

## Ch. 5. Optimization

We have seen that many data analysis tasks are formulated as an optimization problem

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{or} \quad \begin{aligned} & \min_{x \in \mathbb{R}^n} f(x) \\ & \text{s.t. } g_i(x) \leq 0 \quad i=1,2,\dots,p \\ & h_i(x) = 0 \quad i=1,2,\dots,q. \end{aligned}$$

This chapter studies the optimality condition for these optimization problems. We consider optimization in  $\mathbb{R}^n$  for simplicity. Most of the results in this chapter can be extended to Hilbert spaces  $H$ .

### § 5.1 Smooth Unconstrained Optimization.

Consider unconstrained optimization

$$\min_{x \in \mathbb{R}^n} f(x). \quad (\text{OPT})$$

We assume  $f(x)$  is differentiable.

- Solvability of (OPT).

We say  $x^{(*)}$  is a solution of (OPT) if

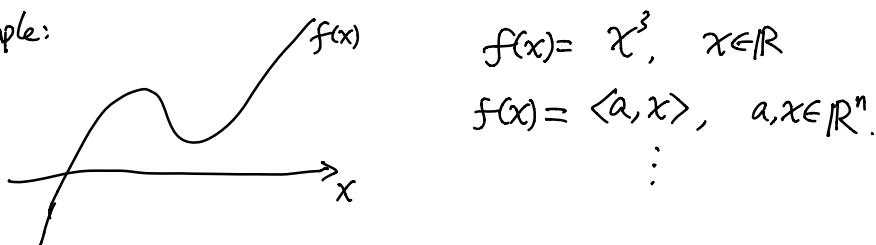
$$f(x^{(*)}) \leq f(x) \quad \forall x \in \mathbb{R}^n.$$

In this case, we write  $x^{(*)} = \arg \min_{x \in \mathbb{R}^n} f(x)$ .

We also call  $x^{(*)}$  a global minimizer of  $f$  in  $\mathbb{R}^n$ .

- The existence of a solution of (OPT) is NOT guaranteed.

Example:



$$f(x) = x^3, \quad x \in \mathbb{R}$$

$$f(x) = \langle a, x \rangle, \quad a, x \in \mathbb{R}^n.$$

⋮

- When there exists a solution of (OPT) ?

Thm: If  $f$  is continuous (i.e.,  $f(x^{(n)}) \rightarrow f(x)$  as  $x^{(n)} \rightarrow x$ )  
and coercive (i.e.,  $f(x^{(n)}) \rightarrow +\infty$  as  $\|x^{(n)}\| \rightarrow +\infty$ )

Then there exists at least a solution  $x^{(*)}$  of (OPT).

proof. Consider the sets  $S_\alpha = \{x \mid f(x) \leq \alpha\}$  for a given  $\alpha \in \mathbb{R}$ .

$S_\alpha$  is closed: Let  $x^{(n)} \rightarrow x$  and  $\{x^{(n)}\} \subset S_\alpha$ . Then  $f(x^{(n)}) \leq \alpha$ .

By the continuity,  $f(x) \leq \alpha \Rightarrow x \in S_\alpha$ .

$S_\alpha$  is bounded: Suppose  $S_\alpha$  is unbounded, i.e.,  $\exists \{x^{(n)}\} \subset S_\alpha$  s.t.  
 $\|x^{(n)}\| \rightarrow +\infty$ . Then coercivity implies  $f(x^{(n)}) \rightarrow +\infty$ .

This contradicts with  $f(x^{(n)}) \leq \alpha$ .

Therefore,  $S_\alpha$  is bounded and closed for any  $\alpha \in \mathbb{R}$ .

We choose  $\alpha$  s.t.  $S_\alpha$  is non-empty. By Weierstrass's theorem,  
any continuous function on a bounded and closed set must have  
a minimizer, So,  $\min_{x \in S_\alpha} f(x)$  has a solution, which  
is also a solution of  $\min_{x \in \mathbb{R}^n} f(x)$ .  $\square$

- Necessary condition for optimality.

Theorem: Assume  $f$  is differentiable at  $x^{(*)}$ . Then,

$$x^{(*)} = \arg \min_{x \in \mathbb{R}^n} f(x) \implies \nabla f(x^{(*)}) = 0$$

proof. By expansion,

$$f(x) = f(x^{(*)}) + \langle \nabla f(x^{(*)}), x - x^{(*)} \rangle + o(\|x - x^{(*)}\|)$$

Suppose  $\nabla f(x^{(*)}) \neq 0$ .

Then choose  $\tilde{x} = x^{(*)} - t \nabla f(x^{(*)})$  with  $t > 0$  gives

$$f(\tilde{x}) = f(x^{(*)}) - t \|\nabla f(x^{(*)})\|^2 + o(|t| \|\nabla f(x^{(*)})\|)$$

Since  $\|\nabla f(x^{(*)})\|$  is a constant of  $t > 0$ ,  $o(|t| \|\nabla f(x^{(*)})\|) = o(t)$ .

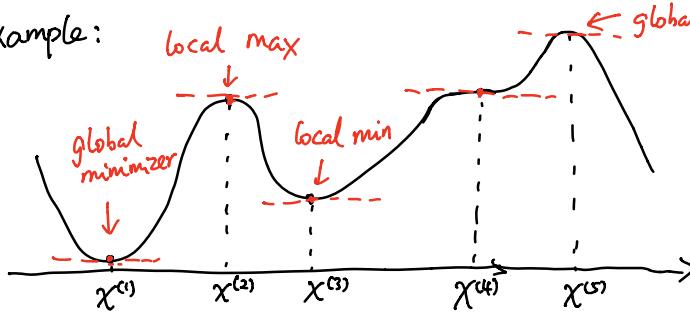
Also,  $t \|\nabla f(x^*)\|^2 = O(t)$  (big  $O$  of  $t$ )  
 $\Rightarrow$  by choosing a sufficiently small  $t$ ,  $t \|\nabla f(x^*)\|^2 > 0$  ( $|t| \|\nabla f(x^*)\|$ ),

This implies  $f(\tilde{x}) < f(x^{**})$ , which contradicts with  
 $x^{**} = \arg \min_{x \in \mathbb{R}^n} f(x)$ .  $\square$

The condition  $\nabla f(x^{**}) = 0$  is only a necessary condition.

The reverse " $\nabla f(x^{**}) = 0 \Rightarrow x^{**} = \arg \min_{x \in \mathbb{R}^n} f(x)$ " is generally NOT true.

Example:



Only  $x^{(1)}$  is a global minimizer, though the gradient at  $x^{(2)}, x^{(3)}, x^{(4)}, x^{(5)}$  are also 0.

From this example, we see that  $x^{**}$  with  $\nabla f(x^{**}) = 0$  can be

- Global minimizer. i.e.,  $x^{**} = \arg \min_{x \in \mathbb{R}^n} f(x)$  (see  $x^{(1)}$ )
- Local minimizer, i.e.,

$\exists \varepsilon, \text{s.t. } f(x^{**}) \leq f(x) \quad \forall x: \|x - x^{**}\| \leq \varepsilon$ . (see  $x^{(2)}$ )

• global max, i.e.,  $\forall x \in \mathbb{R}^n, f(x^{**}) \geq f(x)$ , (See  $x^{(5)}$ )

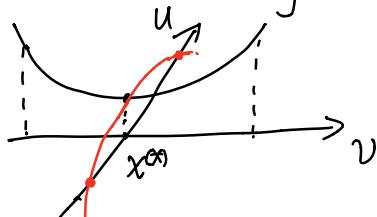
• local max, i.e.,  $\exists \varepsilon. \text{s.t. } f(x^{**}) \geq f(x) \quad \forall x: \|x - x^{**}\| \leq \varepsilon$   
 (See  $x^{(3)}$ )

• Saddle points (only for  $\mathbb{R}^n$  with  $n \geq 2$ ), i.e.,

$\exists u, v \in \mathbb{R}^n \text{ s.t. } f(x^{**}) \geq f(x^{**} + tu) \text{ for all } |t| \leq \varepsilon$ .

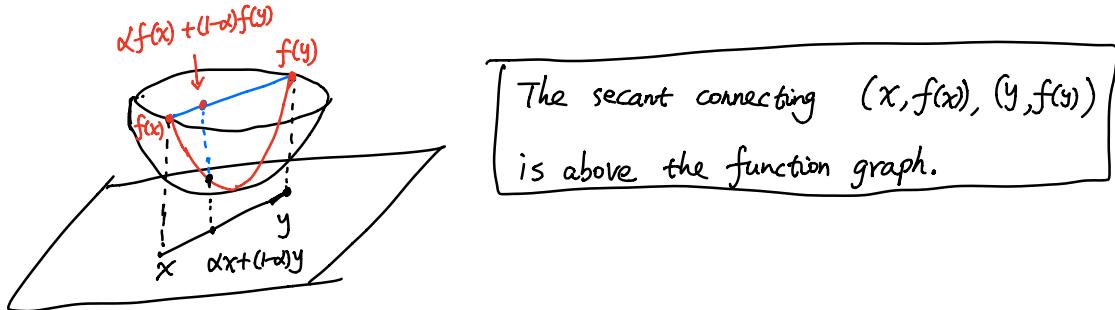
and  $f(x^{**}) \leq f(x^{**} + tv) \text{ for all } |t| \leq \varepsilon$ .

(i.e.,  $f(x^{**})$  is a local min along  $v$  and a local max along  $u$ )



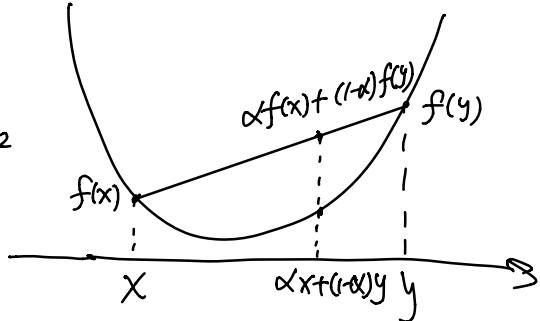
- None of the above (see  $\chi^{(4)}$ )
- Sufficient condition for optimality  
Convexity : A function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is convex if  

$$f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y) \quad \forall x, y \in \mathbb{R}^n, \alpha \in [0, 1].$$



Example 1:  $f(x) = x^2, x \in \mathbb{R}$

$$\begin{aligned}
 & f(\alpha x + (1-\alpha)y) \\
 &= (\alpha x + (1-\alpha)y)^2 \\
 &= \alpha^2 x^2 + 2\alpha(1-\alpha)xy + (1-\alpha)^2 y^2 \\
 &= \alpha x^2 + (1-\alpha)y^2 + (\alpha^2 - \alpha)x^2 \\
 &\quad + (1-\alpha)^2 y^2 + 2\alpha(1-\alpha)xy \\
 &= \alpha x^2 + (1-\alpha)y^2 - \alpha(1-\alpha)(x^2 + y^2 - 2xy) \\
 &= \alpha x^2 + (1-\alpha)y^2 - \alpha(1-\alpha)(x-y)^2 \\
 &\leq \alpha x^2 + (1-\alpha)y^2 = \alpha f(x) + (1-\alpha)f(y)
 \end{aligned}$$



Example 2:  $f(x) = \|x\|^2$

$$\begin{aligned}
 f(\alpha x + (1-\alpha)y) &= \|\alpha x + (1-\alpha)y\|^2 = \alpha^2 \|x\|^2 + (1-\alpha)^2 \|y\|^2 + 2\alpha(1-\alpha) \langle x, y \rangle \\
 &= \alpha \|x\|^2 + (1-\alpha) \|y\|^2 + 2\alpha(1-\alpha) \langle x, y \rangle + (\alpha^2 - \alpha) \|x\|^2 + (\alpha^2 - \alpha) \|y\|^2 \\
 &= \alpha \|x\|^2 + (1-\alpha) \|y\|^2 - \alpha(1-\alpha) (\|x\|^2 + \|y\|^2 - 2\langle x, y \rangle) \\
 &= \alpha f(x) + (1-\alpha)f(y) - \alpha(1-\alpha) \|x-y\|^2 \leq \alpha f(x) + (1-\alpha)f(y).
 \end{aligned}$$

Example 3:  $f(x) = \|x\|$ , where  $\|x\|$  is a norm on  $\mathbb{R}^n$  (e.g.,  $p$ -norm,  $l_1$ -norm, ...)

$$\begin{aligned} f(\alpha x + (1-\alpha)y) &= \|\alpha x + (1-\alpha)y\| \leq \|\alpha x\| + \|(1-\alpha)y\| \\ &\leq \alpha \|x\| + (1-\alpha)\|y\| \quad \forall \alpha \in [0,1], \quad x, y \in \mathbb{R}^n \end{aligned}$$

Therefore, any norm function is convex.

