_i \in \mathbb{R}^d\}_{i=1}^n$  and their labels  $\{l_i \in \mathbb{R}\}_{i=1}^n$ .
- Let the cost function  $f(\cdot)$  be decomposed into summation of  $n$  terms  $\{f_i(\mathbf{x})\}_{i=1}^n$ . Some well-known examples for the cost function terms are:
  - ▶ Least squares error:  $f_i(\mathbf{x}) = 0.5(\mathbf{a}_i^\top \mathbf{x} - l_i)^2$ ,
  - ▶ Absolute error:  $f_i(\mathbf{x}) = \mathbf{a}_i^\top \mathbf{x} - l_i$ ,
  - ▶ Hinge loss (for  $l_i \in \{-1, 1\}$ ):  $f_i(\mathbf{x}) = \max(0, 1 - l_i \mathbf{a}_i^\top \mathbf{x})$ .
  - ▶ Logistic loss (for  $l_i \in \{-1, 1\}$ ):  $\log\left(\frac{1}{1 + \exp(-l_i \mathbf{a}_i^\top \mathbf{x})}\right)$ .
- The optimization problem becomes:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}). \quad (31)$$

In this case, the full gradient is the average gradient, i.e:

$$\nabla f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}), \quad (32)$$

so Eq. (9),  $\Delta \mathbf{x} = -(1/L)\nabla f(\mathbf{x}^{(k)})$ , becomes  $\Delta \mathbf{x} = -(1/(Ln)) \sum_{i=1}^n \nabla f_i(\mathbf{x}^{(k)})$ . This is what gradient descent uses for updating the solution at every iteration.

# Stochastic gradient descent

$$\nabla f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}),$$

- Calculation of this full gradient is time-consuming and inefficient for large values of  $n$ , especially as it needs to be recalculated at every iteration.
- **Stochastic Gradient Descent (SGD)**, also called **stochastic gradient method**, approximates gradient descent stochastically and samples (i.e. bootstraps) one of the points at every iteration for updating the solution. Hence, it uses:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \eta^{(k)} \nabla f_i(\mathbf{x}^{(k)}), \quad (33)$$

rather than Eq. (10),  $\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \eta \nabla f(\mathbf{x}^{(k)})$ .

- The idea of stochastic approximation was first proposed in 1951 [15]. It was first used for machine learning in 1998 [16].
- As Eq. (33) states, SGD often uses an adaptive step size which changes in every iteration. The step size can be decreasing because in initial iterations, where we are far away from the optimal solution, the step size can be large; however, it should be small in the last iterations which is supposed to be close to the optimal solution. Some well-known adaptations for the step size are:

$$\eta^{(k)} := \frac{1}{k}, \quad \eta^{(k)} := \frac{1}{\sqrt{k}}, \quad \eta^{(k)} := \eta. \quad (34)$$

# Stochastic gradient descent

## Theorem (Convergence rates for SGD)

Consider a function  $f(\mathbf{x}) = \sum_{i=1}^n f_i(\mathbf{x})$  and which is bounded below and each  $f_i$  is differentiable. Let the domain of function  $f(\cdot)$  be  $\mathcal{D}$  and its gradient be  $L$ -smooth (see Eq. (2)). Assume  $\mathbb{E}[\|\nabla f_i(\mathbf{x}_k)\|_2^2 | \mathbf{x}_k] \leq \beta^2$  where  $\beta$  is a constant. Depending on the step size, the convergence rate of SGD is:

$$\mathcal{O}\left(\frac{1}{\log t}\right) \quad \text{if } \eta^{(k)} = \frac{1}{k}, \quad (35)$$

$$\mathcal{O}\left(\frac{\log t}{\sqrt{t}}\right) \quad \text{if } \eta^{(k)} = \frac{1}{\sqrt{k}}, \quad (36)$$

$$\mathcal{O}\left(\frac{1}{t} + \eta\right) \quad \text{if } \eta^{(k)} = \eta, \quad (37)$$

where  $t$  denotes the iteration index. If the functions  $f_i$ 's are  $\mu$ -strongly convex, then the convergence rate of SGD is:

$$\mathcal{O}\left(\frac{1}{t}\right) \quad \text{if } \eta^{(k)} = \frac{1}{\mu k}, \quad (38)$$

$$\mathcal{O}\left(\left(1 - \frac{\mu}{L}\right)^t + \eta\right) \quad \text{if } \eta^{(k)} = \eta. \quad (39)$$

# Stochastic gradient descent

- Recall Eqs. (37) and (39):

$$\text{convex: } \mathcal{O}\left(\frac{1}{t} + \eta\right) \quad \text{if } \eta^{(k)} = \eta,$$

$$\text{strongly convex: } \mathcal{O}\left(\left(1 - \frac{\mu}{L}\right)^t + \eta\right) \quad \text{if } \eta^{(k)} = \eta.$$

- These equations show that with a fixed step size  $\eta$ , SGD converges sublinearly for a non-convex function and linearly for a strongly convex function in the initial iterations.
- However, in the late iterations, it stagnates to a neighborhood around the optimal point and never reaches it. Hence, SGD has less accuracy than gradient descent (whose convergence rate is  $\mathcal{O}\left(\frac{1}{t^2}\right)$ ).
- The advantage of SGD over gradient descent is that its every iteration is much faster than every iteration of gradient descent because of less computations for gradient. This faster pacing of every iteration shows off more when  $n$  is huge.
- In summary, SGD has fast convergence to low accurate optimal point.
- It is noteworthy that the full gradient is not available in SGD to use for checking convergence, as discussed before. One can use other criteria or merely check the norm of gradient for the sampled point.
- SGD can be used with the line-search methods, too. SGD can also use a momentum term.

# Mini-batch stochastic gradient descent

- Gradient descent uses the entire  $n$  data points and SGD uses one randomly sampled point at every iteration. For large datasets, gradient descent is very slow and intractable in every iteration while SGD will need a significant number of iterations to roughly cover all data. Besides, SGD has low accuracy in convergence to the optimal point.
- We can have a middle case where we use a batch of  $b$  randomly sampled points at every iteration. This method is named the **mini-batch SGD** or the **hybrid deterministic-stochastic gradient** method. This batch-wise approach is wise for large datasets.
- Usually, before start of optimization, the  $n$  data points are randomly divided into  $\lfloor n/b \rfloor$  batches of size  $b$ . This is equivalent to simple random sampling for sampling points into batches without replacement. We denote the dataset by  $\mathcal{D}$  (where  $|\mathcal{D}| = n$ ) and the  $i$ -th batch by  $\mathcal{B}_i$  (where  $|\mathcal{B}_i| = b$ ). The batches are disjoint:

$$\bigcup_{i=1}^{\lfloor n/b \rfloor} \mathcal{B}_i = \mathcal{D}, \quad (40)$$

$$\mathcal{B}_i \cap \mathcal{B}_j = \emptyset, \quad \forall i, j \in \{1, \dots, \lfloor n/b \rfloor\}, i \neq j. \quad (41)$$

- Another less-used approach for making batches is to sample points for a batch during optimization. This is equivalent to bootstrapping for sampling points into batches with replacement. In this case, the batches are not disjoint anymore and Eqs. (40) and (41) do not hold.

# Mini-batch stochastic gradient descent

## Definition (Epoch)

In mini-batch SGD, when all  $\lfloor n/b \rfloor$  batches of data are used for optimization once, an **epoch** is completed. After completion of an epoch, the next epoch is started and epochs are repeated until convergence of optimization.

- In mini-batch SGD, if the  $k$ -th iteration of optimization is using the  $k'$ -th batch, the update of solution is done as:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \eta^{(k)} \frac{1}{b} \sum_{i \in B_{k'}} \nabla f_i(\mathbf{x}^{(k)}). \quad (42)$$

- The scale factor  $1/b$  is sometimes dropped for simplicity.
- Mini-batch SGD is used significantly in machine learning, especially in neural networks [16, 17].
- Because of dividing data into batches, mini-batch SGD can be solved on parallel servers as a distributed optimization method.

