

Gonghan Xu, SDS 631, hw5

1. (BV Exercise 9.7)(a) :  $\Delta X_{nsd}$  is a descent direction  $\therefore \nabla f(x)^T \Delta X_{nsd} < 0$  ①

By definition of  $\Delta X_{nsd}$ , we have  $\Delta X_{nsd} = \operatorname{argmin}_V \{ \nabla f(x)^T V \mid \|V\| \leq 1 \}$  ②

i. Combining ①, ②, we have  $\nabla f(x)^T \Delta X_{nsd} = \min \{ \nabla f(x)^T V \mid \|V\| \leq 1 \}$  (then use ②)

$$= -\max \{ -\nabla f(x)^T V \mid \|V\| \leq 1 \} = -\sup \{ \nabla f(x)^T (-V) \mid \| -V \| \leq 1 \} \quad (\text{because } \{ V \mid \|V\| \leq 1 \} \text{ is a compact set so } \max = \sup)$$

$$= -\sup \{ \nabla f(x)^T V' \mid \|V'\| \leq 1 \} \quad (\text{where we have set } V' = -V)$$

$$= -\sup \{ \langle \nabla f(x), V' \rangle \mid \|V'\| \leq 1 \} = -\| \nabla f(x) \|_* \quad ③, \text{ where we have used the definition of } \| \cdot \|_*$$

$$(b). \nabla f(x)^T \Delta X_{sd} = \nabla f(x)^T (\| \nabla f(x) \|_* \Delta X_{nsd}) = \| \nabla f(x) \|_* (\nabla f(x)^T \Delta X_{nsd}) \quad (\text{then use 1.(a) ③})$$

$$= \| \nabla f(x) \|_* \{ -\| \nabla f(x) \|_* \} = -\| \nabla f(x) \|_*^2 \quad ①$$

(c) Now we try to find  $V \in \mathbb{R}^n$  that minimize  $\nabla f(x)^T V + \frac{1}{2} \|V\|^2$ . ①

Note that  $V = \|V\| \cdot \frac{V}{\|V\|} \equiv \|V\| \cdot \hat{V}$  ②, where we define  $\hat{V} \equiv \frac{V}{\|V\|}$  ③ Then  $\|\hat{V}\| = 1$  ④

$$\text{Then our objective function } h(V) \equiv \nabla f(x)^T V + \frac{1}{2} \|V\|^2 = \frac{1}{2} \|V\|^2 + \|V\| \nabla f(x)^T \hat{V} \quad ⑤$$

Note that  $\|V\|$  and  $\hat{V}$  are independent on each other. So we can minimize  $h(V)$  with respect to them in a separate-two-step manner. First, with  $\|V\|$  fixed, since  $\|V\| \geq 0$  then we want to minimize  $\nabla f(x)^T \hat{V}$  in order to minimize  $h(V)$ . Then  $\hat{V}^* = \operatorname{argmin}_{\hat{V}} \{ \nabla f(x)^T \hat{V} \mid \|\hat{V}\| = 1 \} = \Delta X_{nsd}$  ⑥

Note that when  $\hat{V} = \hat{V}^* = \Delta X_{nsd}$ ,  $\nabla f(x)^T \hat{V}^* = \nabla f(x)^T \Delta X_{nsd} \leq 0$  ⑦

Then using the quadratic polynomial minimization property  $-\frac{b}{2a} = \operatorname{argmin}_x \{ ax^2 + bx + c \mid x \geq 0 \}$

$$\text{we have: } \|V\|^* = \operatorname{argmin}_{\|V\|} \left\{ \frac{1}{2} \|V\|^2 + \|V\| \nabla f(x)^T \Delta X_{nsd} \mid \|V\| \geq 0 \right\} = -\nabla f(x)^T \Delta X_{nsd} \geq 0 \quad ⑧$$

$$\therefore \|V\|^* = -\nabla f(x)^T \Delta X_{nsd} \quad ⑨ \text{ and } \hat{V}^* = \Delta X_{nsd} \quad ⑩$$

$$\text{Then using 1.(a) ③, we have } \|V\|^* = -\{ \nabla f(x)^T \Delta X_{nsd} \} = \| \nabla f(x) \|_* \quad ⑪$$