Example 4: Any affine function  $f$  is convex, since

$$f(\alpha x + (1-\alpha)y) = \alpha f(x) + (1-\alpha)f(y) \leq \alpha f(x) + (1-\alpha)f(y). \quad \forall \alpha \in [0,1], \quad x, y \in \mathbb{R}^n.$$

Example 5: Let  $f_1, \dots, f_n$  are convex, then  $f = \sum_{i=1}^n c_i f_i$ ,  $c_i \geq 0$ , is convex.

$$\begin{aligned} f(\alpha x + (1-\alpha)y) &= \sum_{i=1}^n c_i f_i(\alpha x + (1-\alpha)y) \leq \sum_{i=1}^n c_i (\alpha f_i(x) + (1-\alpha)f_i(y)) \\ &= \alpha \sum_{i=1}^n c_i f_i(x) + (1-\alpha) \sum_{i=1}^n c_i f_i(y) = \alpha f(x) + (1-\alpha)f(y) \end{aligned}$$

Example 6: Let  $f$  be convex and  $g$  be affine.

Then  $f \circ g$  is convex.

$$\begin{aligned} (f \circ g)(\alpha x + (1-\alpha)y) &= f(g(\alpha x + (1-\alpha)y)) = f(\alpha g(x) + (1-\alpha)g(y)) \\ &\leq \alpha f(g(x)) + (1-\alpha)f(g(y)). \end{aligned}$$

Theorem: If  $f$  is convex and differentiable, then

$$x^{**} = \arg \min_{x \in \mathbb{R}^n} f(x) \iff \nabla f(x^{**}) = 0.$$

To prove the theorem, we need a Lemma, which is also useful later.

Lemma: If  $f$  is differentiable, then

$$f \text{ is convex} \iff f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle \quad \forall x, y.$$

Proof. We prove the one variable case, i.e.,  $f: \mathbb{R} \rightarrow \mathbb{R}$ .

" $\Rightarrow$ ". By convexity,  $f(\alpha x + (1-\alpha)y) = f(x + (1-\alpha)(y-x)) \leq \alpha f(x) + (1-\alpha)f(y)$

This implies  $f(y) \geq f(x) + \frac{f(x + (1-\alpha)(y-x)) - f(x)}{(1-\alpha)(y-x)}(y-x)$

Let  $\alpha \rightarrow 1$ , then  $f(y) \geq f(x) + f'(x) \cdot (y-x)$

" $\Leftarrow$ ". Choose  $x \neq y$ , and  $\alpha \in (0,1)$ . Consider  $z = \alpha x + (1-\alpha)y$ .

$$\text{Then } f(x) \geq f(z) + f'(z)(x-z) \quad (1)$$

$$f(y) \geq f(z) + f'(z)(y-z) \quad (2)$$

$$(1) \times \alpha + (2) \times (1-\alpha) \Rightarrow f(z) \leq \alpha f(x) + (1-\alpha) f(y).$$

Next, we prove the multivariable case, i.e.,  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ .

" $\Rightarrow$ ". Consider a function  $g(t) = f(tx + (1-t)y)$ ,  $t \in \mathbb{R}$ .

So  $g'(t) = \langle \nabla f(tx + (1-t)y), x-y \rangle$  by chain rule.

Since  $f$  is convex,

$$\begin{aligned} g(\alpha s + (1-\alpha)t) &= f((\alpha s + (1-\alpha)t)x + (1-\alpha s - (1-\alpha)t)y) \\ &= f(\alpha(sx + (1-s)y) + (1-\alpha)(tx + (1-t)y)) \\ &\leq \alpha f(sx + (1-s)y) + (1-\alpha) f(tx + (1-t)y) = \alpha g(s) + (1-\alpha)g(t) \end{aligned}$$

which implies  $g(t)$  is convex.

The result from one variable case leads to

$$g(0) \geq g(1) + g'(1) \cdot (-1),$$

$$\text{i.e., } f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle$$

" $\Rightarrow$ ". Choose  $x \neq y$ ,  $\alpha \in [0, 1]$ , consider  $z = \alpha x + (1-\alpha)y$ .

$$\text{Then } f(x) \geq f(z) + \langle \nabla f(z), x-z \rangle \quad (1)$$

$$f(y) \geq f(z) + \langle \nabla f(z), y-z \rangle \quad (2)$$

$$(1) \times \alpha + (2) \times (1-\alpha) \Rightarrow f(z) \leq \alpha f(x) + (1-\alpha) f(y). \quad \blacksquare$$

With the lemma, now we can prove the theorem.

Proof of the theorem: We prove only  $\nabla f(x^{**}) = 0 \Rightarrow x^{**} = \arg \min_x f(x)$

Since  $f$  is convex and differentiable, for any  $x \in \mathbb{R}^n$ ,

$$f(x) \geq f(x^{**}) + \langle \nabla f(x^{**}), x - x^{**} \rangle \quad \}$$

By assumption,  $\nabla f(x^{**}) = 0$

$$\Rightarrow f(x) \geq f(x^{**}), \text{ i.e., } x^{**} = \arg \min_x f(x). \quad \blacksquare$$

- When the global minimizer is unique?

Strictly convex function: A function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is strictly convex if

$$f(\alpha x + (1-\alpha)y) < \alpha f(x) + (1-\alpha)f(y) \quad \forall x \neq y, \alpha \in (0, 1)$$

- Example 1:  $f(x) = x^2, x \in \mathbb{R}$ .

$$f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y) - \alpha(1-\alpha)(x-y)^2 \quad \forall x, y \in \mathbb{R}^n, \alpha \in [0, 1].$$

Therefore,  $f(\alpha x + (1-\alpha)y) < \alpha f(x) + (1-\alpha)f(y)$  if  $x \neq y, \alpha \neq 0, \alpha \neq 1$ .

Thus,  $f$  is strictly convex.

- Example 2:  $f(x) = \|x\|_2^2, x \in \mathbb{R}^n$ , is strictly convex.

- Example 3:  $f(x) = \|x\|_1, x \in \mathbb{R}^n$ , is NOT strictly convex.

Because: Let  $x = e_1, y = e_2$ ,

$$\begin{aligned} f(\alpha x + (1-\alpha)y) &= (\alpha + (1-\alpha)) = \alpha + (1-\alpha) = \alpha \|x\|_1 + (1-\alpha) \|y\|_1, \\ &= \alpha f(x) + (1-\alpha)f(y) \quad \forall \alpha \in [0, 1]. \end{aligned}$$

- Example 4: If  $f_1, f_2, \dots, f_n$  are strictly convex, then

$$f = \sum_{i=1}^n c_i f_i, \text{ where } c_i \geq 0 \text{ for all } i \text{ and not all } c_i \text{'s are 0,}$$

is strictly convex.

Theorem: Assume  $f$  is strictly convex. Then the solution of  $\min_{x \in \mathbb{R}^n} f(x)$  is unique if it exists.

proof. Suppose there are at least two solutions  $x^{(*)}, y^{(*)}$ . Then  $f(x^{(*)}) = f(y^{(*)})$

Consider  $z = \alpha x^{(*)} + (1-\alpha)y^{(*)}$  with  $\alpha \in (0, 1)$ .

Then  $z \neq x^{(*)}, z \neq y^{(*)}$ . Moreover,

$$\begin{aligned} f(z) &= f(\alpha x^{(*)} + (1-\alpha)y^{(*)}) \\ &< \alpha f(x^{(*)}) + (1-\alpha)f(y^{(*)}) = \alpha f(x^{(*)}) + (1-\alpha)f(x^{(*)}) \\ &= f(x^{(*)}), \end{aligned}$$

i.e.,  $f(z) < f(x^{(k)})$ , which contradicts with  $x^{(k)} = \arg \min_{x \in \mathbb{R}^n} f(x)$ . ☒

- Gradient Descent

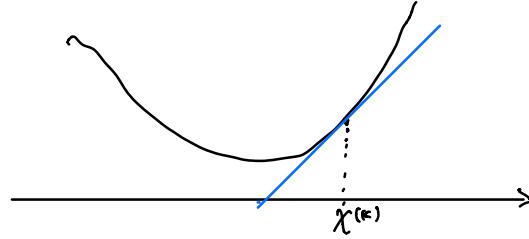
The simplest algorithm for finding a solution of  $\min_{x \in \mathbb{R}^n} f(x)$  is gradient descent in case that  $f$  is differentiable.

Let  $x^{(k)}$  be the current estimation.

Let us find  $x^{(k+1)}$  such that  $f(x^{(k+1)}) \leq f(x^{(k)})$

Given  $x^{(k)}$ , we approximate  $f(x)$  linearly via Taylor's expansion.

$$f(x) \approx f(x^{(k)}) + \langle \nabla f(x^{(k)}), x - x^{(k)} \rangle$$



This approximation is accurate only when  $x$  is close to  $x^{(k)}$ .

Therefore, instead of  $\min_{x \in \mathbb{R}^n} f(x)$ , we minimize the linear approximation and the distance of  $x$  to  $x^{(k)}$  simultaneously.

To achieve this, we solve

$$\min_{x \in \mathbb{R}^n} f(x^{(k)}) + \langle \nabla f(x^{(k)}), x - x^{(k)} \rangle + \frac{1}{2\alpha_k} \|x - x^{(k)}\|_2^2$$

↑ ↑ ↑  
To make the linear approximation small      To make  $x$  close to  $x^{(k)}$

where  $\alpha_k > 0$  is a parameter to balance the two terms.

$$\text{Let } F(x) = f(x^{(k)}) + \langle \nabla f(x^{(k)}), x - x^{(k)} \rangle + \frac{1}{2\alpha_k} \|x - x^{(k)}\|_2^2$$

It is easily checked that  $F$  is convex. Thus,

$$\begin{aligned} \min_{x \in \mathbb{R}^n} F(x) &\iff \nabla F(x) = 0 \\ &\iff \nabla f(x^{(k)}) + \frac{1}{\alpha_k}(x - x^{(k)}) = 0 \\ &\iff x = x^{(k)} - \alpha_k \nabla f(x^{(k)}) \end{aligned}$$

We obtain an algorithm:

$$x^{(k+1)} = x^{(k)} - \alpha_k \nabla f(x^{(k)}), \quad k=0, 1, 2, \dots$$

which is known as **Gradient Descent** algorithm.

- The parameter  $\alpha_k > 0$  is called a **step size** in optimization or a **learning rate** in machine learning.
- How to choose  $\alpha_k$ ?

- ① Exact linear search: choose

$$\alpha_k = \arg \min_{\alpha \geq 0} f(x^{(k)} - \alpha \nabla f(x^{(k)})).$$

- ② Linear search by back-tracking:

- Try a very large  $\alpha_k$ .
- Test whether or not  $\alpha_k$  is good enough:
  - Yes. Go to the next iteration.
  - No. Decrease  $\alpha_k$  and try again.

Example: Armijo-Goldstein backtracking

Choose a large  $\alpha_0$ . Choose  $\beta < 1$  (e.g.  $\beta = 0.9$ ).

At step  $k$ :

$$\left\{ \begin{array}{l} \alpha_k = \alpha_{k-1} \\ x^{(k+1)} = x^{(k)} - \alpha_k \nabla f(x^{(k)}) \\ \text{if } f(x^{(k)}) - f(x^{(k+1)}) < \frac{\alpha_k}{2} \|\nabla f(x^{(k)})\|_2^2 \\ \text{then } \alpha_k = \beta \alpha_k \text{ and goto } \end{array} \right.$$

otherwise go to the next iteration

There are many other criteria to test a good step size

Search "Armijo-Goldstein Condition" in google

- ③ Use a fixed step size.

- Choice of step size  $\alpha_k$ . (i.e.,  $\alpha_k = \alpha, \forall k$ )

For simplicity, we assume  $f$  is convex.

By convexity, the linear approximation gives a lower bound of  $f(x^{(k+1)})$

$$f(x^{(k+1)}) \geq f(x^{(k)}) + \langle \nabla f(x^{(k)}), x^{(k+1)} - x^{(k)} \rangle.$$

To assure  $f(x^{(k+1)}) \leq f(x^{(k)})$ , we need, however, an upper bound of  $f(x^{(k+1)})$ . To this end, we consider linear approximation at  $x^{(k+1)}$  as

$$\begin{aligned} f(x^{(k)}) &\geq f(x^{(k+1)}) + \langle \nabla f(x^{(k+1)}), x^{(k)} - x^{(k+1)} \rangle \\ \text{i.e., } f(x^{(k+1)}) &\leq f(x^{(k)}) - \langle \nabla f(x^{(k+1)}), x^{(k)} - x^{(k+1)} \rangle \\ &= f(x^{(k)}) + \langle \nabla f(x^{(k)}), x^{(k+1)} - x^{(k)} \rangle \\ &\quad + \langle \nabla f(x^{(k+1)}) - \nabla f(x^{(k)}), x^{(k+1)} - x^{(k)} \rangle \end{aligned}$$

Since, by gradient descent algorithm,

$$x^{(k+1)} - x^{(k)} = -\alpha \nabla f(x^{(k)}), \quad \forall k.$$

So,

$$f(x^{(k+1)}) \leq f(x^{(k)}) - \alpha \|\nabla f(x^{(k)})\|_2^2 + \langle \nabla f(x^{(k+1)}) - \nabla f(x^{(k)}), x^{(k+1)} - x^{(k)} \rangle$$

To obtain an upper bound of  $f(x^{(k+1)})$ , we assume

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \leq M \|x - y\|_2^2 \quad \forall x, y \in \mathbb{R}^n$$

(The gradient  $\nabla f(x)$  doesn't change too much w.r.t  $x$ )

Then,

$$\begin{aligned} f(x^{(k+1)}) &\leq f(x^{(k)}) - \alpha \|\nabla f(x^{(k)})\|_2^2 + M \|x^{(k+1)} - x^{(k)}\|_2^2 \\ &= f(x^{(k)}) - \alpha \|\nabla f(x^{(k)})\|_2^2 + M \alpha^2 \|\nabla f(x^{(k)})\|_2^2 \\ &= f(x^{(k)}) - \alpha (1 - M\alpha) \|\nabla f(x^{(k)})\|_2^2. \end{aligned}$$

To make  $f(x^{(k+1)}) < f(x^{(k)})$  for  $\nabla f(x^{(k)}) \neq 0$ , we require

$$\alpha > 0 \quad \text{and} \quad (1 - M\alpha) > 0 \quad (\text{i.e. } \alpha < \frac{1}{M}).$$

Thus, we choose  $\alpha \in (0, \frac{1}{M})$ .