# Mini-batch stochastic gradient descent

## Theorem (Convergence rates for mini-batch SGD)

Consider a function  $f(\mathbf{x}) = \sum_{i=1}^n f_i(\mathbf{x})$  which is bounded below and each  $f_i$  is differentiable. Let the domain of function  $f(\cdot)$  be  $\mathcal{D}$  and its gradient be  $L$ -smooth (see Eq. (2)) and assume  $\eta^{(k)} = \eta$  is fixed. The batch-wise gradient is an approximation to the full gradient with some error  $\mathbf{e}_t$  for the  $t$ -th iteration:

$$\frac{1}{b} \sum_{i \in \mathcal{B}_t} \nabla f_i(\mathbf{x}^{(t)}) = \nabla f(\mathbf{x}^{(t)}) + \mathbf{e}_t. \quad (43)$$

The convergence rate of mini-batch SGD for non-convex and convex functions are:

$$\mathcal{O}\left(\frac{1}{t} + \|\mathbf{e}_t\|_2^2\right), \quad (44)$$

where  $t$  denotes the iteration index. If the functions  $f_i$ 's are  $\mu$ -strongly convex, then the convergence rate of mini-batch SGD is:

$$\mathcal{O}\left(\left(1 - \frac{\mu}{L}\right)^t + \|\mathbf{e}_t\|_2^2\right). \quad (45)$$

Therefore, the convergence rate of mini-batch gets closer to that of gradient descent,  $\mathcal{O}(1/t)$ , if the batch size increases.

# Mini-batch stochastic gradient descent

- If we sample the batches without replacement (i.e., sampling batches by simple random sampling before start of optimization) or with replacement (i.e., bootstrapping during optimization), the expected error is [18, Proposition 3]:

$$\mathbb{E}[\|e_t\|_2^2] = \left(1 - \frac{b}{n}\right) \frac{\sigma^2}{b}, \quad (46)$$

$$\mathbb{E}[\|e_t\|_2^2] = \frac{\sigma^2}{b}, \quad (47)$$

respectively, where  $\sigma^2$  is the variance of whole dataset.

- According to Eqs. (46) and (47), the accuracy of SGD by sampling without and with replacement increases by  $b \rightarrow n$  and  $b \rightarrow \infty$ , respectively.
- However, this increase makes every iteration slower so there is trade-off between accuracy and speed.

## Stochastic Average Gradient Methods

# Stochastic average gradient

- SGD is faster than gradient descent but its problem is its lower accuracy compared to gradient descent. **Stochastic Average Gradient (SAG)**, proposed in 2012 [19], keeps a trade-off between accuracy and speed.
- Let  $\nabla f_i(\mathbf{x}^{(k)})$  be the gradient of  $f_i(\cdot)$ , evaluated in point  $\mathbf{x}^{(k)}$ , at iteration  $k$ . According to Eqs. (10) and (32), gradient descent updates the solution as:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \frac{\eta^{(k)}}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}^{(k)}).$$

SAG randomly samples one of the points and updates its gradient among the gradient terms. If the sampled point is the  $j$ -th one, we have:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \frac{\eta^{(k)}}{n} \left( \nabla f_j(\mathbf{x}^{(k)}) - \nabla f_j(\mathbf{x}^{(k-1)}) + \sum_{i=1}^n \nabla f_i(\mathbf{x}^{(k-1)}) \right). \quad (48)$$

- In other words, we subtract the  $j$ -th gradient from the summation of all  $n$  gradients in previous iteration ( $k-1$ ) by  $\sum_{i=1}^n \nabla f_i(\mathbf{x}^{(k-1)}) - \nabla f_j(\mathbf{x}^{(k-1)})$ ; then, we add back the new  $j$ -th gradient in this iteration by adding  $\nabla f_j(\mathbf{x}^{(k)})$ .

# Stochastic average gradient

## Theorem (Convergence rates for SAG [19, Proposition 1])

Consider a function  $f(\mathbf{x}) = \sum_{i=1}^n f_i(\mathbf{x})$  which is bounded below and each  $f_i$  is differentiable. Let the domain of function  $f(\cdot)$  be  $\mathcal{D}$  and its gradient be  $L$ -smooth (see Eq. (2)). The convergence rate of SAG is:

$$\mathcal{O}\left(\frac{1}{t}\right),$$

where  $t$  denotes the iteration index.

- SAG has the same rate order as gradient descent; although, it usually needs some more iterations to converge.
- Practical experiments have shown that SAG requires many parameter fine-tuning to perform perfectly.

# Stochastic variance reduced gradient

- Another effective first-order method is the **Stochastic Variance Reduced Gradient (SVRG)**, proposed in 2013 [20].

```
1 Initialize  $\tilde{\mathbf{x}}^{(0)}$ 
2 for iteration  $k = 1, 2, \dots$  do
3    $\tilde{\mathbf{x}} := \tilde{\mathbf{x}}^{(k-1)}$ 
4    $\nabla f(\tilde{\mathbf{x}}) \stackrel{(112)}{:=} \frac{1}{n} \sum_{i=1}^n \nabla f_i(\tilde{\mathbf{x}})$ 
5    $\mathbf{x}^{(0)} := \tilde{\mathbf{x}}$ 
6   for iteration  $\tau = 0, 1, \dots, m-1$  do
7     Randomly sample  $j$  from  $\{1, \dots, n\}$ .
8      $\mathbf{x}^{(\tau+1)} := \mathbf{x}^{(\tau)} - \eta^{(\tau)} (\nabla f_j(\mathbf{x}^{(\tau)}) -$ 
9        $\nabla f_j(\tilde{\mathbf{x}}) + \nabla f(\tilde{\mathbf{x}}))$ .
9    $\tilde{\mathbf{x}}^{(k)} := \mathbf{x}^{(m)}$ 
```

- The update of solution is similar to SAG but for every iteration, it updates the solution for  $m$  times.
- SVRG is an efficient method and its convergence rate is similar to that of SAG.
- Both SAG and SVRG reduce the variance of solution to optimization [20].

## Adaptive Learning Rate

# Adaptive Gradient (AdaGrad)

- We can adapt the learning rate in stochastic gradient methods. Three most well-known methods for adapting the learning rate are AdaGrad, RMSProp, and Adam.
- **Adaptive Gradient (AdaGrad)** method, proposed in 2011 [21], updates the solution iteratively as:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \eta^{(k)} \mathbf{G}^{-1} \nabla f_i(\mathbf{x}^{(k)}), \quad (49)$$

where  $\mathbf{G}$  is a  $(d \times d)$  diagonal matrix whose  $(j, j)$ -th element is:

$$\mathbf{G}(j, j) := \sqrt{\varepsilon + \sum_{\tau=0}^k (\nabla_j f_{i_\tau}(\mathbf{x}^{(\tau)}))^2}, \quad (50)$$

where  $\varepsilon \geq 0$  is for stability (making  $\mathbf{G}$  full rank),  $i_\tau$  is the randomly sampled point (from  $\{1, \dots, n\}$ ) at iteration  $\tau$ , and  $\nabla_j f_{i_\tau}(\cdot)$  is the partial derivative of  $f_{i_\tau}(\cdot)$  w.r.t. its  $j$ -th element (note that  $f_{i_\tau}(\cdot)$  is  $d$ -dimensional).

- Putting Eq. (50) in Eq. (49) can simplify AdaGrad to:

$$\mathbf{x}_j^{(k+1)} := \mathbf{x}_j^{(k)} - \frac{\eta^{(k)}}{\sqrt{\varepsilon + \sum_{\tau=0}^k (\nabla_j f_{i_\tau}(\mathbf{x}^{(\tau)}))^2}} \nabla f_j(\mathbf{x}_j^{(k)}). \quad (51)$$

- AdaGrad keeps a history of the sampled points and it takes derivative for them to use. During the iterations so far, if a dimension has changed significantly, it dampens the learning rate for that dimension (see the inverse in Eq. (49)); hence, it gives more weight for changing the dimensions which have not changed noticeably.