$\therefore V^* = \|V\|^* \cdot \hat{V}^* = \| \nabla f(x) \|_* \Delta X_{nsd} = \Delta X_{sd} \quad ⑫ \quad (\text{where we used the definition of } \Delta X_{sd})$

$$\therefore \Delta X_{sd} = V^* \equiv \operatorname{argmin}_V \{ \nabla f(x)^T V + \frac{1}{2} \|V\|^2 \} \quad ⑬$$

2. (AE Exercise 9.1(b))  $\because \gamma > 1$ , i.e.  $x^{(0)} = (\gamma, 1) \in \{(x_1, x_2) \mid |x_2| \leq x_1\}$  ①

Now we want to show  $x_1^{(k)} = \gamma \left(\frac{\gamma-1}{\gamma+1}\right)^k$ ,  $x_2^{(k)} = \left(\frac{1-\gamma}{\gamma+1}\right)^k$  ②. by mathematical induction. Note that  $x_1^{(k)} > |x_2^{(k)}|$  ③ since  $\gamma > 1$ .

Define  $D_1 = \{(x_1, x_2)^T \mid |x_2| \leq x_1\}$  ④. Then when  $x = (x_1, x_2)^T \in D_1$ ,

$$\text{we have: } \frac{\partial f}{\partial x_1} = \frac{1}{2} \{x_1^2 + \gamma x_2^2\}^{-\frac{1}{2}} \cdot 2x_1 = \frac{1}{f(x)} \cdot x_1 \quad ⑤$$

$$\frac{\partial f}{\partial x_2} = \frac{1}{2} \{x_1^2 + \gamma x_2^2\}^{-\frac{1}{2}} \cdot 2\gamma \cdot x_2 = \frac{1}{f(x)} \cdot \gamma x_2 \quad ⑥, \quad \nabla f(x) = \frac{1}{f(x)} \cdot (x_1, \gamma x_2)^T \quad ⑦$$

Now we use mathematical induction:

(i). The base case: when  $k=0$ , we have  $x_1^{(0)} = \gamma \left(\frac{\gamma-1}{\gamma+1}\right)^0 = \gamma$  ⑧

and  $x_2^{(0)} = \left(\frac{1-\gamma}{\gamma+1}\right)^0 = 1$  ⑨ This is indeed our starting point.

(ii). Suppose the  $k$ -th step iterates satisfy  $x_1^{(k)} = \gamma \left(\frac{\gamma-1}{\gamma+1}\right)^k$  ⑩ and

$x_2^{(k)} = \left(\frac{1-\gamma}{\gamma+1}\right)^k$  ⑪, then we want to show  $x_1^{(k+1)} = \gamma \left(\frac{\gamma-1}{\gamma+1}\right)^{k+1}$  and

$x_2^{(k+1)} = \left(\frac{1-\gamma}{\gamma+1}\right)^{k+1}$  ⑫, where  $k \in \mathbb{N}$ . Now given ⑩, ⑪, we have

$$\begin{aligned} \tilde{x}^{(k+1)}(t) &\equiv x^{(k)} - t \nabla f(x^{(k)}) = \left(\gamma \left(\frac{\gamma-1}{\gamma+1}\right)^k, \left(\frac{1-\gamma}{\gamma+1}\right)^k\right)^T - t \cdot \frac{1}{f(x^{(k)})} \left(\gamma \left(\frac{\gamma-1}{\gamma+1}\right)^k, \gamma \left(\frac{1-\gamma}{\gamma+1}\right)^k\right)^T \\ &= \left(1 - \frac{t}{f(x^{(k)})}, \gamma \left(\frac{\gamma-1}{\gamma+1}\right)^k, \left(1 - \frac{t\gamma}{f(x^{(k)})}\right) \cdot \left(\frac{1-\gamma}{\gamma+1}\right)^k\right)^T \quad ⑬ \end{aligned}$$

$$\begin{aligned} \text{Then } \tilde{f}(t) &\equiv f(\tilde{x}^{(k+1)}(t)) = \left\{ \left(1 - \frac{t}{f(x^{(k)})}\right)^2 \gamma^2 \left(\frac{\gamma-1}{\gamma+1}\right)^{2k} + \gamma \left(1 - \frac{t\gamma}{f(x^{(k)})}\right)^2 \left(\frac{1-\gamma}{\gamma+1}\right)^{2k} \right\}^{\frac{1}{2}} \\ &= \left\{ \left[ \gamma^2 \left(\frac{\gamma-1}{\gamma+1}\right)^{2k} + \gamma^3 \left(\frac{1-\gamma}{\gamma+1}\right)^{2k} \right] \frac{t^2}{[f(x^{(k)})]^2} - \frac{2t}{f(x^{(k)})} \left[ \gamma^2 \left(\frac{\gamma-1}{\gamma+1}\right)^{2k} + \gamma^2 \left(\frac{1-\gamma}{\gamma+1}\right)^{2k} \right] + h(\gamma) \right\}^{\frac{1}{2}} \quad ⑭ \end{aligned}$$