Indeed, this estimation can be extended to non-convex  $f$  and

can be improved to  $\alpha \in (0, \frac{2}{M})$ .

Furthermore, we have

Theorem: Assume ①  $f$  is differentiable and  $\exists$  a solution of  $\min_{x \in \mathbb{R}^n} f(x)$ .

$$\text{② } \langle \nabla f(x) - \nabla f(y), x - y \rangle \leq M \|x - y\|_2^2, \quad \forall x, y \in \mathbb{R}^n.$$

$$\text{③ } \alpha \in (0, \frac{2}{M})$$

Then, the sequence  $\{x^{(k)}\}$  generated by

$$x^{(k+1)} = x^{(k)} - \alpha \nabla f(x^{(k)})$$

satisfies: i)  $f(x^{(k+1)}) < f(x^{(k)})$  (the function value decreases)

ii)  $\lim_{k \rightarrow \infty} \|\nabla f(x^{(k)})\|_2 = 0$  (the limit has a vanishing gradient)

- Since a vanishing gradient is not a sufficient condition for a global minimizer, the gradient descent is NOT guaranteed to find a solution of  $\min_{x \in \mathbb{R}^n} f(x)$ . It finds only a vanishing gradient with a decreasing function value of  $f$ .
- In the special case " $f$  is convex", the gradient descent will finds a solution of  $\min_{x \in \mathbb{R}^n} f(x)$ , as the global minimizer is equivalent to a vanishing gradient.

## § 5.2 Case Studies of Gradient Descent

### § 5.2.1. Least Squares

- Recall linear regression:

Given  $(x_1, y_1), \dots, (x_N, y_N)$ ,  $x_i \in \mathbb{R}^n$ ,  $y_i \in \mathbb{R}$ .

We want to find  $a \in \mathbb{R}^n$ ,  $b \in \mathbb{R}$  s.t.

$$\langle x_i, a \rangle + b \approx y_i, \quad i=1, \dots, N.$$

via minimizing the squares error

$$\min_{\substack{a \in \mathbb{R}^n \\ b \in \mathbb{R}}} \frac{1}{2} \sum_{i=1}^N (\langle x_i, a \rangle + b - y_i)^2$$

Let  $X = \begin{bmatrix} x_1^T & | & \\ \vdots & | & \\ x_N^T & | & \end{bmatrix} \in \mathbb{R}^{N \times (n+1)}$   $\beta = \begin{bmatrix} a \\ b \end{bmatrix} \in \mathbb{R}^{n+1}$   $y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \in \mathbb{R}^N$ .

Then we need to solve the least squares (LS) problem

$$\min_{\beta \in \mathbb{R}^{n+1}} \frac{1}{2} \|X\beta - y\|_2^2$$

- We consider the standard LS problem

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - b\|_2^2, \quad \text{where } A \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m.$$

Let  $f(x) = \frac{1}{2} \|Ax - b\|_2^2$ .

- $f(x)$  is convex, because:

Let  $f_1(x) = Ax - b$ ,  $f_2(y) = \frac{1}{2} \|y\|_2^2$ . Then  $f = f_2 \circ f_1$

Since  $f_1$  is affine and  $f_2$  is convex,  $f$  is convex.

- Obviously,  $f$  is differentiable. We obtain  $\nabla f(x)$  by the following:

We approximate  $f(x)$  by: At any  $x \in \mathbb{R}^n$ , for any  $y \in \mathbb{R}^n$

$$\begin{aligned} f(y) &= \frac{1}{2} \|Ay - b\|_2^2 = \frac{1}{2} \|(Ax - b) + A(y - x)\|_2^2 \\ &= \frac{1}{2} \|Ax - b\|_2^2 + \langle Ax - b, A(y - x) \rangle + \frac{1}{2} \|A(y - x)\|_2^2 \\ &= f(x) + \langle A^T(Ax - b), y - x \rangle + \frac{1}{2} \|A(y - x)\|_2^2 \\ &\quad \text{Affine approximation} \qquad \qquad \qquad \text{error} \end{aligned}$$

where we have used the equality  $\langle Ax - b, A(y - x) \rangle = \langle A^T(Ax - b), y - x \rangle$ .

Generally, we have

$$\langle u, Mv \rangle = \sum_i \sum_j m_{ij} u_i v_j = \sum_j \left( \sum_i m_{ij} u_i \right) v_j = \langle M^T u, v \rangle$$

for any vectors  $u, v$  and any matrix  $M$ .

Let us estimate the error:

$$\begin{aligned} \frac{1}{2} \|A(y-x)\|_2^2 &= \frac{1}{2} \|A \frac{(y-x)}{\|y-x\|_2}\|_2^2 \cdot \|y-x\|_2^2 \\ &\leq \frac{1}{2} \left( \max_{\|z\|_2=1} \|Az\|_2 \right) \|y-x\|_2^2 = \frac{1}{2} \left( \max_{\|z\|_2=1} \|Az\|_2 \right)^2 \|y-x\|_2^2 \end{aligned}$$

The quantity  $\max_{\|z\|_2=1} \|Az\|_2$  is finite, because

- $f(z) = \|Az\|_2$  is continuous on  $\mathbb{R}^n$ .
- the set  $\{z \mid \|z\|_2=1\}$  is bounded and closed, and non-empty.

By Weierstrass' theorem,  $\max_{\|z\|_2=1} \|Az\|_2$  exists and is finite.

Actually  $\max_{\|z\|_2=1} \|Az\|_2$  is a norm of  $A \in \mathbb{R}^{m \times n}$ , denoted by  $\|A\|_2$ .

Thus,  $\frac{1}{2} \|A(x-y)\|_2^2 \leq \frac{1}{2} \|A\|_2^2 \|y-x\|_2^2$ , such that

$$0 \leq \lim_{\|y-x\|_2 \rightarrow 0} \frac{\frac{1}{2} \|A(x-y)\|_2^2}{\|y-x\|_2} \leq \lim_{\|y-x\|_2 \rightarrow 0} \frac{\frac{1}{2} \|A\|_2^2 \|y-x\|_2^2}{\|y-x\|_2} = 0$$

Hence,

$$\nabla f(x) = A^T(Ax - b).$$

- Therefore,

$$x^{(*)} = \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - b\|_2^2 \iff A^T A x^{(*)} = A^T b$$

called the Normal equation of Least Squares.

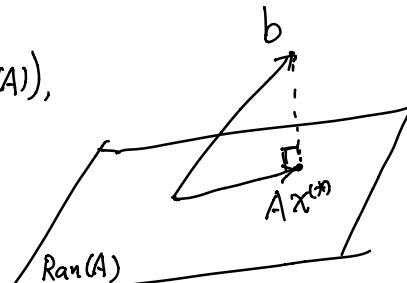
- Geometric explanation:

$Ax$  is always in the range of  $A$  ( $\text{Ran}(A)$ ),

$b$  is NOT necessarily in  $\text{Ran}(A)$

Therefore,

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - b\|_2^2 \iff \min_{y \in \text{Ran}(A)} \frac{1}{2} \|y - b\|_2^2$$



i.e.,  $Ax^{(*)}$  is the projection of  $b$  onto  $\text{Ran}(A)$ . So,

$$b - Ax^{(*)} \perp \text{Ran}(A),$$

which is the same as

$$\langle Ax^{(*)} - b, Ay \rangle = 0 \quad \forall y \in \mathbb{R}^n.$$

$$\Leftrightarrow \langle A^T(Ax^{(*)} - b), y \rangle = 0 \quad \forall y \in \mathbb{R}^n$$

$$\Leftrightarrow A^T(Ax^{(*)} - b) = 0.$$

- It can be shown that: if  $A^T A$  is invertible, then  $f(x)$  is strictly convex, and therefore the solution of least squares is unique.
- To find  $x^{(*)}$ , we may use a gradient descent

$$x^{(k+1)} = x^{(k)} - \alpha_k A^T(Ax^{(k)} - b)$$

To choose a good  $\alpha_k$ , we may use "line search", i.e., we set  $\alpha_k = \arg \min_{\alpha \in \mathbb{R}} f(x^{(k)} - \alpha A^T(Ax^{(k)} - b))$ .

In other words,  $\alpha_k$  is the optimal step size

$$\text{Let } g(\alpha) = f(x^{(k)} - \alpha A^T(Ax^{(k)} - b)).$$

It can be checked  $g(\alpha)$  is convex, and therefore

$$g'(\alpha_k) = 0,$$

$$\text{which gives } \alpha_k = \frac{\|A^T(Ax^{(k)} - b)\|_2^2}{\|A A^T(Ax^{(k)} - b)\|_2^2}.$$

This leads to the steepest descent algorithm for least squares

```

Initialize  $x^{(0)}$ 
for  $k = 0, 1, 2, \dots$ 
     $g^{(k)} = A^T(Ax^{(k)} - b)$ 
     $\alpha_k = \frac{\|g^{(k)}\|_2^2}{\|A g^{(k)}\|_2^2}$ 
     $x^{(k+1)} = x^{(k)} - \alpha_k g^{(k)}$ 
end

```

- When  $A^T A$  is non-invertible, the least squares doesn't have a unique solution.

In this case, we usually use regularization techniques.

- Ridge Regression.** (use  $\|x\|_2^2$  as regularization)

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - b\|_2^2 + \frac{\lambda}{2} \|x\|_2^2,$$

where  $\lambda > 0$  is a regularization parameter.

$$\text{Let } f(x) = \frac{1}{2} \|Ax - b\|_2^2 + \frac{\lambda}{2} \|x\|_2^2$$

- Since  $\frac{\lambda}{2} \|x\|_2^2$  is strictly convex and  $\frac{1}{2} \|Ax - b\|_2^2$  is convex,  $f(x)$  is strictly convex.

—  $\nabla f(x) = A^T(Ax - b) + \lambda x$

- There is a unique solution of Ridge Regression, which is given by the solution of

$$\begin{aligned} \nabla f(x) = 0 &\iff A^T(Ax - b) + \lambda x = 0 \\ &\iff (A^T A + \lambda I)x = A^T b. \end{aligned}$$

- We can solve it by Gradient Descent with exact line search.

- LASSO regression** (use  $\lambda \|x\|_1$  as regularization)

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1,$$

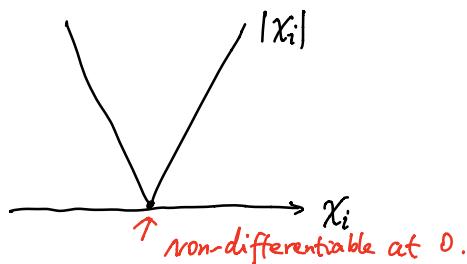
where  $\lambda > 0$  is a regularization parameter.

$$\text{Let } f(x) = \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1,$$

- $f(x)$  is convex, because  $\|x\|_1$  is convex.

- However,  $f(x)$  is NOT differentiable, because  $\|x\|_1$  is NOT.

$$\|x\|_1 = \sum_{i=1}^n |x_i|$$



### § 5.2.2. Neural Network Training

Given  $\{x^{(i)}, y_i\}_{i=1}^m$ , where  $x^{(i)} \in \mathbb{R}^n$ ,  $y_i \in \mathbb{R}$

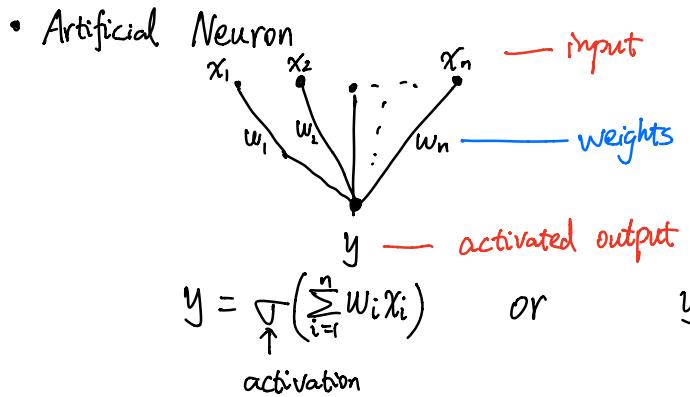
We want to find a function  $f$  such that

$$f(x^{(i)}) \approx y_i, \quad i=1, 2, \dots, m.$$

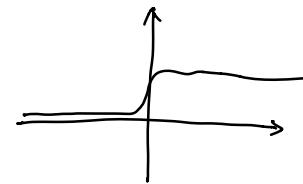
In linear regression, we choose  $f$  be an affine function.