# Root Mean Square Propagation (RMSProp)

- **Root Mean Square Propagation (RMSProp)** was first proposed in 2012 [22] which is unpublished.
- It is an improved version of **Rprop (resilient backpropagation)**, proposed in 1992 [23], which uses the sign of gradient in optimization.
- Inspired by momentum in Eq. (19):

$$(\Delta \mathbf{x})^{(k)} := \alpha(\Delta \mathbf{x})^{(k-1)} - \eta^{(k)} \nabla f(\mathbf{x}^{(k)}),$$

it updates a scalar variable  $v$  as [24]:

$$v^{(k+1)} := \gamma v^{(k)} + (1 - \gamma) \|\nabla f_j(\mathbf{x}^{(k)})\|_2^2, \quad (52)$$

where  $\gamma \in [0, 1]$  is the forgetting factor (e.g.,  $\gamma = 0.9$ ). Then, it uses this  $v$  to weight the learning rate:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \frac{\eta^{(k)}}{\sqrt{\epsilon + v^{(k+1)}}} \nabla f_j(\mathbf{x}_j^{(k)}), \quad (53)$$

where  $\epsilon \geq 0$  is for stability not to have division by zero.

- Comparing Eqs. (51) and (53) shows that RMSProp has a similar form to AdaGrad.

# Adaptive Moment Estimation (Adam)

- **Adam optimizer** [25] improves over RMSProp by adding a momentum term.
- It updates the scalar  $v$  and the vector  $\mathbf{m} \in \mathbb{R}^d$  as:

$$\mathbf{m}^{(k+1)} := \gamma_1 \mathbf{m}^{(k)} + (1 - \gamma_1) \nabla f_i(\mathbf{x}^{(k)}), \quad (54)$$

$$v^{(k+1)} := \gamma_2 v^{(k)} + (1 - \gamma_2) \|\nabla f_i(\mathbf{x}^{(k)})\|_2^2, \quad (55)$$

where  $\gamma_1, \gamma_2 \in [0, 1]$ . It normalizes these variables as:

$$\hat{\mathbf{m}}^{(k+1)} := \frac{1}{1 - \gamma_1^k} \mathbf{m}^{(k+1)}, \quad \hat{v}^{(k+1)} := \frac{1}{1 - \gamma_2^k} v^{(k+1)}.$$

- Then, it updates the solution as:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \frac{\eta^{(k)}}{\sqrt{\varepsilon + \hat{v}^{(k+1)}}} \hat{\mathbf{m}}^{(k+1)}, \quad (56)$$

which is stochastic gradient descent with momentum while using RMSProp.

- The Adam optimizer is one of the mostly used optimizers in **neural networks**.

## Coding a Neural Network

# Neural network: importing packages

## ▼ Importing packages

```
✓ [188] # installation in Google Colab's Jupyter notebook:
```

```
!pip install torch
```

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
```

```
Requirement already satisfied: torch in /usr/local/lib/python3.8/dist-packages (1.13.1+cu116)
```

```
Requirement already satisfied: typing-extensions in /usr/local/lib/python3.8/dist-packages (from torch) (4.4.0)
```

```
[320] # importing packages (libraries):
```

```
import torch
```

```
import torch.nn as nn
```

```
from torch.utils.data import DataLoader, Dataset
```

```
import matplotlib.pyplot as plt
```

```
import numpy as np
```

```
from tqdm import tqdm
```

# Neural network: defining the network

## ▾ Defining the network

```
✓ [216] device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")  
0s
```

```
✓ [339] # defining the structure of neural network:  
0s  
class NeuralNetwork(nn.Module):  
    def __init__(self):  
        super(NeuralNetwork, self).__init__()  
        self.layer1 = nn.Linear(1, 10)  
        self.layer2 = nn.Linear(10, 20)  
        self.layer3 = nn.Linear(20, 10)  
        self.layer4 = nn.Linear(10, 1)  
        self.relu1 = nn.ReLU()  
        self.relu2 = nn.ReLU()  
        self.relu3 = nn.ReLU()  
  
    def forward(self, x):  
        x = self.relu1(self.layer1(x))  
        x = self.relu2(self.layer2(x))  
        x = self.relu3(self.layer3(x))  
        x = self.layer4(x)  
        return x
```

```
✓ [340] # instantiate the class of neural network:  
0s
```

```
net = NeuralNetwork()  
print(net)  
  
NeuralNetwork(  
  (layer1): Linear(in_features=1, out_features=10, bias=True)  
  (layer2): Linear(in_features=10, out_features=20, bias=True)  
  (layer3): Linear(in_features=20, out_features=10, bias=True)  
  (layer4): Linear(in_features=10, out_features=1, bias=True)  
  (relu1): ReLU()  
  (relu2): ReLU()  
  (relu3): ReLU()  
)
```

# Neural network: optimizer

## ▼ Optimizer

```
✓ [341] # define optimizer:
0s optimizer = 'Adam'
    if optimizer == 'SGD':
        optimizer = torch.optim.SGD(net.parameters(), lr=0.02)
    elif optimizer == 'Adam':
        optimizer = torch.optim.Adam(net.parameters(), lr=0.02)

# define the loss function:
loss_func = torch.nn.MSELoss()
```

# Neural network: data loader

## ▼ Data loader

```
✓ [342] class Data(Dataset):  
05     def __init__(self, x, y):  
        self.data = torch.from_numpy(x.reshape(-1,1)).float()  
        self.label = torch.from_numpy(y.reshape(-1,1)).float()  
  
        def __len__(self):  
            return len(self.data)  
  
        def __getitem__(self, item):  
            data_point = self.data[item]  
            label_point = self.label[item]  
            return data_point, label_point
```

```
✓ [343] batch_size = 16  
05     def load_dataset(x_train, y_train, x_test, y_test):  
        # data loader for training data:  
        train_ds = Data(x_train, y_train)  
        train_loader = DataLoader(train_ds, batch_size=batch_size, shuffle=True)  
  
        # data loader for test data:  
        test_ds = Data(x_test, y_test)  
        test_loader = DataLoader(test_ds, batch_size=batch_size, shuffle=False)  
  
        return train_loader, test_loader
```

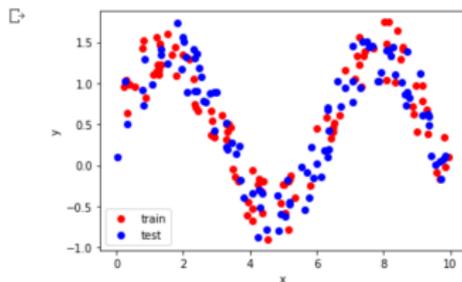
# Neural network: dataset

```
dataset_type = 'nonlinear'

if dataset_type == 'linear':
    # almost linear dataset:
    x_train = np.random.rand(100)
    y_train = np.sin(x_train) * (x_train**3) + 3*x_train + np.random.rand(100)*0.8
    x_test = np.random.rand(100)
    y_test = np.sin(x_test) * (x_test**3) + 3*x_test + np.random.rand(100)*0.8
elif dataset_type == 'nonlinear':
    # dataset:
    x_train = np.random.rand(100) * 10
    y_train = np.sin(x_train) + np.random.rand(100)*0.8
    x_test = np.random.rand(100) * 10
    y_test = np.sin(x_test) + np.random.rand(100)*0.8

# reshape to have samples in rows:
x_train = x_train.reshape((-1, 1))
x_test = x_test.reshape((-1, 1))

# visualize data:
plt.scatter(x_train, y_train, c='r', label='train')
plt.scatter(x_test, y_test, c='b', label='test')
plt.xlabel('x')
plt.ylabel('y')
plt.legend()
plt.show()
```



# Neural network: training

## ▼ Training neural network

```
✓ [345] # load dataset:  
0s train_loader, test_loader = load_dataset(x_train, y_train, x_test, y_test)
```

```
✓ [346] n_epochs = 1000  
12s loss_list = []  
for epoch in tqdm(range(n_epochs), desc='epochs'):  
    loss_list_in_epoch = []  
    for step, (data_batch, label_batch) in enumerate(train_loader):  
        data_batch, label_batch = data_batch.to(device), label_batch.to(device)  
        prediction = net(data_batch)  
        loss = loss_func(prediction, label_batch)  
        loss_list_in_epoch.append(loss.cpu().detach().item())  
        optimizer.zero_grad()  
        loss.backward()  
        optimizer.step()  
    loss_list.append(np.mean(loss_list_in_epoch))
```

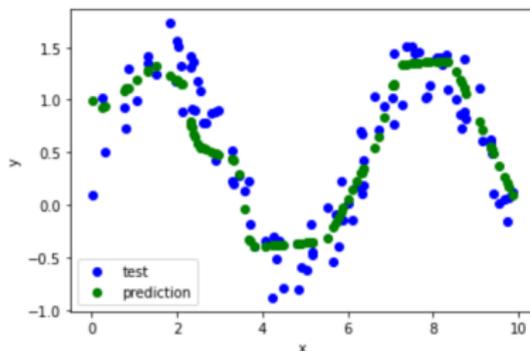
```
epochs: 100%|██████████| 1000/1000 [00:12<00:00, 81.84it/s]
```

# Neural network: test (evaluation)

## ▼ Test (evaluation) phase