where  $h(\gamma)$  is a function of  $\gamma$  only

inside  $\{\gamma\}$  in

Then to minimize the quadratic expression ⑭ of  $t$ , we only need to take

$$t^* = \frac{2 \left[ \gamma^2 \left(\frac{\gamma-1}{\gamma+1}\right)^{2k} + \gamma^3 \left(\frac{1-\gamma}{\gamma+1}\right)^{2k} \right]}{2 \left[ \gamma^2 \left(\frac{\gamma-1}{\gamma+1}\right)^{2k} + \gamma^3 \left(\frac{1-\gamma}{\gamma+1}\right)^{2k} \right]} \cdot f(x^{(k)}) \quad (\text{where we used } x^* = -\frac{b}{2a} = \arg \min_x \{ax^2 + bx + c\} \text{ for } a > 0)$$

$$= \left(\frac{2}{1+\gamma}\right) \cdot f(x^{(k)}) \quad ⑮ \quad \text{Then plugging } t = t^* \text{ into ⑬, we get:}$$

$$x_1^{(k+1)} = \left(1 - \frac{2}{1+\gamma}\right) \cdot \gamma \cdot \left(\frac{\gamma-1}{\gamma+1}\right)^k = \gamma \left(\frac{\gamma-1}{\gamma+1}\right)^{k+1} \quad ⑯$$

$$x_2^{(k+1)} = \left(1 - \gamma \cdot \frac{2}{1+\gamma}\right) \cdot \left(\frac{1-\gamma}{\gamma+1}\right)^k = \left(\frac{1-\gamma}{1+\gamma}\right)^{k+1} \quad ⑰$$

i. We can see that ⑯ and ⑰ do match our target form ⑪.

∴ We have proven the  $k$ -th step iterates  $x_1^{(k)}$  and  $x_2^{(k)}$  satisfy the closed form ② by mathematical induction.

### 3. (Proximal Gradient Descent I)

$$\begin{aligned}
 (a) \text{ We have } \text{prox}_{\alpha, h}(x - \alpha \nabla g(x)) &\equiv \underset{y}{\operatorname{argmin}} \left\{ \frac{1}{2\alpha} \|x - \alpha \nabla g(x) - y\|_2^2 + h(y) \right\} \\
 &= \underset{y}{\operatorname{argmin}} \left\{ \frac{1}{2\alpha} \|(y-x) + \alpha \nabla g(x)\|_2^2 + h(y) \right\} = \underset{y}{\operatorname{argmin}} \left\{ \frac{1}{2\alpha} [(y-x) + \alpha \nabla g(x)]^\top [(y-x) + \alpha \nabla g(x)] + h(y) \right\} \\
 &= \underset{y}{\operatorname{argmin}} \left\{ \frac{1}{2\alpha} \|y-x\|_2^2 + \nabla g(x)^\top (y-x) + \frac{\alpha}{2} \|\nabla g(x)\|_2^2 + h(y) \right\} \\
 &= \underset{y}{\operatorname{argmin}} \left\{ U(x, y) + \left[ \frac{\alpha}{2} \|\nabla g(x)\|_2^2 - g(x) \right] \right\} \quad (\text{Then since } [\frac{\alpha}{2} \|\nabla g(x)\|_2^2 - g(x)] \text{ does not depend on } y, \text{ we can drop it out of the } \underset{y}{\operatorname{argmin}} \{ \dots \}) \\
 &= \underset{y}{\operatorname{argmin}} \{ U(x, y) \} \quad \textcircled{1}
 \end{aligned}$$