In kernel regression, we choose  $f$  be a linear function on the feature space.

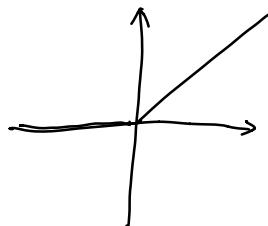
In deep learning, we choose  $f$  be a function generated from a deep neural network.



Activation function  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$



classical choice



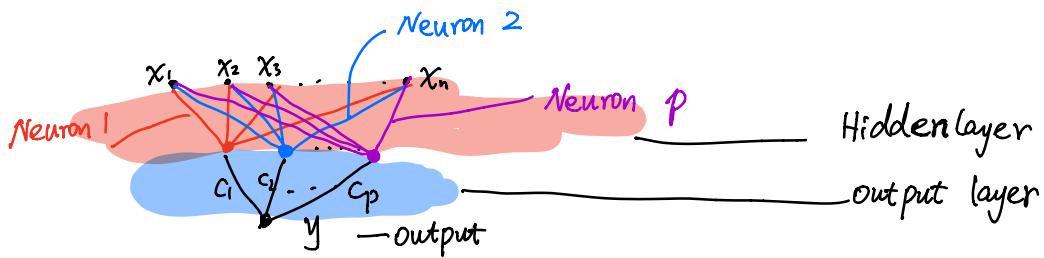
ReLU (popular recently)

... many others

- Neural Networks: Place many neurons together.

For simplicity, we first consider the following

2-layer neural network:



Let  $w^{(j)} \in \mathbb{R}^n$  be the weights in neuron  $j$ ,  $j = 1, \dots, p$ .

Denote  $W = [w^{(1)} \dots w^{(p)}] \in \mathbb{R}^{n \times p}$

Let  $c = \begin{bmatrix} c_1 \\ \vdots \\ c_p \end{bmatrix} \in \mathbb{R}^p$  be the weights in the output layer.

Then the input-output relation is

$$\begin{aligned} y &= \sum_{j=1}^p c_j \sigma(w^{(j)} \cdot x) \\ &= \langle c, \sigma(W^T x) \rangle \equiv f_{w,c}(x), \text{ where } x \in \mathbb{R}^n. \end{aligned}$$

- Neural Network training:

- Find weights  $W, c$  to make the neural network function

$f_{w,c} : \mathbb{R}^n \rightarrow \mathbb{R}$  fit the training data  $\{x^{(i)}, y_i\}_{i=1}^m$  the best.

- There are many different quantitative definitions of "the best".

We may use the least squares error:

$$\min_{\substack{W \in \mathbb{R}^{n \times p} \\ C \in \mathbb{R}^p}} \sum_{i=1}^m (f_{w,c}(x^{(i)}) - y_i)^2$$

- Let  $F_i(W, c) = (f_{w,c}(x^{(i)}) - y_i)^2$

and  $F(W, c) = \sum_{i=1}^m F_i(W, c)$ .

- Then  $F, F_i : \mathbb{R}^{n \times p} \times \mathbb{R}^p (\equiv \mathbb{R}^{(n+p)}) \rightarrow \mathbb{R}$

- Although  $F$  is NOT convex, we still apply gradient descent algorithm (or its variants) to train the neural network.

- A mysterious fact is that: despite of non-convexity, gradient-type algorithms always train a good neural network. This fact inspires the recent burst of research and applications of AI in the real world.

- We need to compute  $\nabla F(W, p)$ .
  - Since  $F$  is a function on a Euclidean space  $\mathbb{R}^{(n+1)p}$ , we have
- $$\nabla F(W, p) = \begin{bmatrix} \frac{\partial F}{\partial w^{(1)}}(W, c) \\ \vdots \\ \frac{\partial F}{\partial w^{(p)}}(W, c) \\ \frac{\partial F}{\partial c}(W, c) \end{bmatrix}, \text{ where } \frac{\partial F}{\partial w^{(j)}} \text{ stands for the gradient of } F \text{ w.r.t } w^{(j)} \text{ when other variables are fixed.}$$
- By direct calculation,

$$\begin{aligned} \frac{\partial F}{\partial w^{(j)}}(W, c) &= \frac{\partial}{\partial w^{(j)}} \sum_{i=1}^m (f_{w,c}(x^{(i)}) - y_i)^2 \\ &= \sum_{i=1}^m \left[ \frac{\partial}{\partial w^{(j)}} (f_{w,c}(x^{(i)}) - y_i)^2 \right] \end{aligned}$$

- Let  $g_1(t) = (t - y_i)^2$ ,  $g_2(w^{(j)}) = f_{w,c}(x^{(i)})$   
So,  $g_1: \mathbb{R} \rightarrow \mathbb{R}$ ,  $g_2: \mathbb{R}^n \rightarrow \mathbb{R}$ , and  $(g_1 \circ g_2)(w^{(j)}) = (f_{w,c}(x^{(i)}) - y_i)^2$

By the chain rule,

$$\begin{aligned} \frac{\partial}{\partial w^{(j)}} (f_{w,c}(x^{(i)}) - y_i)^2 &= g_1'(g_2(w^{(j)})) \cdot \nabla g_2(w^{(j)}) \\ &= 2(f_{w,c}(x^{(i)}) - y_i) \cdot \frac{\partial}{\partial w^{(j)}} f_{w,c}(x^{(i)}) \end{aligned}$$

- Because  $f_{w,c}(x^{(i)}) = \sum_{k=1}^p c_k \sigma(\langle w^{(k)}, x^{(i)} \rangle)$   
 $= c_j \sigma(\langle w^{(j)}, x^{(i)} \rangle) + \sum_{k \neq j} c_k \sigma(\langle w^{(k)}, x^{(i)} \rangle)$  Constant of  $w^{(j)}$ , denoted by A.  
 $= c_j \sigma(\langle w^{(j)}, x^{(i)} \rangle) + A$

Let  $g_3(w^{(j)}) = \langle w^{(j)}, x^{(i)} \rangle$  ( $g_3: \mathbb{R}^n \rightarrow \mathbb{R}$ )

$$g_4(t) = c_j \sigma(t) + A \quad (\text{So, } g_4: \mathbb{R} \rightarrow \mathbb{R})$$

$$\text{Obviously, } (g_4 \circ g_3)(w^{(j)}) = g_4(g_3(w^{(j)})) = c_j \sigma(\langle w^{(j)}, x^{(i)} \rangle) + A = f_{w,c}(x^{(i)})$$

By the chain rule,

$$\frac{\partial}{\partial w^{(j)}} f_{w,c}(x^{(i)}) = \nabla (g_4 \circ g_3)(w^{(j)}) = \underbrace{g_4'(g_3(w^{(j)}))}_{\text{Chain rule}} \cdot \nabla g_3(w^{(j)})$$

Direct calculation:

$$g_4'(t) = c_j \sigma'(t), \quad \nabla g_3(w^{(j)}) = x^{(i)}$$

Thus,  $\frac{\partial}{\partial w^{(i)}} f_{w,c}(x^{(i)}) = c_j \sigma'(\langle w^{(j)}, x^{(i)} \rangle) x^{(i)}$ .

- Altogether,

$$\begin{aligned} \frac{\partial F}{\partial w^{(i)}}(W, c) &= \sum_{i=1}^m \left[ \frac{\partial}{\partial w^{(i)}} (f_{w,c}(x^{(i)}) - y_i)^2 \right] \\ &= \sum_{i=1}^m 2(f_{w,c}(x^{(i)}) - y_i) c_j \sigma'(\langle w^{(j)}, x^{(i)} \rangle) x^{(i)} \\ &= 2c_j \sum_{i=1}^m [(f_{w,c}(x^{(i)}) - y_i) \cdot \sigma'(\langle w^{(j)}, x^{(i)} \rangle)] x^{(i)} \end{aligned}$$

- For  $\frac{\partial F}{\partial c}$ ,

$$\begin{aligned} \frac{\partial F}{\partial c}(W, c) &= \sum_{i=1}^m \frac{\partial}{\partial c} (f_{w,c}(x^{(i)}) - y_i)^2 \\ &= \sum_{i=1}^m 2(f_{w,c}(x^{(i)}) - y_i) \cdot \frac{\partial}{\partial c} f_{w,c}(x^{(i)}) \end{aligned}$$

Since  $\frac{\partial}{\partial c} f_{w,c}(x^{(i)}) = \frac{\partial}{\partial c} \left( \sum_{k=1}^p c_k \sigma(\langle w^{(k)}, x^{(i)} \rangle) \right)$

$$= \frac{\partial}{\partial c} \langle c, \sigma(W^T x^{(i)}) \rangle \quad \text{constant of } c.$$

$$= \sigma(W^T x^{(i)})$$

Thus,  $\frac{\partial F}{\partial c}(W, c) = 2 \sum_{i=1}^m (f_{w,c}(x^{(i)}) - y_i) \cdot \sigma(W^T x^{(i)})$

The gradient descent for neural network training is

$$\begin{cases} w^{(j,k+1)} = w^{(j,k)} - \alpha_k \cdot 2 c_j \sum_{i=1}^m [(f_{w,c^{(k)}}(x^{(i)}) - y_i) \cdot \sigma'(\langle w^{(j,k)}, x^{(i)} \rangle)] x^{(i)}, & j=1, \dots, p \\ c^{(k+1)} = c^{(k)} - \alpha_k \cdot 2 \sum_{i=1}^m (f_{w,c^{(k)}}(x^{(i)}) - y_i) \cdot \sigma((W^{(k)})^T x^{(i)}) \end{cases}$$

where  $W^{(k)} = [w^{(1,k)} \dots w^{(p,k)}]$  and  $c^{(k)}$  are the current weights.

In the computation, in addition to  $x^{(i)}$ ,  $y_i$ , the other vectors are all output of the neural networks at different stages

- $f_{w^{(k)}, c^{(k)}}(x^{(i)})$ , output of the neural network in the end
- $\sigma((W^{(k)})^T x^{(i)})$ , output of the hidden neurons
- $\langle w^{(j,k)}, x^{(i)} \rangle$ , output of the  $j$ -th hidden neuron before activation.

This algorithm is also known as **Back Propagation**.

- Stochastic Gradient descent (SGD)

In big data application,  $\{x^{(i)}, y_i\}_{i=1}^m$  with a huge  $m$ .

Therefore, the " $\sum_{i=1}^m$ " in the gradient may be too expensive.

So, we don't use ALL training data in one batch. Instead, we randomly sample some data  $\{x^{(i)}, y_i\}_{i \in I}$ , where  $I \subset \{1, 2, \dots, m\}$ .

Then, we apply gradient descent to

$$\min_{W, C} \sum_{i \in I} (f_{W, C}(x^{(i)}) - y_i)^2.$$

The resulting algorithm is known as SGD or mini-batch SGD.

for  $k = 1, 2, \dots$

randomly choose  $I \subset \{1, \dots, m\}$  with  $|I|$  small.

$$w^{(j, k+1)} = w^{(j, k)} - \alpha_k \cdot 2 \sum_{i \in I} [(f_{W^{(k)}, C^{(k)}}(x^{(i)}) - y_i) \cdot \sigma'((w^{(k)})^T x^{(i)})] x^{(i)},$$

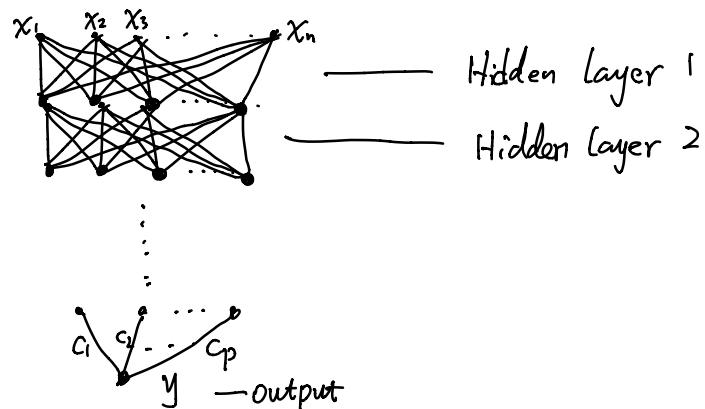
$j = 1, 2, \dots, p$ .

$$C^{(k+1)} = C^{(k)} - \alpha_k \cdot 2 \sum_{i \in I} (f_{W^{(k)}, C^{(k)}}(x^{(i)}) - y_i) \cdot \sigma((W^{(k)})^T x^{(i)})$$

end for

- Deep learning (deep neural network)

Use many hidden layers ,



Similar to one hidden layer case.

