

Name: Zheng Mingren

Stu no: 21126390

1. To solve $\nabla f(x_1, x_2) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \vec{0}$, we just need to solve $\begin{cases} \frac{\partial f}{\partial x_1} = 0 \\ \frac{\partial f}{\partial x_2} = 0 \end{cases}$

then all the solutions (x_1, x_2) is what we want to find.

$$(a) \begin{cases} \frac{\partial f}{\partial x_1} = 2(4x_1^2 - x_2) \cdot (8x_1) = 0 & ① \\ \frac{\partial f}{\partial x_2} = 2(4x_1^2 - x_2) \cdot (-1) = 0 & ② \end{cases}$$

$$① \Leftrightarrow x_2 = 4x_1^2 \text{ or } x_1 = 0$$

$$② \Leftrightarrow x_2 = 4x_1^2, \text{ for } \forall x_1 \in \mathbb{R} (\text{which certainly contains } x_1=0)$$

when $x_2 = 4x_1^2 \checkmark$, ① is always true. i.e. \checkmark

those vectors in the form of $(a, 4a^2)$ for $\forall a \in \mathbb{R}$ are all critical points.

$$(b) \begin{cases} \frac{\partial f}{\partial x_1} = 6x_1x_2 = 0 & ① \\ \frac{\partial f}{\partial x_2} = 6x_2^2 - 12x_2 + 3x_1^2 = 0 & ② \end{cases}$$

$$① \Leftrightarrow x_1 = 0 \text{ or } x_2 = 0$$

\checkmark when $x_1 = 0$, ② can be reformulated as $x_2^2 = 2x_2$,

solve it then we have $x_2 = 0$ or $x_2 = 2$.

i.e. $(0, 0), (0, 2)$ are critical points. \checkmark

\checkmark when $x_2 = 0$, ② can be reformulated as $3x_1^2 = 0$

solve it then we have $x_1 = 0$, i.e. $(0, 0)$ is a critical point.

Altogether, critical points are $(0, 0)$ and $(0, 2)$

$$(c) \begin{cases} \frac{\partial f}{\partial x_1} = 4(x_1 - 2x_2)^3 + 64x_2 = 0 & ① \\ \frac{\partial f}{\partial x_2} = 4(x_1 - 2x_2)^3 \cdot (-2) + 64x_1 = 0 & ② \end{cases}$$

$$\begin{cases} \frac{\partial f}{\partial x_1} = 4(x_1 - 2x_2)^3 + 64x_2 = 0 & ① \\ \frac{\partial f}{\partial x_2} = 4(x_1 - 2x_2)^3 \cdot (-2) + 64x_1 = 0 & ② \end{cases}$$

(in next page)

$$\textcircled{1} \Leftrightarrow 16x_2 = -(x_1 - 2x_2)^3$$

$$\textcircled{2} \Leftrightarrow -8x_1 = -(x_1 - 2x_2)^3$$

thus $16x_2 = -8x_1$ i.e. $x_1 = -2x_2$

Substitute x_1 in \textcircled{1}, then we have:

$$4(-4x_2)^3 + 64x_2 = 0$$

$$64x_2 = 4 \cdot 64x_2^3$$

$$\Rightarrow x_2 = 0 \quad \text{or} \quad x_2 = \frac{1}{2} \quad \text{or} \quad x_2 = -\frac{1}{2}$$

$$\begin{aligned} x_1 &= -2x_2 \\ \Rightarrow x_1 &= 0 \quad \text{or} \quad x_1 = -1 \quad \text{or} \quad x_1 = 1 \end{aligned}$$

i.e. $(0, 0)$, $(-1, \frac{1}{2})$, $(1, -\frac{1}{2})$ are all the critical points.

$$(d) \left\{ \begin{array}{l} \frac{\partial f}{\partial x_1} = 2x_1 + 4x_2 + 1 = 0 \quad \textcircled{1} \\ \frac{\partial f}{\partial x_2} = 4x_1 + 2x_2 - 1 = 0 \quad \textcircled{2} \end{array} \right.$$

$$\Rightarrow \begin{cases} x_1 = \frac{1}{2} \\ x_2 = -\frac{1}{2} \end{cases}$$

i.e. $(\frac{1}{2}, -\frac{1}{2})$ is the only critical point.