$$\begin{aligned}
 (b) \underset{v \in C}{\operatorname{argmin}} \{ g(x) + \langle \nabla g(x), v-x \rangle + \frac{1}{2\alpha} \|v-x\|_2^2 \} &\rightarrow \textcircled{1} \\
 &= \underset{v \in \mathbb{R}^n}{\operatorname{argmin}} \{ g(x) + \langle \nabla g(x), v-x \rangle + \frac{1}{2\alpha} \|v-x\|_2^2 + h(v) \} \quad \text{where } h(v) \equiv \begin{cases} \infty & \text{if } v \notin C \\ 0 & \text{if } v \in C \end{cases} \\
 &\quad (\text{then using the result from part (a)}) \\
 &= \text{prox}_{\alpha, h}(x - \alpha \nabla g(x)) \equiv \underset{v \in \mathbb{R}^n}{\operatorname{argmin}} \left\{ \frac{1}{2\alpha} \|(x - \alpha \nabla g(x)) - v\|_2^2 + h(v) \right\} \quad (\text{then use definition } \textcircled{1}) \\
 &= \underset{v \in C}{\operatorname{argmin}} \left\{ \frac{1}{2\alpha} \|(x - \alpha \nabla g(x)) - v\|_2^2 \right\} = \underset{v \in C}{\operatorname{argmin}} \left\{ \|[(x - \alpha \nabla g(x)) - v]\|_2^2 \right\} \\
 &= \Pi_C(x - \alpha \nabla g(x)) \quad \textcircled{2} \quad \text{because } \Pi_C(u) \equiv \underset{v \in C}{\operatorname{argmin}} \{ \|u-v\|_2^2 \} \quad \textcircled{3} \\
 &= y \quad \textcircled{4}
 \end{aligned}$$

(c) Since  $g(x)$  is  $M$ -smooth, i.e.  $\|\nabla g(v) - \nabla g(y)\|_2 \leq M \|v-y\|_2$  for  $\forall x, y \in \mathbb{R}^n$ , then from lecture notes 16 page 5, we have proven the result:

$$|g(y) - [g(x) + \nabla g(x)^\top (y-x)]| \leq \frac{M}{2} \|x-y\|_2^2 \quad \textcircled{1}$$

$$\therefore g(y) - [g(x) + \nabla g(x)^\top (y-x)] \leq \frac{M}{2} \|x-y\|_2^2$$

$$\begin{aligned}
 \therefore g(y) &\leq g(x) + \nabla g(x)^\top (y-x) + \frac{M}{2} \|x-y\|_2^2 \quad (\text{then use } M \leq \frac{1}{\alpha}) \\
 &\leq g(x) + \nabla g(x)^\top (y-x) + \frac{M}{2} \|y-x\|_2^2 \stackrel{\textcircled{2}}{=} U(x, y) - h(y) \quad \textcircled{3}
 \end{aligned}$$

$$\therefore g(y) \leq U(x, y) - h(y)$$

$$\therefore U(x, y) \geq g(y) + h(y) = f(y) \quad \textcircled{4}$$

(d) For  $\forall x^*$  such that  $G_\alpha(x^*) = 0$ , we have  $x^* = \text{prox}_{\alpha, h}(x^* - \alpha \nabla g(x^*))$  ①  
 $\equiv \arg \min_y \left\{ \frac{1}{2\alpha} \|x^* - \alpha \nabla g(x^*) - y\|_2^2 + h(y) \right\}$  ②

$$\text{Now define } p(y) \equiv \frac{1}{2\alpha} \|x^* - \alpha \nabla g(x^*) - y\|_2^2 + h(y)$$

$$= \frac{1}{2\alpha} (x^* - \alpha \nabla g(x^*) - y)^T (x^* - \alpha \nabla g(x^*) - y) + h(y)$$

$$= \frac{1}{2\alpha} \left\{ y^T y - 2(x^* - \alpha \nabla g(x^*))^T y + \|x^* - \alpha \nabla g(x^*)\|_2^2 \right\} + h(y) \quad ③$$

$$\therefore x^* = \arg \min_y \{ p(y) \} \quad ④ \quad \text{Note that } p(y) \text{ is a differentiable function}$$

since  $h(y)$  is differentiable. Then using the first-order optimality condition on

$$④, \text{ we have } \langle \nabla p(y) \rangle_{y=x^*}, y - x^* \geq 0 \text{ for } \forall y \in \mathbb{R}^n \quad ⑤$$

$$\text{Note that } \nabla p(y)|_{y=x^*} = \left[ \frac{1}{2\alpha} \left\{ 2y - 2(x^* - \alpha \nabla g(x^*)) \right\} + \nabla h(y) \right] |_{y=x^*} \text{ (from ③)}$$