## § 6.2 Unconstrained Non-Smooth Convex Optimization

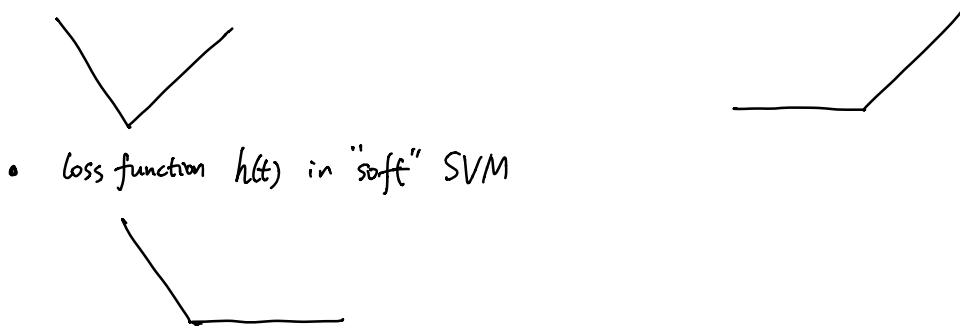
We consider

$$\min_{x \in \mathbb{R}^n} g(x)$$

where  $g(x)$  is non-differentiable but convex.

Examples of non-differentiable convex functions.

- 1-norm (in LASSO and sparsity recovery)
- ReLU function (in deep learning)



### § 6.2.1. Sub-differential / Sub-gradient. and Optimality

- To give an optimality condition, we need to extend differentiation to non-differentiable convex functions.

To this end, we first prove:

- Theorem: Assume  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is convex and differentiable. Then, for any given vector  $x \in \mathbb{R}^n$ ,

$$\{\nabla f(x)\} = \{u \in \mathbb{R}^n \mid f(y) \geq f(x) + \langle u, y-x \rangle \quad \forall y \in \mathbb{R}^n\}.$$

Proof: Since, if  $f$  is convex and differentiable,  $\nabla f(x)$  satisfies

$$f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle \quad \forall y \in \mathbb{R}^n,$$

we have

$$\{\nabla f(x)\} \subseteq \{u \in \mathbb{R}^n \mid f(y) \geq f(x) + \langle u, y-x \rangle, \forall y \in \mathbb{R}^n\}$$

It remains to prove

$$\{\nabla f(x)\} \supseteq \{u \in \mathbb{R}^n \mid f(y) \geq f(x) + \langle u, y-x \rangle, \forall y \in \mathbb{R}^n\}$$

To this end, let  $u \in \mathbb{R}^n$  be satisfying  $f(y) \geq f(x) + \langle u, y-x \rangle \quad \forall y \in \mathbb{R}^n$

Since  $f$  is convex,  $f(2x-y) \geq f(x) + \langle \nabla f(x), x-y \rangle$

$$\text{Therefore, } f(y) + f(2x-y) - 2f(x) \geq \langle u - \nabla f(x), y-x \rangle \quad \forall y \in \mathbb{R}^n. \quad \dots (1)$$

Similarly,  $f(2x-y) \geq f(x) + \langle u, x-y \rangle$

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

$$\Rightarrow f(y) + f(2x-y) - 2f(x) \geq - \langle u - \nabla f(x), y-x \rangle \quad \forall y \in \mathbb{R}^n. \quad \dots (2)$$

$$\text{Since } f(2x-y) + f(y) - 2f(x) = 2(f(2x-y) + f(y) - f(\frac{1}{2}(2x-y) + \frac{1}{2}y)) \geq 0,$$

$$\text{combining (1) \& (2) gives } |\langle u - \nabla f(x), y-x \rangle| \leq |f(y) + f(2x-y) - 2f(x)|.$$

$$\begin{aligned} \text{Also, } f(y) &= f(x) + \langle \nabla f(x), y-x \rangle + o(\|x-y\|) \\ f(2x-y) &= f(x) + \langle \nabla f(x), x-y \rangle + o(\|x-y\|) \end{aligned} \quad \Rightarrow f(y) + f(2x-y) - 2f(x) = o(\|x-y\|),$$

$$\text{i.e., } \lim_{y \rightarrow x} \frac{|f(y) + f(2x-y) - 2f(x)|}{\|x-y\|} = 0.$$

which implies

$$\begin{aligned} 0 &\leq \lim_{y \rightarrow x} \frac{|\langle u - \nabla f(x), y-x \rangle|}{\|y-x\|} \leq \lim_{y \rightarrow x} \frac{|f(y) + f(2x-y) - 2f(x)|}{\|x-y\|} = 0 \\ \Rightarrow \lim_{y \rightarrow x} \frac{|\langle u - \nabla f(x), y-x \rangle|}{\|y-x\|} &= 0. \end{aligned}$$

Now, if we take  $y = x+t(u-\nabla f(x))$  with  $t \in \mathbb{R}$ ,

$$\text{then } \lim_{t \rightarrow 0} \frac{|\langle u - \nabla f(x), y-x \rangle|}{\|y-x\|} = \lim_{t \rightarrow 0} \frac{t \|u - \nabla f(x)\|^2}{t \|u - \nabla f(x)\|} = \|u - \nabla f(x)\|$$

$$\lim_{y \rightarrow x} \frac{|\langle u - \nabla f(x), y-x \rangle|}{\|y-x\|} = 0$$

Therefore  $u = \nabla f(x)$ . \(\blacksquare\)

- Therefore, we can use the set  $\{u \mid f(y) \geq f(x) + \langle u, y-x \rangle, \forall y \in \mathbb{R}^n\}$  to define gradient of a convex differentiable function  $f$ .

Notice that there is no limit in this equivalent definition. So, we can extend it to convex non-differentiable function.

- Sub-differential / sub-gradient:

Given a convex function  $g: \mathbb{R}^n \rightarrow \mathbb{R}$  (differentiable or non-differentiable), its sub-differential at  $x \in \mathbb{R}^n$ , denoted by  $\partial g(x)$ , is defined by

$$\partial g(x) = \left\{ u \in \mathbb{R}^n \mid g(y) \geq g(x) + \langle u, y-x \rangle \quad \forall y \in \mathbb{R}^n \right\}.$$

Any element in  $\partial g(x)$  is called a sub-gradient.

- From previous theorem, if  $g$  is differentiable at  $x$ , then

$$\partial g(x) = \{\nabla f(x)\}.$$

- What if  $g$  is not differentiable at  $x \in \mathbb{R}^n$ ?

Theorem: If  $g$  is NOT differentiable at  $x \in \mathbb{R}^n$ , then  $\partial g(x)$  contains more than one element.

Proof. We prove only the 1-D case, i.e.,  $g: \mathbb{R} \rightarrow \mathbb{R}$  is convex.

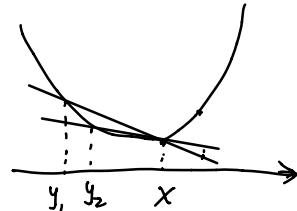
Let  $y_1 \leq y_2 < x$ . Then,  $\exists t \in (0,1)$  s.t.

$$y_2 = ty_1 + (1-t)x = x + t(y_1 - x)$$

$$\text{Let } S(y) = \frac{g(y) - g(x)}{y - x}.$$

Since  $g$  is convex,

$$\begin{aligned} S(y_2) - S(y_1) &= \frac{g(y_2) - g(x)}{y_2 - x} - \frac{g(y_1) - g(x)}{y_1 - x} \\ &= \frac{g(ty_1 + (1-t)x) - g(x)}{ty_1 + (1-t)x - x} - \frac{g(y_1) - g(x)}{y_1 - x} \\ &\stackrel{(y_2 - x < 0)}{\geq} \frac{t(g(y_1) + (1-t)g(x)) - g(x)}{t(y_1 - x)} - \frac{g(y_1) - g(x)}{y_1 - x} \\ &= \frac{t(g(y_1) - g(x))}{t(y_1 - x)} - \frac{(y_1 - g(x))}{y_1 - x} = 0 \end{aligned}$$



Therefore,  $S(y)$  is monotonically non-increasing on  $[y_1, x]$ .

Furthermore, let  $z > x$ . Then, due to the convexity,

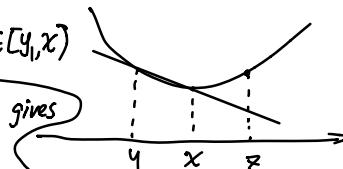
$$(A) \dots g(z) \geq g(x) + S(y)(z-x), \quad \forall y \in [y_1, x]$$

because otherwise letting  $t = \frac{z-x}{z-y} \in (0,1)$  gives

$$x = z + t(y-z) = ty + (1-t)z$$

$$(x) > g(z) - \frac{g(y) - g(x)}{y - x} (z-x)$$

$$= g(z) - \frac{g(y) - g(x)}{(t-1)(z-y)} \cdot t(z-y) = g(z) + \frac{t(g(y) - g(x))}{1-t}$$



$$\begin{aligned} &\Rightarrow (1-t)g(x) > (1-t)g(z) + t g(y) - t g(x) \\ \Rightarrow g(x) &= g(ty + (1-t)z) > (1-t)g(z) + t g(y). \end{aligned}$$

contradiction with convexity of  $g$ .

Eq. (A) implies  $s(y) \leq \frac{g(z)-g(x)}{z-x} \equiv s(z)$ .

To sum up,  $s(y)$  is monotonically non-increasing on  $[y, x]$   
and  $s(y) \leq s(z) \quad \forall z > x$ .

Therefore,  $\lim_{y \rightarrow x^-} s(y) \equiv D_- g(x)$  exists and

Similarly,  $\lim_{z \rightarrow x^+} s(z) \equiv D_+ g(x)$  exists.

Also,  $s(y) \leq D_- g(x) \leq D_+ g(x) \leq s(z) \quad \forall y < x \text{ and } z > x$ .

Choose  $u \in [D_- g(x), D_+ g(x)]$ . Then:

$$\begin{aligned} \text{Case } y < x: \quad g(y) &= g(x) + \left( \frac{g(y)-g(x)}{y-x} \right) (y-x) \\ &= g(x) + s(y)(y-x) \geq (x) + u(y-x) \end{aligned}$$

$$\begin{aligned} \text{Case } y > x: \quad g(y) &= g(x) + \left( \frac{g(y)-g(x)}{y-x} \right) (y-x) \\ &= g(x) + s(y)(y-x) \geq (x) + u(y-x) \end{aligned}$$

$$\text{Case } y = x: \quad g(y) = g(x) + u(y-x)$$

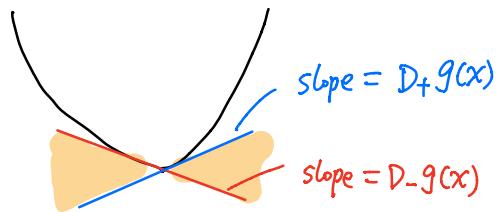
Thus,  $\partial g(x) = [D_- g(x), D_+ g(x)]$ , which is not empty.

When  $g$  is differentiable,  $D_- g(x) = D_+ g(x)$ , i.e., there is only one element in  $\partial g(x)$ .

When  $g$  is Not differentiable,  $D_- g(x) < D_+ g(x)$ ,  $\partial g(x)$  is an interval and contains infinitely many elements.  $\blacksquare$

- From the proof, we see that: the one-sided derivative of a convex function  $g: \mathbb{R} \rightarrow \mathbb{R}$  always exists. Also,

$$\partial g(x) = [D_- g(x), D_+ g(x)]. \quad \forall x \in \mathbb{R}.$$



Example 1:  $g(x) = |x| \quad x \in \mathbb{R}$ .

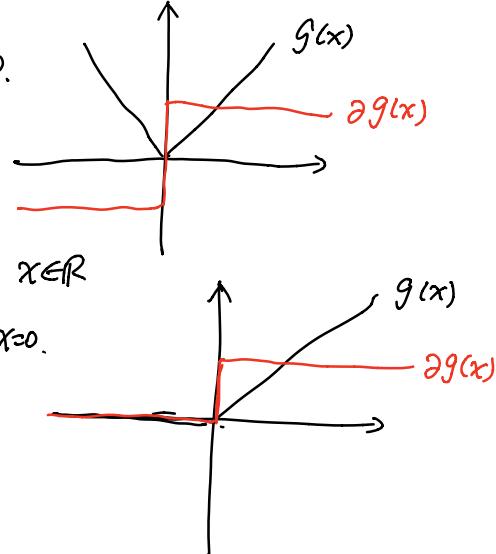
$g(x)$  is NOT differentiable at only  $x=0$ .

$$\partial g(x) = \begin{cases} \{-1\} & \text{if } x < 0 \\ [-1, 1] & \text{if } x = 0 \\ \{1\} & \text{if } x > 0 \end{cases}$$

Example 2:  $g(x) = \begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{if } x \geq 0 \end{cases} \quad x \in \mathbb{R}$

$g(x)$  is NOT differentiable at only  $x=0$ .

$$\partial g(x) = \begin{cases} \{0\} & \text{if } x < 0 \\ [0, 1] & \text{if } x = 0 \\ \{1\} & \text{if } x > 0 \end{cases}$$



Example 3:  $g(x) = \|x\|_2, \quad x \in \mathbb{R}^n$ .