(in next page)

2.

$$(a) \checkmark \text{let } g(x) = \frac{1}{2} \|Ax - b\|_2^2$$

$$\begin{aligned} g(y) &= \frac{1}{2} \|Ay - b\|_2^2 = \frac{1}{2} \|A(x + y - x) - b\|_2^2 \\ &= \frac{1}{2} \|Ax - b\|_2^2 + \langle Ax - b, Ay - Ax \rangle + \frac{1}{2} \|Ay - Ax\|_2^2 \end{aligned}$$

$$\begin{aligned} &\langle u, Av \rangle \\ &= u^T Av \\ &= (A^T u)^T v \\ &= \langle A^T u, v \rangle \end{aligned}$$

$$\begin{aligned} &= g(x) + \langle Ax - b, Ay - Ax \rangle + \frac{1}{2} \|Ay - Ax\|_2^2 \\ &= g(x) + \underbrace{\langle A^T(Ax - b), y - x \rangle}_{\text{(*)}} + \frac{1}{2} \|Ay - Ax\|_2^2 \end{aligned}$$

$\lim_{\|y-x\|_2 \rightarrow 0} \frac{|g(y) - (g(x) + \langle A^T(Ax - b), y - x \rangle)|}{\|y-x\|_2}$

$$\stackrel{(*)}{=} \lim_{\|y-x\|_2 \rightarrow 0} \frac{\frac{1}{2} \|Ay - Ax\|_2^2}{\|y-x\|_2}$$

by Cauchy-Schwarz inequality

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

$$\Rightarrow \nabla g(x) = \boxed{A^T(Ax - b)}$$



$$\checkmark \text{let } h(x) = \lambda \|x\|_2^2$$

$$\nabla h(x) = \lambda \cdot (2x) = 2\lambda x$$

$$\begin{aligned} \text{so, } \nabla f(x) &= \nabla(g(x) + h(x)) = \nabla g(x) + \nabla h(x) \\ &= A^T(Ax - b) + 2\lambda x \end{aligned}$$

(b) Assume the (i,j) -entry of X is x_{ij} , denote i -th entry of c as c_i

$$\begin{aligned} f(X) &= b^T X c = \langle b, Xc \rangle = \sum_{i=1}^n \left(b_i \cdot \sum_{j=1}^n x_{ij} \cdot c_j \right) \\ &\quad \text{---} \rightarrow (Xc)_i \end{aligned}$$

$$= \sum_{i=1}^n \sum_{j=1}^n b_i x_{ij} c_j$$

$\frac{\partial f(X)}{\partial x_{ij}} = b_i \cdot c_j$, which is exactly the (i,j) -entry of the matrix bc^T

i.e. $\frac{\partial f(X)}{\partial x_{ij}} = b_i c_j = (bc^T)_{ij}$

as $\frac{\partial f(X)}{\partial X} = \left[\frac{\partial f(X)}{\partial x_{ij}} \right]$,

we have $\frac{\partial f(X)}{\partial X} = \left[\frac{\partial f(X)}{\partial x_{ij}} \right] = \left[(bc^T)_{ij} \right] = bc^T$ ✓

($A = [a_{ij}]$ means the (i,j) -entry of A is a_{ij})

(c) $f(X) = b^T X^T X c = (Xb)^T (Xc)$

$= \langle Xb, Xc \rangle$

$$= \sum_{i=1}^n \left(\left(\sum_{j=1}^n x_{ij} \cdot b_j \right) \cdot \left(\sum_{j=1}^n x_{ij} \cdot c_j \right) \right)$$

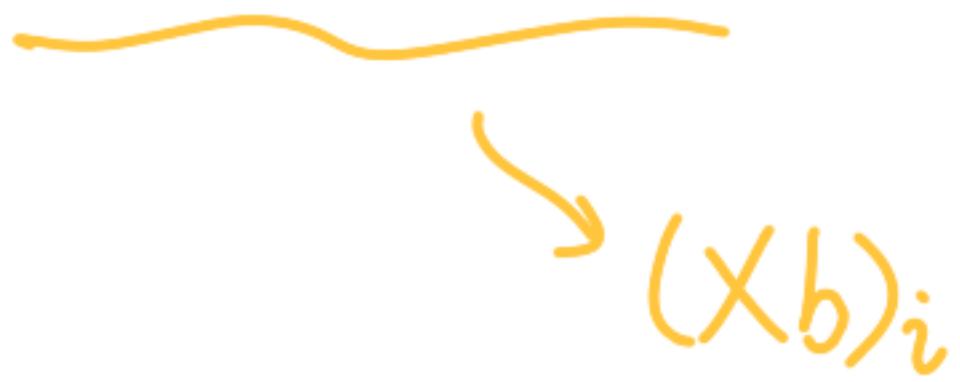
(Xb)_i (Xc)_i

Then $\frac{\partial f(X)}{\partial x_{ij}} = \frac{\partial (Xb)_i}{\partial x_{ij}} \cdot (Xc)_i + (Xb)_i \cdot \frac{\partial (Xc)_i}{\partial x_{ij}}$

$$= b_j \cdot \underbrace{\left(\sum_{j=1}^n x_{ij} \cdot c_j \right)}_{\text{denote as } I_1} + c_j \cdot \underbrace{\left(\sum_{j=1}^n x_{ij} \cdot b_j \right)}_{\text{denote as } I_2}$$