$$= \frac{1}{2\alpha} \left\{ 2x^* - 2x^* + 2\alpha \nabla g(x^*) \right\} + \nabla h(x^*)$$

$$= \nabla g(x^*) + \nabla h(x^*) = \nabla f(x^*) \quad ⑥ \quad (\text{since } f(x) \equiv g(x) + h(x))$$

$$\text{Then plugging ⑥ in ⑤, we have } \langle \nabla f(x^*), y - x^* \rangle \geq 0 \quad ⑦$$

for  $\forall y \in \mathbb{R}^n$ . Then since  $y \in \mathbb{R}^n$ , we have  $y - x^* \in \mathbb{R}^n$  can be any arbitrary vector in  $\mathbb{R}^n$ .

$\therefore$  For ⑦ to hold for  $\forall y \in \mathbb{R}^n$ , we must

have  $\nabla f(x^*) = 0$  ⑧.  $\therefore$  We have proven for  $\forall x^*$  such that  $G_\alpha(x^*) = 0$ ,

we must have  $\nabla f(x^*) = 0$ , i.e.  $x^*$  is a stationary point of  $f$ .

4. (a) We can simply redo the entire procedure as in problem 3 part (d), except that 3(d) ① becomes  $x^{k+1} = \text{prox}_{\alpha, h}\{x^k - \alpha \nabla g(x^k)\}$  ①

$$\text{Then 3(d) ⑥ becomes: } \nabla p(x^{k+1})|_{y=x^{k+1}} = \frac{1}{2\alpha} \left\{ 2x^{k+1} - 2x^k + 2\alpha \nabla g(x^k) \right\} + \nabla h(x^{k+1}) \quad ②$$

Then using the same reasoning as in 3(d), we have:

$$\nabla p(x^{k+1}) = \frac{1}{2} \left\{ x^{k+1} - x^k + \alpha \nabla g(x^k) \right\} + \nabla h(x^{k+1}) = 0 \quad ③$$

$$\therefore \text{Equivalently, we have } \nabla g(x^k) + \frac{1}{2} (x^{k+1} - x^k) + \nabla h(x^{k+1}) = 0 \quad ④$$

$$\therefore \nabla g(x^k) - G_\alpha(x^k) + \nabla h(x^{k+1}) = 0 \quad ⑤$$

$$\begin{aligned} \text{Note that ① trivially comes from } x^{k+1} &= x^k - \alpha G_\alpha(x^k) \\ &= x^k - \alpha \cdot \frac{1}{2} \left\{ x^k - \text{prox}_{\alpha, h}(x^k - \alpha \nabla g(x^k)) \right\} \\ &= -\text{prox}_{\alpha, h}(x^k - \alpha \nabla g(x^k)) \quad ⑥ \end{aligned}$$

(b)  $\because h(x)$  is convex and differentiable,  $\therefore$  By first-order condition for convex functions, we have  $h(x^k) \geq h(x^{k+1}) + \langle \nabla h(x^{k+1}), x^k - x^{k+1} \rangle$  ①

And by definition of  $x^{k+1}$  for proximal gradient updates, we have:

$$x^k - x^{k+1} = \alpha G_\alpha(x^k) \quad \text{②} \quad (\text{see definition (5) in the prompt}) \quad \text{Then plugging}$$

$$\text{② into ①, we get: } h(x^{k+1}) + \alpha \langle \nabla h(x^{k+1}), G_\alpha(x^k) \rangle \leq h(x^k) \quad \text{③}$$

(c) First we prove  $f(x^{k+1}) \leq U(x^k, x^{k+1})$  ① This is very straightforward.

We can simply plug  $x = x^k$  and  $y = x^{k+1}$  into ③(1), then we directly get:  $U(x^k, x^{k+1}) \geq g(x^{k+1}) + h(x^{k+1}) = f(x^{k+1})$  ②

Next, we prove  $U(x^k, x^{k+1}) \leq f(x^k) - \frac{\alpha}{2} \|G_\alpha(x^k)\|_2^2$  ③