$g(x)$  is differentiable at  $x \neq 0$ , and  $\nabla g(x) = \frac{x}{\|x\|_2}$ . if  $x \neq 0$ .

Therefore  $\partial g(x) = \left\{ \frac{x}{\|x\|_2} \right\}$  if  $x \neq 0$ .

Let us find  $\partial g(0)$  by definition.

$$\begin{aligned} \partial g(0) &= \{ u \in \mathbb{R}^n \mid g(y) \geq g(0) + \langle u, y \rangle, \forall y \in \mathbb{R}^n \} \\ &= \{ u \in \mathbb{R}^n \mid \|y\|_2 \geq \langle u, y \rangle, \forall y \in \mathbb{R}^n \} \equiv S \end{aligned}$$

Claim:  $S = U \equiv \{ u \in \mathbb{R}^n \mid \|u\|_2 \leq 1 \}$ .

Proof. For any  $u \in U$ ,  $\|u\|_2 \leq 1$ , which implies

$$\langle u, y \rangle \leq \|u\|_2 \|y\|_2 \leq \|y\|_2, \quad \forall y \in \mathbb{R}^n \implies u \in S.$$

$$\text{So. } U \subseteq S.$$

For any  $u \notin U$ , i.e.,  $\|u\|_2 > 1$ , we choose  $y = u$ .

$$\text{Then } \langle u, y \rangle = \|u\|_2^2 = \|u\|_2 \|y\|_2 > \|y\|_2$$

$$\text{So. } u \notin S. \quad \text{Thus, } U^c \subseteq S^c \implies S \subseteq U.$$

Therefore  $S = U$ .  $\square$

Altogether,  $\partial \|x\|_2 = \begin{cases} \left\{ \frac{x}{\|x\|_2} \right\} & \text{if } x \neq 0 \\ \{u \in \mathbb{R}^n \mid \|u\|_2 \leq 1\} & \text{if } x = 0. \end{cases}$

- Sub-differential calculus rules:

- $\partial(\alpha g) = \alpha \partial g \quad \forall \alpha \in \mathbb{R}$ .
- $\partial(g_1 + g_2) = \partial g_1 + \partial g_2$
- If  $g(x) = f(Ax+b)$ , then  $\partial g(x) = A^T \partial f(Ax+b)$

Example 4: If  $g(x) = \sum_{i=1}^n g_i(x_i)$ , where  $g_i : \mathbb{R} \rightarrow \mathbb{R}$  is convex,

then  $\partial g(x) = \left\{ \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{pmatrix} \mid u_i \in \partial g_i(x_i) \right\}$

proof. Let  $\tilde{g}_i(x) = g_i(x_i)$ .

$$\text{Then } \partial g(x) = \partial \left( \sum_{i=1}^n \tilde{g}_i(x) \right) = \sum_{i=1}^n \partial \tilde{g}_i(x).$$

Now let's show  $\partial \tilde{g}_i(x) = \{ae_i \mid a \in \partial g_i(x_i)\}$  by direct calculation

$$\begin{aligned} \partial \tilde{g}_i(x) &= \{u \in \mathbb{R}^n \mid \tilde{g}_i(y) \geq \tilde{g}_i(x) + \langle u, y-x \rangle, \forall y \in \mathbb{R}^n\} \\ &= \{u \in \mathbb{R}^n \mid g_i(y_i) \geq g_i(x_i) + u_i \cdot (y_i - x_i) + \sum_{j \neq i} u_j (y_j - x_j), \forall y \in \mathbb{R}^n\} \end{aligned}$$

- $\partial \tilde{g}_i(x) \supseteq \{u \in \mathbb{R}^n \mid u_i \in \partial g_i(x_i), u_j = 0 \quad \forall j \neq i\} = \{ae_i \mid a \in \partial g_i(x_i)\}$

- If  $u \neq ae_i$ , where  $a \in \partial g_i(x_i)$ , and  $u \in \partial \tilde{g}_i(x)$ , then,

- If  $u_j = 0 \quad \forall j \neq i$ , then

$$g_i(y) \geq g_i(x_i) + u_i(y_i - x_i) + \sum_{j \neq i} u_j (y_j - x_j) = g_i(x_i) + u_i(y_i - x_i) \quad \forall y \in \mathbb{R}^n$$

$$\Rightarrow u_i \in \partial g_i(x_i) \Rightarrow u = ae_i, a \in \partial g_i(x_i)$$

- Therefore,  $\exists j \neq i$  s.t.  $u_j \neq 0$ . Choose  $y = (c u_j + e_j)_j$  with  $c \in \mathbb{R}$

$$g_i(y) \geq g_i(x_i) + u_i(y_i - x_i) + \sum_{j \neq i} u_j (y_j - x_j)$$

$$= g_i(x_i) + c |u_j|^2 \rightarrow +\infty \quad \text{as } c \rightarrow +\infty. \text{ contradiction.}$$

Thus,  $\partial \tilde{g}_i(x) = \{ae_i \mid a \in \partial g_i(x_i)\}$

We obtain  $\partial g(x) = \left\{ \sum_{i=1}^n a_i e_i \mid a_i \in \partial g_i(x_i) \right\}$

$$= \left\{ \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix} \mid u_i \in \partial g_i(x_i) \right\}$$

Example 5:  $\partial \|x\|_1 = \left\{ \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix} \mid u_i \in \partial |x_i| \right\}$ , because  $\|x\|_1 = \sum_{i=1}^n |x_i|$

Example 6:  $\partial \|Dx\|_1 = \left\{ D^T \begin{bmatrix} u_1 \\ \vdots \\ u_m \end{bmatrix} \mid u_i \in \partial |\Phi x|_1 \right\}$ , where  $D \in \mathbb{R}^{m \times n}$ .

- Optimality:

Fermat's Lemma: Assume  $g: \mathbb{R}^n \rightarrow \mathbb{R}$  is convex. Then

$$x^{(*)} = \arg \min_{x \in \mathbb{R}^n} g(x) \iff 0 \in \partial g(x^{(*)})$$

$$\begin{aligned} \text{Proof. } x^{(*)} = \arg \min_{x \in \mathbb{R}^n} g(x) &\iff g(x) \geq g(x^{(*)}) \quad \forall x \in \mathbb{R}^n \\ &\iff g(x) \geq g(x^{(*)}) + \langle 0, x - x^{(*)} \rangle \quad \forall x \in \mathbb{R}^n \\ &\iff 0 \in \partial g(x^{(*)}) \quad \blacksquare \end{aligned}$$

Example 1: If  $g$  is convex and differentiable, then  $\partial g(x) = \{\nabla g(x)\}$ .

$$\text{Therefore, } x^{(*)} = \arg \min_{x \in \mathbb{R}^n} g(x) \iff \nabla g(x^{(*)}) = 0.$$

Example 2:  $g(x) = \|x\|_1$ ,

$$0 \in \partial \|x^{(*)}\|_1 \iff 0 \in \left\{ \begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix} \mid u_i \in \partial |x_i| \right\} \iff x^{(*)} = 0.$$

Therefore,  $g(x) = \|x\|_1$  is minimized at 0.

Example 3:  $g(x) = \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1$ , where  $A \in \mathbb{R}^{m \times n}$ ,  $b \in \mathbb{R}^m$ .

$$\partial g(x) = A^T(Ax - b) + \lambda \partial \|x\|_1$$

$$\text{Therefore, } x^{(*)} = \arg \min_{x \in \mathbb{R}^n} g(x) \iff 0 \in A^T(Ax^{(*)} - b) + \lambda \partial \|x^{(*)}\|_1$$

$$\iff \frac{1}{\lambda} A^T(b - Ax^{(*)}) \in \partial \|x^{(*)}\|_1$$

$$\text{Therefore, } [\frac{1}{\lambda} A^T(b - Ax^{(*)})]_i \in [-1, 1].$$

$$\text{If } \left| [\frac{1}{\lambda} A^T(b - Ax^{(*)})]_i \right| < 1, \text{ then } [x^{(*)}]_i = 0.$$

## § 6.2.2. Sub-Gradient Descent

Consider  $\min_{x \in \mathbb{R}^n} g(x)$ , where  $g$  is convex but non-differentiable.

- Given  $x^{(k)} \in \mathbb{R}^n$ , let  $u^{(k)} \in \partial g(x^{(k)})$ . Then, by convexity,

$$g(x^{(k+1)}) \geq g(x^{(k)}) + \langle u^{(k)}, x^{(k+1)} - x^{(k)} \rangle$$

$$\Rightarrow g(x^{(k+1)}) - g(x^{(k)}) \geq \langle u^{(k)}, x^{(k+1)} - x^{(k)} \rangle$$

The right-hand side is minimized if

$$x^{(k+1)} = x^{(k)} - \alpha_k u^{(k)}, \quad \text{where } \alpha_k > 0.$$

This is called sub-gradient descent.

- However, the right-hand side is only a lower bound, and there might be a gap as big as  $\mathcal{O}(||x^{(k+1)} - x^{(k)}||)$ , which makes  $g(x^{(k+1)}) > g(x^{(k)})$  for any  $\alpha_k > 0$ .

Example:  $g(x) = |x|$ ,  $x^{(k)} = 0$ . Choose  $u^{(k)} = -1 \in \partial g(x^{(k)})$ .

$$\text{Then } x^{(k+1)} = x^{(k)} - \alpha_k (-1) = \alpha_k.$$

$$\text{left hand side} = g(x^{(k+1)}) - g(x^{(k)}) = \alpha_k$$

$$\text{right hand side} = \langle u^{(k)}, x^{(k+1)} - x^{(k)} \rangle = \langle -1, \alpha_k \rangle = -\alpha_k.$$

the gap is

$$g(x^{(k+1)}) - g(x^{(k)}) - \langle u^{(k)}, x^{(k+1)} - x^{(k)} \rangle = 2\alpha_k$$

$$\text{which is } 2||x^{(k+1)} - x^{(k)}||. \sim \mathcal{O}(||x^{(k+1)} - x^{(k)}||),$$

Consequently,

$$\begin{aligned} g(x^{(k+1)}) - g(x^{(k)}) &= \langle u^{(k)}, x^{(k+1)} - x^{(k)} \rangle + 2\alpha_k \\ &= -\alpha_k ||u^{(k)}||^2 + 2\alpha_k = \alpha_k > 0. \end{aligned}$$

In general, to make sub-gradient work, we have to choose  $\alpha_k \rightarrow 0$ , as  $k \rightarrow +\infty$ .

In summary, sub-gradient converges very slowly.

- We may use backward sub-gradient descent.

Let  $x^{(k)} \in \mathbb{R}^n$  be current iteration, and  $x^{(k+1)} \in \mathbb{R}^n$  be the next iteration.

Instead lower bound of  $g(x^{(k+1)}) - g(x^{(k)})$ , we estimate an upper bound of  $g(x^{(k+1)}) - g(x^{(k)})$ , i.e.,

$$g(x^{(k+1)}) - g(x^{(k)}) \leq \dots$$

So we need  $g(x^{(k)}) \geq g(x^{(k+1)}) \dots$

The idea is to expand  $g$  at  $x^{(k+1)}$ . For this, we need sub-gradient at  $x^{(k+1)}$ . Let  $u^{(k+1)} \in \partial g(x^{(k+1)})$ . Then, the convexity implies

$$g(x^{(k)}) \geq g(x^{(k+1)}) + \langle u^{(k+1)}, x^{(k)} - x^{(k+1)} \rangle,$$

$$\text{i.e., } g(x^{(k+1)}) - g(x^{(k)}) \leq \langle u^{(k+1)}, x^{(k+1)} - x^{(k)} \rangle \quad \dots \quad (\text{A})$$

The right hand side is minimized when  $x^{(k+1)} - x^{(k)} = -\alpha_k u^{(k+1)}$ , i.e.,

$$x^{(k+1)} = x^{(k)} - \alpha_k u^{(k+1)}, \quad \alpha_k \geq 0 \quad \dots \quad (\text{C})$$

This is called backward sub-gradient descent.

Then, (A) implies

$$g(x^{(k+1)}) - g(x^{(k)}) \leq -\alpha_k \|u^{(k+1)}\|^2 \quad \dots \quad (\text{B})$$

Thus,

(i). if  $\alpha_k > 0$ , then  $g(x^{(k+1)}) \leq g(x^{(k)})$ , i.e., the backward sub-gradient descent gives a monotonically non-increasing  $\{g(x^{(k)})\}_{k=0}^{+\infty}$ .

Moreover, if  $\exists$  a solution of  $\min_{x \in \mathbb{R}^n} g(x)$ , (i.e.  $\exists x^{(*)} = \arg \min_{x \in \mathbb{R}^n} g(x)$ ), then  $\lim_{k \rightarrow +\infty} g(x^{(k)}) = C$  exists.