```
✓ 0s ▶ prediction_list = []  
with torch.no_grad():  
    for step, (data_batch, label_batch) in enumerate(test_loader):  
        prediction = net(data_batch)  
        prediction_list.extend(prediction)
```

```
✓ [350] # visualize the predicted and actual data:  
plt.scatter(x_test, y_test, c='b', label='test')  
plt.scatter(x_test, prediction_list, c='g', label='prediction')  
plt.xlabel('x')  
plt.ylabel('y')  
plt.legend()  
plt.show()
```



## Proximal Methods

# Proximal Mapping



## Definition (Proximal mapping/operator [26])

The proximal mapping or proximal operator of a convex function  $g(\cdot)$  is:

$$\lambda=1 \rightarrow \text{prox}_{g(\cdot)}(x) = \arg \min_u \left( g(u) + \frac{1}{2} \|u-x\|_2^2 \right). \quad (57)$$

In case the function  $g(\cdot)$  is scaled by a scalar  $\lambda$ , the proximal mapping is defined as:

$$\star \text{prox}_{\lambda g(\cdot)}(x) := \arg \min_u \left( g(u) + \frac{1}{2\lambda} \|u-x\|_2^2 \right). \quad (58)$$

$\lambda(u-x)$

# The Moreau-Yosida regularization

- The proximal mapping:

$$\star \quad \boxed{\text{prox}_{\lambda g}(\mathbf{x}) := \underset{\mathbf{u}}{\text{argmin}} \left( g(\mathbf{u}) + \frac{1}{2\lambda} \|\mathbf{u} - \mathbf{x}\|_2^2 \right),}$$

is related to the Moreau-Yosida regularization defined below.

Definition (Moreau-Yosida regularization or Moreau envelope [27, 28])

The Moreau-Yosida regularization or the Moreau envelope of function  $g(\cdot)$  is:

$$\boxed{M_{\lambda g}(\mathbf{x}) := \inf_{\mathbf{u}} \left( g(\mathbf{u}) + \frac{1}{2\lambda} \|\mathbf{u} - \mathbf{x}\|_2^2 \right).} \quad (59)$$

This Moreau-Yosida regularized function has the same minimizer as the function  $g(\cdot)$  [29].

# The Moreau decomposition

## Lemma (Moreau decomposition [30])

We always have the following decomposition, named the Moreau decomposition:

$$\lambda \geq 1 \rightarrow \boxed{x = \text{prox}_g(x) + \text{prox}_{g^*}(x)}, \quad (60)$$

$$\boxed{x = \text{prox}_{\lambda g}(x) + \lambda \text{prox}_{\frac{1}{\lambda} g^*}\left(\frac{x}{\lambda}\right)}, \quad (61)$$

where  $g(\cdot)$  is a function in a space and  $g^*(\cdot)$  is its corresponding function in the dual space (e.g., if  $g(\cdot)$  is a norm,  $g^*(\cdot)$  is its dual norm or if  $g(\cdot)$  is projection onto a cone,  $g^*(\cdot)$  is projection onto the dual cone).

## Projection onto set

- In optimization, **indicator function**  $\mathbb{I}(\cdot)$  is zero if its condition is satisfied and is infinite otherwise.

$$\mathbb{I}(x \in S) = \begin{cases} 0 & \text{if } x \in S \\ \infty & \text{if } x \notin S \end{cases} \quad (62)$$

~~$\mathbb{I}(x \in S) = \begin{cases} 1 & \text{if } x \in S \\ 0 & \text{if } x \notin S \end{cases}$~~

$\mathbb{I}(x \in S) = \begin{cases} 0 & \text{if } x \in S \\ \infty & \text{if } x \notin S \end{cases}$

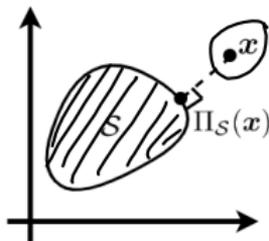
$\min f \rightarrow f = \infty$

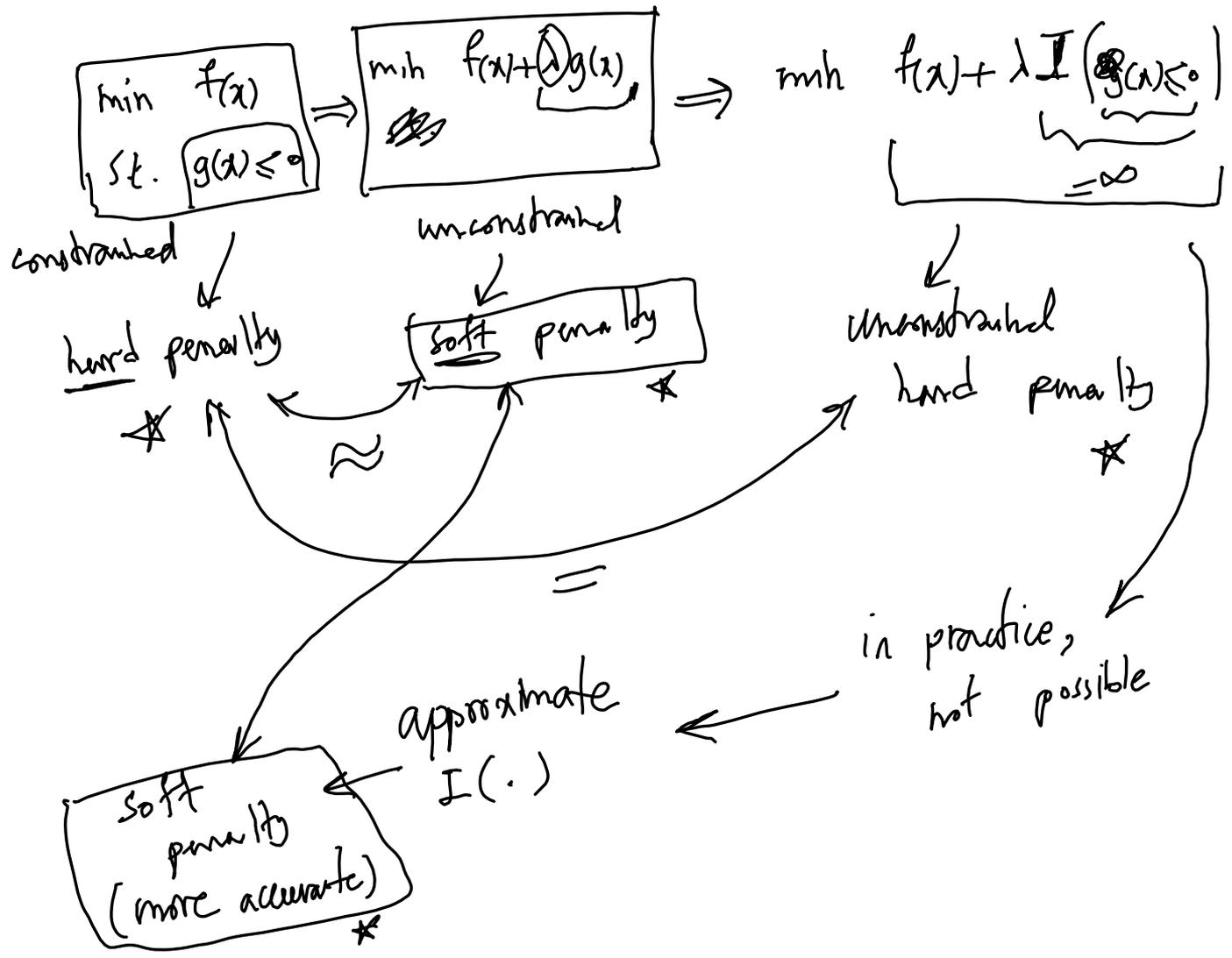
## Lemma (Projection onto set)

The proximal mapping of the indicator function to a convex set  $S$ , i.e.  $\mathbb{I}(x \in S)$ , is projection of the point  $x$  onto the set  $S$ . Hence, projection of  $x$  onto set  $S$ , denoted by  $\Pi_S(x)$ , is defined as:

$$\Pi_S(x) := \text{prox}_{\mathbb{I}(\cdot \in S)}(x) = \arg \min_{u \in S} \frac{1}{2} \|u - x\|_2^2 \quad (63)$$

- This projection simply means projecting the point  $x$  onto the closest point of set from the point  $x$ . Hence, the vector connecting the points  $x$  and  $\Pi_S(x)$  is orthogonal to the set  $S$ .





# Projection onto set

Lemma:

$$\star \quad \Pi_S(x) := \text{prox}_{\mathbb{I}(\cdot \in S)}(x) \stackrel{\text{def}}{=} \arg \min_{u \in S} \left( \frac{1}{2} \|u - x\|_2^2 \right).$$

The proximal mapping:

$$\star \quad \text{prox}_g(x) := \arg \min_u \left( g(u) + \frac{1}{2} \|u - x\|_2^2 \right),$$

Proof.

$$\star \quad \text{prox}_{\mathbb{I}(\cdot \in S)}(x) \stackrel{(57)}{=} \arg \min_u \left( \mathbb{I}(x \in S) + \frac{1}{2} \|u - x\|_2^2 \right)$$

*if  $x \in S$   
if  $x \notin S$*

$$\stackrel{(a)}{=} \arg \min_{u \in S} \left( \frac{1}{2} \|u - x\|_2^2 \right),$$

where (a) is because  $\mathbb{I}(x \in S)$  becomes infinity if  $x \notin S$  (see Eq. (62)). □

# Moreau decomposition for norm

Recall Eq. (61) from the Moreau decomposition:

$$x = \text{prox}_{\lambda g}(x) + \lambda \text{prox}_{\frac{1}{\lambda} g^*}\left(\frac{x}{\lambda}\right) \Rightarrow \text{prox}_{\lambda g}(x) = x - \lambda \text{prox}_{\frac{1}{\lambda} g^*}\left(\frac{x}{\lambda}\right).$$

$\text{prox}_{\frac{1}{\lambda} g^*}\left(\frac{x}{\lambda}\right) = \frac{\text{prox}_{\|\cdot\|}(\frac{x}{\lambda})}{\lambda}$   
 $\text{prox}_{\|\cdot\|}(\frac{x}{\lambda}) = \frac{\lambda}{\lambda} \text{prox}_{\|\cdot\|}(\frac{x}{\lambda})$   
 $\star \star$

## Corollary (Moreau decomposition for norm)

If the function is a scaled norm,  $g(\cdot) = \lambda \|\cdot\|$ , we have from this equation:

$$\star \text{prox}_{\lambda \|\cdot\|}(x) = x - \lambda \Pi_{\mathcal{B}}\left(\frac{x}{\lambda}\right), \quad (64)$$

$g^* = \lambda \|\cdot\|$

where  $\mathcal{B}$  is the unit ball of dual norm.

- Derivation of proximal operator for various  $g(\cdot)$  functions are available in [31, Chapter 6]. Here, we review the proximal mapping of some mostly used functions.