Use " $\Leftrightarrow$ " to denote equivalence. Then using the definition of  $U(x, y)$ ,

$$\text{③} \Leftrightarrow g(x^k) + \langle \nabla g(x^k), x^{k+1} - x^k \rangle + \frac{1}{2\alpha} \|x^{k+1} - x^k\|_2^2 + h(x^{k+1}) \leq f(x^k) - \frac{\alpha}{2} \|G_\alpha(x^k)\|_2^2$$

$$\Leftrightarrow g(x^k) - \alpha \langle \nabla g(x^k), G_\alpha(x^k) \rangle + \frac{1}{2\alpha} \|\alpha^2 G_\alpha(x^k)\|_2^2 + h(x^{k+1}) \leq g(x^k) + h(x^k) - \frac{\alpha}{2} \|G_\alpha(x^k)\|_2^2$$

(where we have used  $x^{k+1} = x^k - \alpha G_\alpha(x^k)$  and  $f \equiv g + h$ )

$$\Leftrightarrow -\alpha \langle \nabla g(x^k), G_\alpha(x^k) \rangle + \alpha \|G_\alpha(x^k)\|_2^2 + h(x^{k+1}) \leq h(x^k) \quad \text{④}$$

$$\Leftrightarrow -\alpha \langle \nabla g(x^k), G_\alpha(x^k) \rangle + \alpha \langle G_\alpha(x^k), G_\alpha(x^k) \rangle + h(x^{k+1}) \leq h(x^k)$$

$$\Leftrightarrow \alpha \langle [G_\alpha(x^k) - \nabla g(x^k)], G_\alpha(x^k) \rangle + h(x^{k+1}) \leq h(x^k) \quad \text{⑤}$$

(Then plugging  $\nabla h(x^{k+1}) = G_\alpha(x^k) - \nabla g(x^k)$  from 4.(a) ⑤ into ⑤)

$$\Leftrightarrow \alpha \langle \nabla h(x^{k+1}), G_\alpha(x^k) \rangle + h(x^{k+1}) \leq h(x^k) \quad \text{⑥}$$

Notice that inequality ⑥ is exactly the result of part (b) that we have proven. Then by equivalence, we can immediately see that ③ is true.

Then combining ② and ③, we have:  $f(x^{k+1}) \leq U(x^k, x^{k+1}) \leq f(x^k) - \frac{\alpha}{2} \|G_\alpha(x^k)\|_2^2$  ⑦

$$5. (a) \text{prox}_h(x) \equiv \underset{y}{\operatorname{argmin}} \left\{ \frac{1}{2} \|x-y\|_2^2 + h(y) \right\} \quad (1)$$

$$\begin{aligned} \text{Then define } P(y) &\equiv \frac{1}{2} \|x-y\|_2^2 + h(y) \quad (2), \text{ we have } \nabla P(y) = \nabla_y \left\{ \frac{1}{2} (y-x)^T (y-x) + h(y) \right\} \\ &= \nabla_y \left\{ \frac{1}{2} y^T y - x^T y + \frac{1}{2} x^T x + h(y) \right\} = y - x + \nabla h(y) \quad (3) \end{aligned}$$

Then similar to what we did in problem 3(d), by optimality condition,

$$\text{we have } \nabla P(y) \Big|_{y=\text{prox}_h(x)} = 0 \quad (4) \quad \therefore \text{plugging (3) into (4), we have:}$$

$$\text{prox}_h(x) - x + \nabla h(\text{prox}_h(x)) = 0 \quad (5) \quad \therefore \nabla h(\text{prox}_h(x)) = x - \text{prox}_h(x) \quad (6)$$

Now using the result of BV Exercise 3.11, which is homework 2 Problem 6,

$$\text{we know that } \langle \nabla h(x) - \nabla h(y), x - y \rangle \geq 0 \quad (7) \text{ for } \forall x, y \in \mathbb{R}^n \text{ since}$$

$h$  is a differentiable convex function. Then plugging  $x = \text{prox}_h(x_1)$ ,

$$y = \text{prox}_h(x_2) \text{ into (7), we have: } \langle \nabla h(\text{prox}_h(x_1)) - \nabla h(\text{prox}_h(x_2)),$$

$$\text{prox}_h(x_1) - \text{prox}_h(x_2) \rangle \geq 0 \quad (8) \quad \text{Then plugging (6) into (8), we}$$

$$\text{have: } \langle [x_1 - \text{prox}_h(x_1)] - [x_2 - \text{prox}_h(x_2)], \text{prox}_h(x_1) - \text{prox}_h(x_2) \rangle \geq 0 \quad (9)$$