(ii). if  $\alpha_k \geq \alpha > 0$ , then summing (B) over  $k$  gives

$$g(x^{(k+1)}) - g(x^{(0)}) \leq -\sum_{k=0}^K \alpha_k \|u^{(k+1)}\|^2 \leq -\alpha \sum_{k=0}^K \|u^{(k+1)}\|^2$$

Sending  $K \rightarrow +\infty$ ,

$$C - g(x^{(0)}) \leq -\alpha \sum_{k=1}^{\infty} \|u^{(k)}\|^2, \quad \text{i.e.,}$$

$$\sum_{k=1}^{\infty} \|u^{(k)}\|^2 \leq \frac{g(x^{(0)}) - C}{\alpha} < +\infty$$

Therefore,  $\lim_{k \rightarrow +\infty} \|u^{(k)}\| = 0$ , where  $u^{(k)} \in \partial g(x^{(k)})$ .

If  $\{\chi^{(k)}\}_{k=1}^{+\infty}$  is bounded, i.e.,  $\exists M > 0$  s.t.  $\|\chi^{(k)}\|_2 \leq M \quad \forall k$ ,  
then  $g(\chi^{(x)}) \geq g(\chi^{(k)}) + \langle u^{(k)}, \chi^{(k)} - \chi^{(x)} \rangle \geq g(\chi^{(k)}) - \|u^{(k)}\| \|\chi^{(k)} - \chi^{(x)}\|$   
 $\geq g(\chi^{(k)}) - (M + \|\chi^{(k)}\|) \|u^{(k)}\|$

Sending  $k \rightarrow +\infty$ , we obtain

$$\min_{x \in \mathbb{R}^n} g(x) = g(\chi^{(x)}) \geq C, \text{ which implies } C = \min_{x \in \mathbb{R}^n} g(x).$$

$$\text{Therefore, } \lim_{k \rightarrow +\infty} g(\chi^{(k)}) = \min_{x \in \mathbb{R}^n} g(x).$$

i.e., backward sub-gradient converges to a global minimum  
as long as  $\alpha_k \geq \alpha > 0$ .

Theorem: Consider  $\min_{x \in \mathbb{R}^n} g(x)$ , where  $g: \mathbb{R}^n \rightarrow \mathbb{R}$  is convex. Let  $\{\chi^{(k)}\}_k$   
be generated by

$$\chi^{(k+1)} = \chi^{(k)} - \alpha_k u^{(k+1)}, \text{ where } u^{(k+1)} \in \partial g(\chi^{(k+1)}).$$

Then, if  $\alpha_k \geq \alpha > 0 \quad \forall k$  and  $\min_{x \in \mathbb{R}^n} g(x)$  has a solution, then,

(i)  $\{g(\chi^{(k)})\}_k$  is monotonically non-increasing;

(ii)  $\lim_{k \rightarrow +\infty} \|u^{(k)}\|_2 = 0$ ;

(iii) If  $\|\chi^{(k)}\|_2 \leq M \quad \forall k$ , then  $\lim_{k \rightarrow +\infty} g(\chi^{(k)}) = \min_{x \in \mathbb{R}^n} g(x)$ .

- However, given  $\chi^{(k)} \in \mathbb{R}^n$ , we cannot obtain  $\chi^{(k+1)}$  directly by the iteration :

$$\chi^{(k+1)} = \chi^{(k)} - \alpha_k u^{(k+1)}, \text{ where } u^{(k+1)} \in \partial g(\chi^{(k+1)})$$

$\uparrow$   
depends on  $\chi^{(k+1)}$

We need to solve the above equation, which is equivalent to

$$\chi^{(k+1)} \in \chi^{(k)} - \alpha_k \partial g(\chi^{(k+1)})$$



$$0 \in \chi^{(k+1)} - \chi^{(k)} + \alpha_k \partial g(\chi^{(k+1)})$$

$\updownarrow \leftarrow$  Fermat's Lemma

$$\chi^{(k+1)} = \arg \min_{x \in \mathbb{R}^n} F_{\alpha_k}(x), \text{ where } F_{\alpha_k}(x) = \frac{1}{2} \|x - \chi^{(k)}\|_2^2 + \alpha_k g(x)$$

- $F_{\alpha_k}$  is continuous (because all convex function is continuous).

- $F_{\alpha_k}$  is coercive if  $\min_{x \in \mathbb{R}^n} g(x)$  exists a solution.

Therefore,  $\min_{x \in \mathbb{R}^n} F_{\alpha_k}(x)$  exists at least a solution.

- Furthermore,  $F_{\alpha_k}$  is strictly convex.

Therefore,  $\min_{x \in \mathbb{R}^n} F_{\alpha_k}(x)$  has a unique solution.

Altogether,  $\chi^{(k+1)} = \arg \min_{x \in \mathbb{R}^n} F_{\alpha_k}(x)$  is well-defined.

Thus, backward sub-gradient descent is rewritten as

$$\chi^{(k+1)} = \arg \min_{x \in \mathbb{R}^n} \left( \frac{1}{2} \|x - \chi^{(k)}\|_2^2 + \alpha_k g(x) \right), \quad k=0, 1, 2, \dots$$

- Proximity Operator

Definition: Let  $g: \mathbb{R}^n \rightarrow \mathbb{R}$  be convex. Let  $\lambda > 0$  be a parameter.

The proximity operator  $\text{prox}_{\lambda g}: \mathbb{R}^n \rightarrow \mathbb{R}^n$  of  $g$  is defined by

$$\text{prox}_{\lambda g}(y) = \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|x - y\|_2^2 + \lambda g(x), \quad \forall y \in \mathbb{R}^n \quad \blacksquare$$

- Then, the backward sub-gradient descent is

$$\chi^{(k+1)} = \text{prox}_{\alpha_k g}(\chi^{(k)}), \quad k=0, 1, 2, \dots$$

For this reason, the backward sub-gradient descent is also called the proximal algorithm.

- Example 1:  $g(x) = \|x\|_1$

$$\text{Then, } \chi = \text{prox}_{\lambda g}(y) = \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|x - y\|_2^2 + \lambda \|x\|_1$$

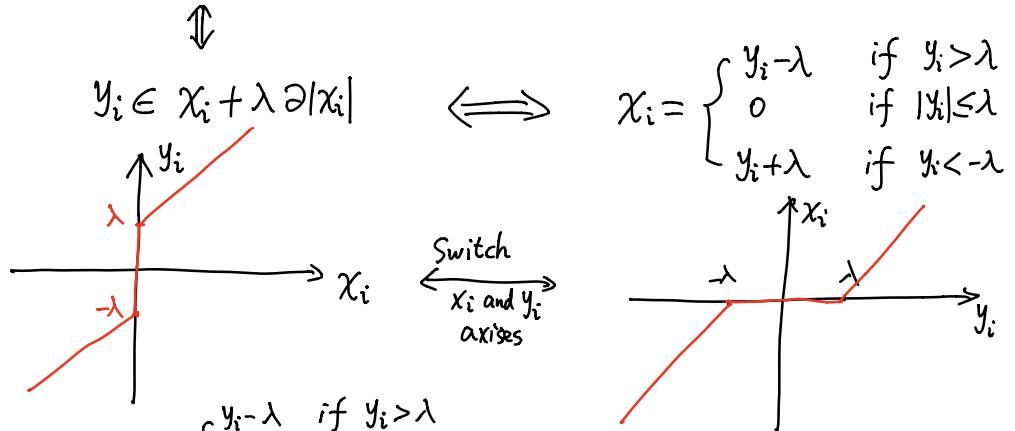


$$0 \in x - y + \lambda \partial \|x\|_1$$



$$(Recall \quad \partial \|x\|_1 = \left\{ \begin{bmatrix} u_i \\ \vdots \\ u_n \end{bmatrix} \mid u_i \in \partial |x_i| \right\})$$

$$0 \in x_i - y_i + \lambda \partial |x_i|, \quad i=1, 2, \dots, n.$$



We define  $T_\lambda(y_i) = \begin{cases} y_i - \lambda & \text{if } y_i > \lambda \\ 0 & \text{if } |y_i| \leq \lambda \\ y_i + \lambda & \text{if } y_i < -\lambda \end{cases}$

Then,  $\text{prox}_{\lambda \|\cdot\|_1}(y) = T_\lambda(y) \equiv \begin{bmatrix} T_\lambda(y_1) \\ T_\lambda(y_2) \\ \vdots \\ T_\lambda(y_n) \end{bmatrix}$  — called soft-thresholding operator.

The proximal algorithm for  $\min_{x \in \mathbb{R}^n} \|x\|_1$  gives

$$x^{(k+1)} = T_{\alpha_k}(x^{(k)}) \quad , \quad k=0,1,2, \dots$$

At each iteration, components of  $x^{(k)}$  are shrunk by  $\alpha_k$  until it becomes 0. Therefore,  $x^{(k)} \rightarrow 0 \equiv \arg \min_{x \in \mathbb{R}^n} \|x\|_1$ . \(\blacksquare\)

- Example 2:  $g(x) = \|x\|_2^2$

Then  $x \equiv \text{prox}_{\lambda g}(y) = \arg \min_{x \in \mathbb{R}^n} \left( \frac{1}{2} \|x-y\|_2^2 + \lambda \|x\|_2^2 \right)$

$$x-y + 2\lambda x = 0$$

$\Downarrow$

$$x = \frac{1}{1+2\lambda} y$$

The proximal algorithm for  $\min_{x \in \mathbb{R}^n} \|x\|_2^2$  is

$$x^{(k+1)} = \frac{1}{1+2\alpha_k} x^{(k)},$$

which obviously converges to 0, the arg min of  $\|x\|_2^2$ .

- Example 3:  $g(x) = \|x\|_2$

Then  $x \equiv \text{prox}_{\lambda g}(y) = \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|x-y\|_2^2 + \lambda \|x\|_2$

$$\begin{array}{c}
 \uparrow \\
 0 \in x - y + \lambda \partial \|x\|_2 \\
 \Downarrow \\
 y \in x + \lambda \partial \|x\|_2 \quad \dots \quad (\star\star)
 \end{array}$$

Let us find  $x$ .

If  $x \neq 0$ , then  $y = x + \lambda \frac{x}{\|x\|_2} \Rightarrow x = cy$  with  $c = \frac{1}{1 + \frac{\lambda}{\|x\|_2}} \geq 0$ .

If  $x = 0$ , then  $\Rightarrow x = cy$  with  $c = 0$ .

In any cases, the solution  $x = cy$  with  $c \geq 0$ .

So the original problem  $\Leftrightarrow \min_{c \geq 0} \frac{1}{2}(1-c)^2 \|y\|_2^2 + \lambda |c| \|y\|_2$

If  $\|y\|_2 = 0$ , obviously  $c = 0$ .

If  $\|y\|_2 > 0$ , then  $\min_{c \geq 0} \frac{1}{2}(1-c)^2 + \lambda |c|$

$$\Leftrightarrow \min_{c \geq 0} \frac{1}{2}(c-1)^2 + \frac{\lambda}{\|y\|_2} |c|$$

Since  $T_{\frac{\lambda}{\|y\|_2}}(1)$  is the solution of  $\min_{c \in \mathbb{R}} \frac{1}{2}(c-1)^2 + \frac{\lambda}{\|y\|_2} |c|$ ,

it is also the solution  $\min_{c \geq 0} \frac{1}{2}(c-1)^2 + \frac{\lambda}{\|y\|_2} |c|$ .

$$\begin{aligned}
 \text{Therefore, } c &= T_{\frac{\lambda}{\|y\|_2}}(1) = \begin{cases} 1 - \frac{\lambda}{\|y\|_2} & \text{if } 1 \geq \frac{\lambda}{\|y\|_2}, \\ 0 & \text{if } 1 \leq \frac{\lambda}{\|y\|_2}. \end{cases} \\
 &= \begin{cases} \frac{\|y\|_2 - \lambda}{\|y\|_2} & \text{if } \|y\|_2 \geq \lambda, \\ 0 & \text{if } \|y\|_2 \leq \lambda. \end{cases}
 \end{aligned}$$

$$\text{Finally, } \text{prox}_{\lambda g}(y) = x = \begin{cases} \frac{\|y\|_2 - \lambda}{\|y\|_2} y & \text{if } \|y\|_2 \geq \lambda, \\ 0 & \text{if } \|y\|_2 \leq \lambda. \end{cases}$$

The proximal algorithm for  $\min_{x \in \mathbb{R}^n} \|x\|_2$  is

$$x^{(k+1)} = \begin{cases} \frac{\|x^{(k)}\|_2 - \alpha_k}{\|x^{(k)}\|_2} x^{(k)} & \text{if } \|x^{(k)}\|_2 \geq \alpha_k, \\ 0 & \text{if } \|x^{(k)}\|_2 \leq \alpha_k \end{cases}$$

which obviously converges to 0.

- In general, it is difficult to find  $\text{prox}_{\lambda g}(\cdot)$  for a given convex

function  $g$ . Therefore, the proximal algorithm (Backward sub-gradient descent) is NOT practical.

### § 6.2.3. Case Study: LASSO regression

LASSO regression model:

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1 \quad (\text{LASSO})$$

↑ fit the given data.      ↑ promote Sparsity of the solution

- The objective function is convex.
- The objective function is continuous

The objective function is coercive, because

$$\text{when } \|x\|_2 \rightarrow \infty, \quad \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1 \geq \lambda \|x\|_1 \geq \lambda \|x\|_2 \rightarrow \infty.$$

Therefore, there exists at least one solution.