- If  $g(x) = 0$ , proximal mapping becomes an identity mapping:

$$\text{prox}_{\lambda 0}(x) \stackrel{(58)}{=} \arg \min_u \left( 0 + \frac{1}{2\lambda} \|u - x\|_2^2 \right) = \arg \min_u \left( \frac{1}{2\lambda} \|u - x\|_2^2 \right) = x$$

$\|u - x\|_2^2 = 0 \rightarrow u = x$

# Proximal mapping of $\ell_2$ norm

## Lemma (Proximal mapping of $\ell_2$ norm [31, Example 6.19])

The proximal mapping of the  $\ell_2$  norm is:

$$\text{prox}_{\lambda \|\cdot\|_2}(x) = \begin{cases} \left(1 - \frac{\lambda}{\|x\|_2}\right)x & \text{if } \|x\|_2 \geq \lambda \\ 0 & \text{if } \|x\|_2 < \lambda \end{cases} \quad (65)$$

### Proof.

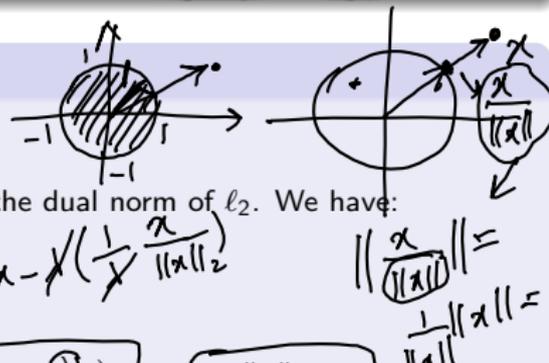
Recall Eq. (64):  $\text{prox}_{\lambda \|\cdot\|}(x) = x - \lambda \Pi_{\mathcal{B}}\left(\frac{x}{\lambda}\right)$ .

Let  $g(\cdot) = \|\cdot\|_2$  and  $\mathcal{B}$  be the unit  $\ell_2$  ball because  $\ell_2$  is the dual norm of  $\ell_2$ . We have:

$$\star \Pi_{\mathcal{B}}\left(\frac{x}{\lambda}\right) = \begin{cases} \frac{x/\lambda}{\|x/\lambda\|_2} & \text{if } \|x/\lambda\|_2 \geq 1 \\ \frac{x}{\lambda} & \text{if } \|x/\lambda\|_2 < 1 \end{cases}$$

$$\Rightarrow \text{prox}_{\lambda \|\cdot\|_2}(x) \stackrel{(64)}{=} x - \lambda \Pi_{\mathcal{B}}\left(\frac{x}{\lambda}\right) = \begin{cases} \left(1 - \frac{\lambda}{\|x\|_2}\right)x & \text{if } \|x\|_2 \geq \lambda \\ x - \lambda \left(\frac{x}{\lambda}\right) = 0 & \text{if } \|x\|_2 < \lambda \end{cases}$$

$\frac{x}{\lambda} \geq 1 \rightarrow x \geq \lambda$   $x - \lambda \frac{x}{\lambda} = x - x = 0$   $\square$



# Proximal mapping of $\ell_1$ norm

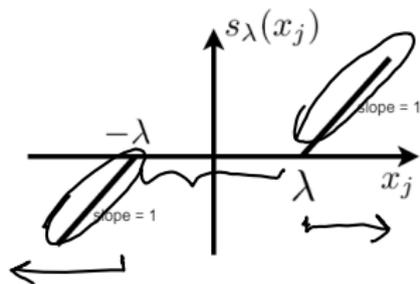
## Lemma (Proximal mapping of $\ell_1$ norm [31, Example 6.8])

Let  $x_j$  denote the  $j$ -th element of  $\mathbf{x} = [x_1, \dots, x_d]^T \in \mathbb{R}^d$ , and let  $[\mathbf{prox}_{\lambda \|\cdot\|_1}(\mathbf{x})]_j$  denote the  $j$ -th element of the  $d$ -dimensional  $\mathbf{prox}_{\lambda \|\cdot\|_1}(\mathbf{x})$  mapping. The  $j$ -th element of proximal mapping of the  $\ell_1$  norm is:

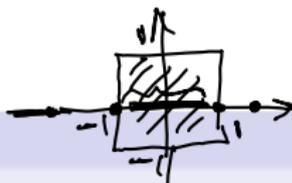
$$[\mathbf{prox}_{\lambda \|\cdot\|_1}(\mathbf{x})]_j = \max(0, |x_j| - \lambda) \text{sign}(x_j) = s_\lambda(x_j) = \begin{cases} x_j - \lambda & \text{if } x_j > \lambda \\ 0 & \text{if } |x_j| < \lambda \\ x_j + \lambda & \text{if } x_j \leq -\lambda, \end{cases} \quad (66)$$

*Handwritten annotations:* An arrow points from the  $\max$  term to the  $0$  case in the piecewise function. Another arrow points from the  $\text{sign}(x_j)$  term to the  $x_j - \lambda$  and  $x_j + \lambda$  cases. A third arrow points from the  $s_\lambda(x_j)$  term to the entire piecewise function.

for all  $j \in \{1, \dots, d\}$ . Eq. (66) is called the soft-thresholding function, denoted here by  $s_\lambda(\cdot)$ .



# Proximal mapping of $l_1$ norm



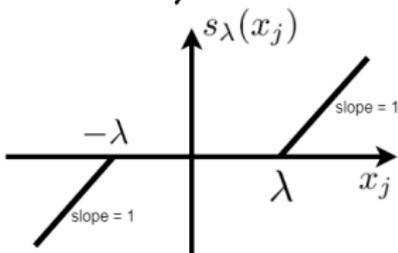
## Proof.

Let  $g(\cdot) = \|\cdot\|_1$  and  $\mathcal{B}$  be the unit  $l_\infty$  ball because  $l_\infty$  is the dual norm of  $l_1$ . The  $j$ -th element of projection is:

$$[\Pi_{\mathcal{B}}(\mathbf{x})]_j = \begin{cases} 1 & \text{if } x_j \geq 1 \\ x_j & \text{if } |x_j| < 1 \\ -1 & \text{if } x_j \leq -1 \end{cases}$$

$$\Rightarrow [\text{prox}_{\lambda \|\cdot\|_1}(\mathbf{x})]_j \stackrel{(64)}{=} x_j - \lambda [\Pi_{\mathcal{B}}(\frac{x}{\lambda})]_j = \begin{cases} x_j - \lambda & \text{if } x_j \geq \lambda \\ 0 & \text{if } |x_j| < \lambda \\ x_j + \lambda & \text{if } x_j \leq -\lambda \end{cases}$$

$$\frac{x_j}{\lambda} - \lambda \left( \frac{x_j}{\lambda} \right) = 0$$



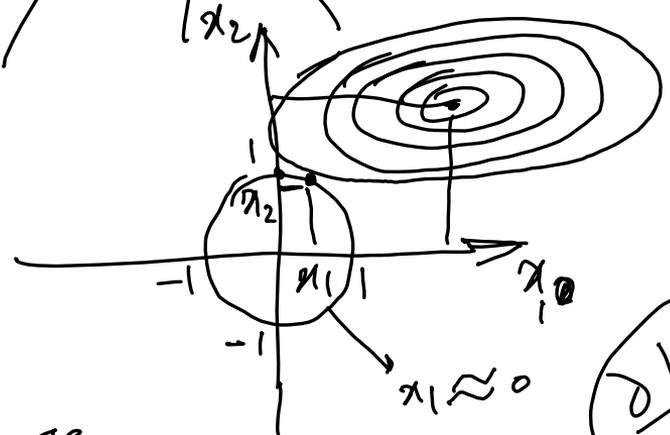
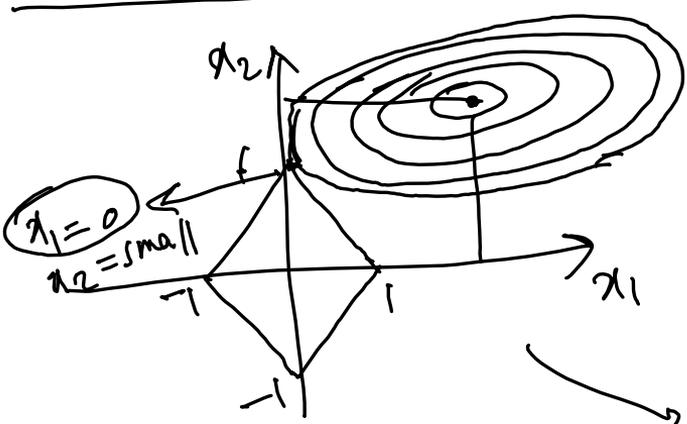
$$\min_{\alpha} f(\alpha) + \lambda \|\alpha\|_2$$

$$\min_{\alpha} f(\alpha) + \lambda \|\alpha\|_1$$

Lasso

$$\min_{\alpha} f(\alpha) \rightarrow \alpha^* = \begin{bmatrix} 10^6 \\ 10^5 \\ \vdots \\ 10^{11} \end{bmatrix}$$

Ridge



sparse

Lasso regression

$$\frac{\partial \|\alpha\|_1}{\partial \alpha}$$

$$\frac{\partial \|\alpha\|_2}{\partial \alpha} = 2\alpha$$

# Proximal point algorithm

Recall Eq. (58):

$$\text{prox}_{\lambda g}(x) := \arg \min_{\mathbf{u}} \left( g(\mathbf{u}) + \frac{1}{2\lambda} \|\mathbf{u} - \mathbf{x}\|_2^2 \right).$$

The term  $g(\mathbf{u}) + 1/(2\lambda)\|\mathbf{u} - \mathbf{x}\|_2^2$  in this equation is strongly convex; hence, the proximal point,  $\text{prox}_{\lambda g}(x)$ , is unique. ★

## Lemma

The point  $\mathbf{x}^*$  minimizes the function  $f(\cdot)$  if and only if  $\mathbf{x}^* = \text{prox}_{\lambda f}(\mathbf{x}^*)$ .

$$\min_{\lambda} f(x) \rightarrow x^*$$

# Proximal point algorithm

- Consider the following optimization problem:

$$\boxed{\underset{x}{\text{minimize}} \quad f(x).} \quad \star$$

$$\nabla f(x) \stackrel{\text{set}}{=} 0 \rightarrow \textcircled{x^*}$$

- Proximal point algorithm**, also called **proximal minimization**, was proposed in 1976 [32]. It finds the optimal point of this problem by iteratively updating the solution as:

*In each iter*  $\rightarrow$   $\boxed{x^{(k+1)} := \text{prox}_{\lambda f}(x^{(k)})} \stackrel{(58)}{=} \boxed{\arg \min_u \left( f(u) + \frac{1}{2\lambda} \|u - x^{(k)}\|_2^2 \right)}, \quad \star \quad (67)$

until convergence.