$$\therefore \langle (x_1 - x_2) - [\text{prox}_h(x_1) - \text{prox}_h(x_2)], \text{prox}_h(x_1) - \text{prox}_h(x_2) \rangle \geq 0$$

$$\therefore \langle x_1 - x_2, \text{prox}_h(x_1) - \text{prox}_h(x_2) \rangle \geq \langle \text{prox}_h(x_1) - \text{prox}_h(x_2), \text{prox}_h(x_1) \rangle$$

$$\therefore \langle \text{prox}_h(x_1) - \text{prox}_h(x_2), x_1 - x_2 \rangle \geq \|\text{prox}_h(x_1) - \text{prox}_h(x_2)\|_2^2 \quad (10) \quad \text{--- prox}_h(x_2) \rangle$$

(b) Using Cauchy-Schwarz inequality, we have:

$$\langle \text{prox}_h(x_1) - \text{prox}_h(x_2), x_1 - x_2 \rangle \leq \|\text{prox}_h(x_1) - \text{prox}_h(x_2)\|_2 \cdot \|x_1 - x_2\|_2 \quad (1)$$

Then combining (1) and 5.(a) (10), we have:

$$\|\text{prox}_h(x_1) - \text{prox}_h(x_2)\|_2 \cdot \|x_1 - x_2\|_2 \geq \|\text{prox}_h(x_1) - \text{prox}_h(x_2)\|_2^2 \quad (2)$$

$$\therefore \|x_1 - x_2\|_2 \geq \|\text{prox}_h(x_1) - \text{prox}_h(x_2)\|_2 \quad (3)$$

Note that if  $\|\text{prox}_h(x_1) - \text{prox}_h(x_2)\|_2 = 0$ , (3) will trivially hold.

If  $\|\text{prox}_h(x_1) - \text{prox}_h(x_2)\|_2 \neq 0$ , then we can simply divide both sides of (2) by this term to get (3).

$$6. (a) V' \equiv \Pi_C(v) \equiv \underset{u \in C}{\operatorname{argmin}} \|v - u\|_2 = \underset{u \in C}{\operatorname{argmin}} \|v - u\|_2^2 \quad \textcircled{1}$$

Then if  $s < n$ , we know for  $u \in C \equiv \{x \mid \|x\|_0 \leq s\}$ , there must be at least  $(n-s)$  zero entries in  $u$ . Then for  $\forall u \in C$ , let set  $J$  contains the indices of  $(n-s)$  zero entries in  $u$ . If there are more than  $(n-s)$  zero entries in  $u$ , we can pick the indices of  $(n-s)$  of them arbitrarily.

$$\text{Then } \|v - u\|_2^2 = \sum_{i=1}^n (v_i - u_i)^2 \geq \sum_{i \in J} (v_i - 0)^2 = \sum_{i \in J} |v_i|^2 \geq \sum_{i \notin V} |v_i|^2 \quad \textcircled{2}$$

where the last inequality comes from the definition of  $V$  such that  $V$  contains the indices of the largest  $s$  entries in absolute value of  $v$ . And notice that when  $u \in C$  satisfies  $u_j = \begin{cases} v_j, & \text{if } j \in V \\ 0, & \text{if } j \notin V \end{cases} \quad \textcircled{3}$ , we have  $\|v - u\|_2^2 = \sum_{i \notin V} |v_i|^2 \quad \textcircled{4}$

Then combining  $\textcircled{2}$  and  $\textcircled{4}$ , we know that

$$\min_{u \in C} \|v - u\|_2^2 = \sum_{i \notin V} |v_i|^2 \quad \textcircled{2}, \text{ and } \Pi_C(v) \equiv \underset{u \in C}{\operatorname{argmin}} \|v - u\|_2^2 = V' \text{ such that}$$

$$V'_j = \begin{cases} v_j, & \text{if } j \in V \\ 0, & \text{if } j \notin V \end{cases} \quad \textcircled{5}$$

If  $s \geq n$ , then  $V' \equiv \Pi_C(v) = V \quad \textcircled{6}$  because  $V \in C \equiv \{x \mid \|x\|_0 \leq s\}$  and  $\|v - V\|_2^2 = 0$ , which is the least possible value for  $L_2$  norm. In this case, since  $s \geq n$ , we have  $V = \{1, 2, 3, \dots, n\}$ . Then  $\textcircled{6}$  will still hold, trivially.

(b) See notebook implementation (attached in next page)

(c) See notebook output cells (attached in next page)