- The uniqueness of the solution is NOT guaranteed, unless  $A$  satisfies some assumptions (e.g.  $A$  is invertible)
- The solution of (LASSO) is sparse.

To see this, let  $x \in \mathbb{R}^n$  be a solution. Then

$$0 \in A^T(Ax - b) + \lambda \partial \|x\|_1$$

$\Downarrow$

$$0 \in \alpha A^T(Ax - b) + \lambda \alpha \partial \|x\|_1 \quad \forall \alpha > 0$$

$\Updownarrow$

$$0 \in x - (x - \alpha A^T(Ax - b)) + \lambda \alpha \partial \|x\|_1 \quad \forall \alpha > 0$$

$\Updownarrow$

$$x = \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|x - y\|_2^2 + \lambda \alpha \partial \|x\|_1, \quad \text{where}$$

$\Updownarrow$

$$y = x - \alpha A^T(Ax - b)$$

$$x = \text{prox}_{\lambda \alpha \|\cdot\|_1} (x - \alpha A^T(AX - b))$$

↓

$$x = T_{\lambda \alpha} (x - \alpha A^T(AX - b))$$

Recall that  $T_{\lambda \alpha}(y)$  set  $y_i$  to 0 if  $|y_i| \leq \lambda \alpha$ .

Therefore,  $x$  is a sparse if  $\lambda$  is large enough.

- Numerical solver for (LASSO).

The objective function in (LASSO) is NON-SMOOTH.

The gradient descent can not be applied.

But if the backward sub-gradient (proximal algorithm) is applied,

then we need the proximity operator, which solves

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} \|x - y\|_2^2 + \alpha (\frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1)$$

This is the same difficulty as (LASSO).

We will use a "mixed" forward and backward (sub-)gradient algorithm.

Let  $f(x) = \frac{1}{2} \|Ax - b\|_2^2$  — convex and smooth

$g(x) = \|x\|_1$  — convex but non-smooth

$$(LASSO) \iff \min_{x \in \mathbb{R}^n} f(x) + \lambda g(x)$$

↑  
forward gradient      ↑  
backward sub-gradient

We use the algorithm

$$x^{(k+1)} = x^{(k)} - \alpha_k (\nabla f(x^{(k)}) + \lambda u^{(k+1)}), \quad \text{with } u^{(k+1)} \in \partial g(x^{(k+1)})$$

It is rewritten as

$$\begin{aligned} x^{(k+1)} &\in x^{(k)} - \alpha_k \nabla f(x^{(k)}) - \alpha_k \lambda \partial g(x^{(k+1)}) \\ 0 &\in x^{(k+1)} - (x^{(k)} - \alpha_k \nabla f(x^{(k)})) + \alpha_k \lambda \partial g(x^{(k+1)}) \end{aligned} \quad \cdots \quad (A)$$

Recall that

$$\begin{aligned} x = \text{prox}_{\beta g}(y) &\iff x = \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|x - y\|_2^2 + \beta g(x) \\ &\iff 0 \in x - y + \beta \partial g(x) \end{aligned} \quad \cdots \quad (B)$$

By comparing (A) and (B), we obtain

$$x^{(k+1)} = \text{prox}_{\alpha_k \lambda g}(x^{(k)} - \alpha_k \nabla f(x^{(k)}))$$

Since  $\nabla f(x) = \nabla(\frac{1}{2} \|Ax - b\|_2^2) = A^T(Ax - b)$

$$\text{prox}_{\alpha_k \lambda g}(x) = T_{\alpha_k \lambda}(x)$$

the algorithm is

$$x^{(k+1)} = T_{\alpha_k \lambda}(x^{(k)} - \alpha_k A^T(Ax^{(k)} - b))$$

This algorithm is known as **iterative soft-thresholding algorithm**

- From the derivation, for a generic optimization

$$\min_{x \in \mathbb{R}^n} f(x) + \lambda g(x),$$

where  $f(x)$  is convex and smooth

$g(x)$  is convex but non-smooth,

We can use the **forward-backward splitting (FBS) algorithm**

$$x^{(k+1)} = x^{(k)} - \alpha_k (\nabla f(x^{(k)}) + \lambda u^{(k+1)}), \quad u^{(k+1)} \in \partial g(x^{(k+1)})$$

which is the same as

$$x^{(k+1)} = \text{prox}_{\lambda \alpha_k g}(x^{(k)} - \alpha_k \nabla f(x^{(k)}))$$

This is also known as **proximal gradient algorithm**.

- Convergence

Since  $f$  is convex and smooth,

$$(C) \dots f(x^{(k+1)}) = f(x^{(k)}) + \langle \nabla f(x^{(k)}), x^{(k+1)} - x^{(k)} \rangle + \frac{1}{2} \langle \nabla^2 f(x^{(k)})(x^{(k+1)} - x^{(k)}), x^{(k+1)} - x^{(k)} \rangle$$

(1) Since  $g$  is convex,

$$g(x^{(k)}) \geq g(x^{(k+1)}) + \langle u^{(k+1)}, x^{(k)} - x^{(k+1)} \rangle, \quad u^{(k+1)} \in \partial g(x^{(k+1)})$$

which is equivalent to

$$(D) \dots g(x^{(k+1)}) \leq g(x^{(k)}) + \langle u^{(k+1)}, x^{(k+1)} - x^{(k)} \rangle$$

(C) +  $\lambda(D)$  gives

$$[f(x^{(k+1)}) + \lambda g(x^{(k+1)})]$$

$$\leq [f(x^{(k)}) + \lambda g(x^{(k)})] + \langle \nabla f(x^{(k)}) + \lambda u^{(k+1)}, x^{(k+1)} - x^{(k)} \rangle \\ + \frac{1}{2} \langle \nabla^2 f(x^{(k)}) (x^{(k+1)} - x^{(k)}), x^{(k+1)} - x^{(k)} \rangle$$

Assuming  $\langle \nabla^2 f(x) u, u \rangle \leq M \|u\|_2^2 \quad \forall x, u \in \mathbb{R}^n$ .

$$\rightarrow \leq [f(x^{(k)}) + \lambda g(x^{(k)})] - \alpha_k \|\nabla f(x^{(k)}) + \lambda u^{(k+1)}\|_2^2 + \frac{M\alpha_k^2}{2} \|\nabla f(x^{(k)}) + \lambda u^{(k+1)}\|_2^2 \\ = [f(x^{(k)}) + \lambda g(x^{(k)})] - \alpha_k \left(1 - \frac{M\alpha_k}{2}\right) \|\nabla f(x^{(k)}) + \lambda u^{(k+1)}\|_2^2$$

Assume  $0 < \alpha_1 \leq \alpha_k \leq \alpha_u < \frac{2}{M}$  for some  $\alpha_1, \alpha_u \in \mathbb{R}$ .

Then,  $\exists \beta$  s.t.  $\alpha_k \left(1 - \frac{M}{2}\alpha_k\right) \geq \beta > 0$ .

$$\text{So, } [f(x^{(k+1)}) + \lambda g(x^{(k+1)})] - [f(x^{(k)}) + \lambda g(x^{(k)})] \leq -\beta \|\nabla f(x^{(k)}) + \lambda u^{(k+1)}\|_2^2 \dots (C)$$

$[f(x^{(k)}) + \lambda g(x^{(k)})]_k$  is a monotonically non-increasing sequence.

(2)

Assume  $\min_{x \in \mathbb{R}^n} f(x) + \lambda g(x)$  has a solution, i.e.,

$$\exists x^{(k)}, \text{ s.t. } f(x^{(k)}) + \lambda g(x^{(k)}) = \min_{x \in \mathbb{R}^n} f(x) + \lambda g(x).$$

Then,  $\lim_{k \rightarrow \infty} f(x^{(k)}) + \lambda g(x^{(k)}) = C$  exists.

Furthermore, summing (C) over  $K$ , we obtain

$$[f(x^{(K+1)}) + \lambda g(x^{(K+1)})] - [f(x^{(0)}) + \lambda g(x^{(0)})] \leq -\beta \sum_{k=0}^K \|\nabla f(x^{(k)}) + \lambda u^{(k+1)}\|_2^2$$

$$\text{So } \sum_{k=0}^K \|\nabla f(x^{(k)}) + \lambda u^{(k+1)}\|_2^2 \leq \frac{1}{\beta} [f(x^{(0)}) + \lambda g(x^{(0)})] - [f(x^{(K+1)}) + \lambda g(x^{(K+1)})]$$

$K \rightarrow +\infty$  gives  $\sum_{k=0}^{+\infty} \|\nabla f(x^{(k)}) + \lambda u^{(k+1)}\|_2^2 < +\infty$

Consequently,  $\lim_{k \rightarrow +\infty} \|\nabla f(x^{(k)}) + \lambda u^{(k+1)}\|_2 = 0$

Since  $\alpha_k (\nabla f(x^{(k)}) + \lambda u^{(k+1)}) = x^{(k)} - x^{(k+1)}$

$$0 \leq \lim_{k \rightarrow \infty} \|x^{(k)} - x^{(k+1)}\|_2 = \lim_{k \rightarrow \infty} \|\alpha_k (\nabla f(x^{(k)}) + \lambda u^{(k+1)})\|_2 \\ \leq \alpha_u \lim_{k \rightarrow \infty} \|\nabla f(x^{(k)}) + \lambda u^{(k+1)}\|_2 = 0$$

which implies  $\lim_{k \rightarrow \infty} \|x^{(k)} - x^{(k+1)}\|_2 = 0$ .

This together with the continuity of  $\nabla f(x)$  implies  $0 = \lim_{k \rightarrow \infty} \|\nabla f(x^{(k)}) - \nabla f(x^{(k+1)})\|_2$

Therefore,

$$0 \leq \lim_{k \rightarrow \infty} \|\nabla f(x^{(k+1)}) + \lambda u^{(k+1)}\|_2 \leq \lim_{k \rightarrow \infty} (\|\nabla f(x^{(k)}) + \lambda u^{(k+1)}\|_2 + \|\nabla f(x^{(k+1)}) - \nabla f(x^{(k)})\|_2) = 0,$$

i.e.,  $\boxed{\lim_{k \rightarrow \infty} \|\nabla f(x^{(k)}) + \lambda u^{(k)}\|_2 = 0}$

③ Moreover,

$$\begin{aligned} [f(x^{(k)}) + \lambda g(x^{(k)})] &\geq [f(x^{(k)}) + \lambda g(x^{(k)})] + \langle \nabla f(x^{(k)}) + \lambda u^{(k)}, x^{(k)} - x^{(k)} \rangle \\ &\geq [f(x^{(k)}) + \lambda g(x^{(k)})] - \|\nabla f(x^{(k)}) + \lambda u^{(k)}\|_2 (\|x^{(k)}\|_2 + \|x^{(k)}\|_2) \end{aligned}$$

Assume  $\exists B > 0$  s.t.  $\|x^{(k)}\|_2 \leq B \quad \forall k$ .

Then, sending  $k \rightarrow +\infty$ , we obtain

$$\min_{x \in \mathbb{R}^n} f(x) + \lambda g(x) \geq C. \quad (\text{Obviously, } C \geq \min_{x \in \mathbb{R}^n} f(x) + \lambda g(x)).$$

Thus,  $\lim_{k \rightarrow \infty} f(x^{(k)}) + \lambda g(x^{(k)}) = \min_{x \in \mathbb{R}^n} f(x) + \lambda g(x)$

To sum up,

Theorem: Consider  $\min_{x \in \mathbb{R}^n} f(x) + \lambda g(x)$ , where  $\lambda > 0$ ,  $f, g$  are convex.

Assume: ①  $\langle \nabla^2 f(x) u, u \rangle \leq M \|u\|_2^2, \forall x, u \in \mathbb{R}^n$ . (i.e.,  $f$  is smooth and has a bounded Hessian.)

②  $0 < \alpha_l \leq \alpha_k \leq \alpha_u < \frac{2}{M}$  for some  $\alpha_l, \alpha_u \in \mathbb{R}$ .

③  $\min_{x \in \mathbb{R}^n} f(x) + \lambda g(x)$  has a solution.

Then, the sequence  $\{x^{(k)}\}_k$  generated by

$$x^{(k+1)} = \text{prox}_{\lambda \alpha_k g}(x^{(k)} - \alpha_k \nabla f(x^{(k)}))$$

satisfies:

(i)  $\{f(x^{(k)}) + \lambda g(x^{(k)})\}_k$  is monotonically non-increasing,

(ii)  $\exists u^{(\infty)} \in \partial g(x^{(\infty)})$ , s.t.  $\lim_{k \rightarrow \infty} \|\nabla f(x^{(k)}) + u^{(k)}\|_2 = 0$ .

(iii) If  $\exists B > 0$  s.t.  $\|x^{(k)}\|_2 \leq B \quad \forall k$ , then

$$\lim_{k \rightarrow \infty} f(x^{(k)}) + \lambda g(x^{(k)}) = \min_{x \in \mathbb{R}^n} f(x) + \lambda g(x).$$

Remark: The condition  $\|x^{(k)}\|_2 \leq B$  is satisfied if  $f(x) + \lambda g(x)$  is coercive. (Why?).