- $\lambda$  can be seen as a parameter related to the step size.
- In other words, proximal gradient method applies gradient descent on the Moreau envelope  $M_{\lambda f}(x)$ , recall Eq. (59):

$$\boxed{M_{\lambda g}(x) := \inf_u \left( g(u) + \frac{1}{2} \|u - x\|_2^2 \right),}$$

rather than on the function  $f(\cdot)$  itself.

*a ~~full~~ gradient descent until convergence*

# Composite problems

## Definition (Composite objective function [14])

In optimization, if a function is stated as a summation of two terms,  $f(x) + g(x)$  it is called a composite function and its optimization is a composite optimization problem.

- Consider the following optimization problem:

$$\underset{x}{\text{minimize}} \quad f(x) + g(x) \quad (68)$$

where  $f(x)$  is a smooth function and  $g(x)$  is a convex function which is not smooth necessarily. This is a composite optimization problem.

- Composite problems are widely used in machine learning and **regularized** problems because  $f(x)$  can be the cost function to be minimized while  $g(x)$  is the **penalty or regularization** term [33].

# Proximal gradient method

- For solving problem (68),  $\underset{x}{\text{minimize}} f(x) + g(x)$ , we can approximate the function  $f(\cdot)$  by its quadratic approximation around point  $x$  because it is smooth (differentiable):

$$f(\mathbf{u}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{u} - \mathbf{x}) + \frac{1}{2\eta} \|\mathbf{u} - \mathbf{x}\|_2^2,$$

where we have replaced  $\nabla^2 f(\mathbf{x})$  with scaled identity matrix,  $(1/\eta)I$ .

- Hence, the solution of problem (68) can be approximated as:

$$\begin{aligned} \mathbf{x} &= \underset{\mathbf{u}}{\text{arg min}} (f(\mathbf{u}) + g(\mathbf{u})) \rightarrow \underset{\text{chair}}{\text{arg min}} (f(\text{chair}) + g(\text{chair})) \\ &\approx \underset{\mathbf{u}}{\text{arg min}} (f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{u} - \mathbf{x}) + \frac{1}{2\eta} \|\mathbf{u} - \mathbf{x}\|_2^2 + g(\mathbf{u})) \\ &\stackrel{(a)}{=} \underset{\mathbf{u}}{\text{arg min}} (\underbrace{\|\nabla f(\mathbf{x})\|_2^2}_{\text{constant}} + \nabla f(\mathbf{x})^\top (\mathbf{u} - \mathbf{x}) + \frac{1}{2\eta} \|\mathbf{u} - \mathbf{x}\|_2^2 + g(\mathbf{u})) \\ &\stackrel{(b)}{=} \underset{\mathbf{u}}{\text{arg min}} \left( \frac{1}{2\eta} \underbrace{\|(\mathbf{u} - \mathbf{x}) + \eta \nabla f(\mathbf{x})\|_2^2}_{\text{binomial square term}} + g(\mathbf{u}) \right) \\ &= \underset{\mathbf{u}}{\text{arg min}} \left( \frac{1}{2\eta} \|\mathbf{u} - (\mathbf{x} - \eta \nabla f(\mathbf{x}))\|_2^2 + g(\mathbf{u}) \right), \end{aligned} \tag{69}$$

where (a) is because  $f(\mathbf{x})$  is a constant in minimization w.r.t.  $\mathbf{u}$ , so it can be simply replaced with another constant  $\|\nabla f(\mathbf{x})\|_2^2$ , and (b) is because of writing the three first terms as the binomial square term.  $\rightarrow \|a+b\|_2^2 = a^2 + b^2 + 2ab$

# Proximal gradient method

as step A gradient descent for  $f$

- We found Eq. (69):

$$x = \arg \min_{\mathbf{u}} \left( \frac{1}{2\eta} \|\mathbf{u} - (\mathbf{x} - \eta \nabla f(\mathbf{x}))\|_2^2 + g(\mathbf{u}) \right).$$

- The first term in this equation keeps the solution close to the solution of gradient descent for minimizing the function  $f(\cdot)$  (see Eq. (10)) and the second term in this equation makes the function  $g(\cdot)$  small.
- **Proximal gradient method**, also called **proximal gradient descent**, uses Eq. (69) for solving the composite problem (68). It was first proposed in 2013 [14] and also in [34] for  $g = \|\cdot\|_1$ . It finds the optimal point by iteratively updating the solution as:

$$\begin{aligned} \mathbf{x}^{(k+1)} &:= \arg \min_{\mathbf{u}} \left( \frac{1}{2\eta^{(k)}} \|\mathbf{u} - (\mathbf{x}^{(k)} - \eta^{(k)} \nabla f(\mathbf{x}^{(k)}))\|_2^2 + g(\mathbf{u}) \right) \\ &\stackrel{(58)}{=} \text{prox}_{\eta^{(k)} g}(\mathbf{x}^{(k)} - \eta^{(k)} \nabla f(\mathbf{x}^{(k)})), \end{aligned} \quad (70)$$

until convergence, where  $\eta^{(k)}$  is the step size which can be fixed or found by line-search.

- In Eq. (68), minimize  $f(\mathbf{x}) + g(\mathbf{x})$ , the function  $g(\cdot)$  can be a regularization term such as  $\ell_2$  or  $\ell_1$  norm. In these cases, we use Lemmas 12 and 13 for calculating Eq. (70).

## Constrained First-order Optimization

# Projected gradient method

- Projected gradient method**, proposed in 2003 [35], also called **gradient projection method** and **projected gradient descent**, considers  $g(x)$  to be the indicator function  $\mathbb{I}(x \in S)$  in problem (68), minimize  $f(x) + g(x)$ .

- This optimization problem is a constrained problem which can be restated to:

$$\begin{array}{l} \text{minimize} \\ \text{s.t. } x \in S \end{array} \quad \text{Pr} \quad \Rightarrow \quad \text{minimize}_x \quad f(x) + \mathbb{I}(x \in S), \quad (71)$$

because the indicator function becomes infinity if its condition is not satisfied.

- Recall Eq. (70) in proximal gradient method:

$$\mathbf{x}^{(k+1)} := \text{prox}_{\eta^{(k)} \mathbb{I}(S)}(\mathbf{x}^{(k)} - \eta^{(k)} \nabla f(\mathbf{x}^{(k)})).$$

According to this equation, the solution is updated as:

$$\left\{ \begin{array}{l} \mathbf{x}^{(k+1)} \stackrel{(70)}{:=} \text{prox}_{\eta^{(k)} \mathbb{I}(S)}(\mathbf{x}^{(k)} - \eta^{(k)} \nabla f(\mathbf{x}^{(k)})) \\ \stackrel{(63)}{=} \Pi_S(\mathbf{x}^{(k)} - \eta^{(k)} \nabla f(\mathbf{x}^{(k)})). \end{array} \right. \quad (72)$$



- So, projected gradient method performs a step of gradient descent and then projects the solution onto the set of constraint. This procedure is repeated until convergence of solution.

# Projected gradient method

- Although most often projected gradient method is used for Eq. (72), there are few other variants of projected gradient methods such as the following proposed in 2004 [36]:

$$\left\{ \begin{array}{l} \mathbf{y}^{(k)} := \Pi_{\mathcal{S}}(\mathbf{x}^{(k)} - \eta^{(k)} \nabla f(\mathbf{x}^{(k)})), \\ \mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} + \gamma^{(k)} (\mathbf{y}^{(k)} - \mathbf{x}^{(k)}), \end{array} \right. \quad (73)$$

(74) ★

where  $\eta^{(k)}$  and  $\gamma^{(k)}$  are positive step sizes at iteration  $k$ .

- In this alternating approach, we find an additional variable  $\mathbf{y}$  by gradient descent followed by projection and then update  $\mathbf{x}$  to get close to the found  $\mathbf{y}$  while staying close to the previous solution by line-search.

# Projection onto the cone of orthogonal matrices

- Singular Value Decomposition (SVD) decomposes the matrix  $X \in \mathbb{R}^{d_1 \times d_2}$  as:

$$d_1 \times d_2 \leftarrow \boxed{X} = \underbrace{U}_{d_1 \times d_1} \underbrace{\Sigma}_{d_1 \times d_1} \underbrace{V^T}_{d_1 \times d_2} \rightarrow d_2 \times d_2 \rightarrow \boxed{d_1 \times d_2}$$

where  $U \in \mathbb{R}^{d_1 \times d_1}$  and  $V \in \mathbb{R}^{d_2 \times d_2}$  are the matrices of left and right singular vectors of  $X$ , respectively.

- Consider the constraint for projection onto the cone of orthogonal matrices, i.e.,  $X^T X = I$ . In this constraint, the constraint deals with the singular values of  $X$ , because:

constraint  $\swarrow$

$$\begin{aligned} X \stackrel{\text{SVD}}{=} U \Sigma V^T &\implies X^T X = \underbrace{U \Sigma V^T}_{\substack{\uparrow I \\ \text{set } I}} \underbrace{V \Sigma U^T}_{\substack{\uparrow I \\ \text{set } I}} \stackrel{(a)}{=} U \Sigma^2 U^T \stackrel{\text{set } I}{=} I \\ &\implies U \Sigma^2 \underbrace{U^T U}_{\substack{\uparrow I \\ \text{set } I}} \stackrel{(b)}{=} U \Sigma^2 \stackrel{\text{set } I}{=} U \implies \Sigma = I \end{aligned}$$

where (a) and (b) are because  $U$  and  $V$  are orthogonal matrices.

- Therefore, the constraint  $X^T X = I$  (i.e., projecting onto the cone of orthogonal matrices) can be modeled by setting all singular values of  $X$  to one:

$$X = U \Sigma V^T$$

$$\boxed{\text{prox}_{\lambda, g}(X) = \Pi_{\mathcal{O}} = U I V^T,}$$

$$X = U \underbrace{\Sigma}_{I} V^T \quad (75)$$

where  $I \in \mathbb{R}^{d_1 \times d_2}$  is a rectangular identity matrix and  $\mathcal{O}$  denotes the cone of orthogonal matrices.

- If the constraint is scaled orthogonality, i.e.,  $X^T X = \lambda I$  with  $\lambda$  as the scale, the projection is setting all singular values to  $\lambda$  by  $U(\lambda I)V^T = \lambda U I V^T$ .

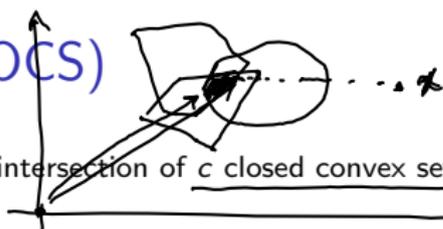
$$\min_X f(X) = \text{trace}(X^T A X)$$

$$\text{s.t. } X^T X = I \rightarrow X \in O$$

$$\min_X \text{trace}(X^T A X) + I(X \in O)$$

use projected gradient method

# Projection onto convex sets (POCS)



- Assume we want to project a point onto the intersection of  $c$  closed convex sets, i.e.,  $\bigcap_{j=1}^c \mathcal{S}_j$ ,
- We can model this by an optimization problem with a fake objective function:

$$\begin{aligned}
 & \underset{\mathbf{x}}{\text{minimize}} && \mathbf{x} \in \mathbb{R}^d, \quad \mathbb{1} \\
 & \text{subject to} && \mathbf{x} \in \mathcal{S}_1, \dots, \mathbf{x} \in \mathcal{S}_c.
 \end{aligned} \tag{76}$$

- Projection Onto Convex Sets (POCS)** solves this problem, similar to projected gradient method, by projecting onto the sets one-by-one [37]:

$$\star \quad \mathbf{x}^{(k+1)} := \underbrace{\Pi_{\mathcal{S}_1}(\underbrace{\Pi_{\mathcal{S}_2}(\dots \Pi_{\mathcal{S}_c}(\mathbf{x}^{(k)}) \dots))}_{\text{one-by-one}})}_{\text{one-by-one}} \tag{77}$$

and repeating it until convergence.

- Another similar method for solving problem (76) is the **averaged projections** which updates the solution as:

$$\star \quad \mathbf{x}^{(k+1)} := \underbrace{\frac{1}{c} \left( \underbrace{\Pi_{\mathcal{S}_1}(\mathbf{x}^{(k)})}_{\text{one-by-one}} + \dots + \underbrace{\Pi_{\mathcal{S}_c}(\mathbf{x}^{(k)})}_{\text{one-by-one}} \right)}_{\text{one-by-one}} \tag{78}$$

# Acknowledgement

- Some slides of this slide deck are inspired by the lectures of Prof. Kimon Fountoulakis at the University of Waterloo.
- Some slides of this slide deck are inspired by the lectures of Prof. Stephen Boyd at the Stanford University.
- Our tutorial also has the materials of this slide deck: [38]

# References

- [1] C. Lemaréchal, "Cauchy and the gradient method," *Doc Math Extra*, vol. 251, no. 254, p. 10, 2012.
- [2] J. Hadamard, *Mémoire sur le problème d'analyse relatif à l'équilibre des plaques élastiques encastrées*, vol. 33.  
Imprimerie nationale, 1908.
- [3] H. B. Curry, "The method of steepest descent for non-linear minimization problems," *Quarterly of Applied Mathematics*, vol. 2, no. 3, pp. 258–261, 1944.
- [4] L. Armijo, "Minimization of functions having Lipschitz continuous first partial derivatives," *Pacific Journal of mathematics*, vol. 16, no. 1, pp. 1–3, 1966.
- [5] J. Nocedal and S. Wright, *Numerical optimization*.  
Springer Science & Business Media, 2 ed., 2006.
- [6] P. Wolfe, "Convergence conditions for ascent methods," *SIAM review*, vol. 11, no. 2, pp. 226–235, 1969.
- [7] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," *nature*, vol. 323, no. 6088, pp. 533–536, 1986.
- [8] E. K. Chong and S. H. Zak, *An introduction to optimization*.  
John Wiley & Sons, 2004.

## References (cont.)

- [9] S. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.
- [10] Y. Nesterov, "A method for solving the convex programming problem with convergence rate  $O(1/k^2)$ ," in *Dokl. Akad. Nauk SSSR*, vol. 269, pp. 543–547, 1983. ★
- [11] Y. Nesterov, *Introductory lectures on convex optimization: A basic course*, vol. 87. Springer Science & Business Media, 2003.
- [12] Y. Nesterov, "On an approach to the construction of optimal methods of minimization of smooth convex functions," *Ekonomika i Mateaticheskie Metody*, vol. 24, no. 3, pp. 509–517, 1988.
- [13] Y. Nesterov, "Smooth minimization of non-smooth functions," *Mathematical programming*, vol. 103, no. 1, pp. 127–152, 2005.
- [14] Y. Nesterov, "Gradient methods for minimizing composite functions," *Mathematical Programming*, vol. 140, no. 1, pp. 125–161, 2013.
- [15] H. Robbins and S. Monro, "A stochastic approximation method," *The annals of mathematical statistics*, pp. 400–407, 1951.
- [16] L. Bottou et al., "Online learning and stochastic approximations," *On-line learning in neural networks*, vol. 17, no. 9, p. 142, 1998.

## References (cont.)

- [17] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*. MIT press, 2016.
- [18] B. Ghojogh, H. Nekoei, A. Ghojogh, F. Karray, and M. Crowley, "Sampling algorithms, from survey sampling to Monte Carlo methods: Tutorial and literature review," *arXiv preprint arXiv:2011.00901*, 2020.
- [19] N. L. Roux, M. Schmidt, and F. Bach, "A stochastic gradient method with an exponential convergence rate for finite training sets," in *Advances in Neural Information Processing Systems*, vol. 25, 2012.
- [20] R. Johnson and T. Zhang, "Accelerating stochastic gradient descent using predictive variance reduction," *Advances in neural information processing systems*, vol. 26, pp. 315–323, 2013.
- [21] J. Duchi, E. Hazan, and Y. Singer, "Adaptive subgradient methods for online learning and stochastic optimization," *Journal of machine learning research*, vol. 12, no. 7, 2011.
- [22] T. Tieleman and G. Hinton, "Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude," *COURSERA: Neural networks for machine learning*, vol. 4, no. 2, pp. 26–31, 2012.
- [23] M. Riedmiller and H. Braun, "Rprop-a fast adaptive learning algorithm," in *Proceedings of the International Symposium on Computer and Information Science VII*, 1992.

## References (cont.)

- [24] G. Hinton, N. Srivastava, and K. Swersky, "Neural networks for machine learning lecture 6a overview of mini-batch gradient descent," tech. rep., Department of Computer Science, University of Toronto, 2012.
- [25] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [26] N. Parikh and S. Boyd, "Proximal algorithms," *Foundations and Trends in optimization*, vol. 1, no. 3, pp. 127–239, 2014.
- [27] J. J. Moreau, "Proximité et dualité dans un espace hilbertien," *Bulletin de la Société mathématique de France*, vol. 93, pp. 273–299, 1965.
- [28] K. Yosida, *Functional analysis*. Springer Berlin Heidelberg, 1965.
- [29] C. Lemaréchal and C. Sagastizábal, "Practical aspects of the Moreau–Yosida regularization: Theoretical preliminaries," *SIAM journal on optimization*, vol. 7, no. 2, pp. 367–385, 1997.
- [30] J. J. Moreau, "Décomposition orthogonale d'un espace Hilbertien selon deux cônes mutuellement polaires," *Comptes rendus hebdomadaires des séances de l'Académie des sciences*, vol. 255, pp. 238–240, 1962.
- [31] A. Beck, *First-order methods in optimization*. SIAM, 2017.

## References (cont.)

- [32] R. T. Rockafellar, "Monotone operators and the proximal point algorithm," *SIAM journal on control and optimization*, vol. 14, no. 5, pp. 877–898, 1976.
- [33] B. Ghojogh and M. Crowley, "The theory behind overfitting, cross validation, regularization, bagging, and boosting: tutorial," *arXiv preprint arXiv:1905.12787*, 2019.
- [34] A. Beck and M. Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems," *SIAM journal on imaging sciences*, vol. 2, no. 1, pp. 183–202, 2009.
- [35] A. N. Iusem, "On the convergence properties of the projected gradient method for convex optimization," *Computational & Applied Mathematics*, vol. 22, pp. 37–52, 2003.
- [36] L. G. Drummond and A. N. Iusem, "A projected gradient method for vector optimization problems," *Computational Optimization and applications*, vol. 28, no. 1, pp. 5–29, 2004.
- [37] H. H. Bauschke and J. M. Borwein, "On projection algorithms for solving convex feasibility problems," *SIAM review*, vol. 38, no. 3, pp. 367–426, 1996.
- [38] B. Ghojogh, A. Ghodsi, F. Karray, and M. Crowley, "KKT conditions, first-order and second-order optimization, and distributed optimization: Tutorial and survey," *arXiv preprint arXiv:2110.01858*, 2021.