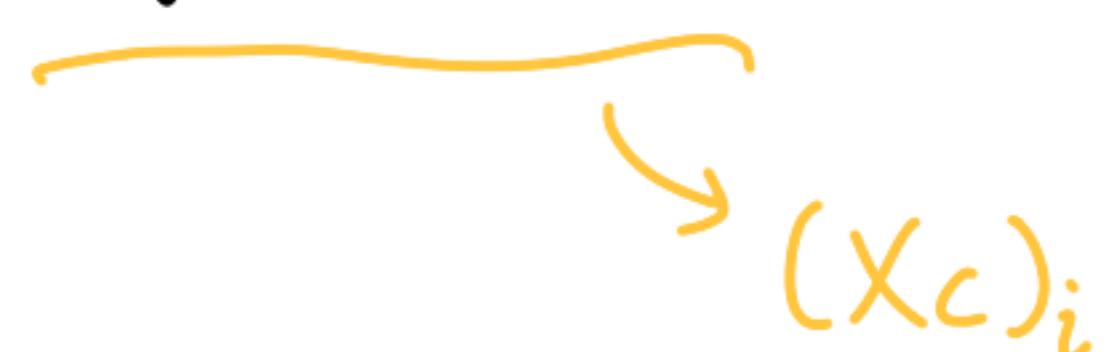
Note that

$$(Xbc^\top)_{ij} = \left(\sum_{j=1}^n x_{ij} \cdot b_j \right) \cdot c_j = I_1,$$



$$(Xb)_i$$

$$(Xcb^\top)_{ij} = \left(\sum_{j=1}^n x_{ij} \cdot c_j \right) \cdot b_j = I_2$$



$$(Xc)_i$$

Then we have

$$\begin{aligned}\frac{\partial f(x)}{\partial x_{ij}} &= I_1 + I_2 = (Xbc^\top)_{ij} + (Xcb^\top)_{ij} \\ &= (Xbc^\top + Xcb^\top)_{ij} \\ &= \left(X(bc^\top + cb^\top) \right)_{ij}\end{aligned}$$

So

$$\begin{aligned}\frac{\partial f(x)}{\partial X} &= \left[\frac{\partial f(x)}{\partial x_{ij}} \right] = \left[\left(X(bc^\top + cb^\top) \right)_{ij} \right] \\ &= X(bc^\top + cb^\top)\end{aligned}$$

✓

3.

(a)

Proof:

Part 1: Prove: For $\forall b \in R^n$, b can be decomposed as

$$b = b_s + \sum_{i=1}^N c_i x_i$$

where $c = \begin{bmatrix} c_1 \\ \vdots \\ c_N \end{bmatrix} \in R^N$ and $b_s \in R^n$ satisfying $\langle b_s, x_i \rangle = 0$ for $\forall i = 1, \dots, N$.

Proof: For simplicity, we prove only the case $N=1$.

$$S = \{ v \in R^n \mid \langle v, x_1 \rangle = 0 \}$$

Then S is a hyperplane and a subspace.

so b can be decomposed as

$$b = P_S b + (b - P_S b)$$

For $P_S b$: Since $P_S b$ is the projection,

$$\left. \begin{array}{l} b - P_S b \perp S \\ P_S b \in S \end{array} \right\} \Rightarrow \langle b - P_S b, P_S b \rangle = 0$$

For $b - P_S b$, By the explicit formula of projection,

$$b - P_S b = \frac{\langle x_1, b \rangle}{\|x_1\|^2} \cdot x_1 \equiv c_1 x_1$$

Define $P_S b = b_s$, then

$$b = b_s + c_1 x_1$$

$$\text{and } \langle b - P_S b, P_S b \rangle = \langle c_1 x_1, b_s \rangle = 0$$

$\Rightarrow \langle x_1, b_s \rangle = 0$. Proved.

the conclusion can be generalized to $N=2, 3, \dots$ similarly.

Part 2: Then the function can be rewritten as:

$$\sum_{i=1}^N (\langle a, x_i \rangle - y_i)^2 + \lambda \|a\|_2^2$$

$$= \sum_{i=1}^N \left(\langle a_s + \sum_{j=1}^N c_j x_j, x_i \rangle - y_i \right)^2 + \lambda \left\| a_s + \sum_{j=1}^N c_j x_j \right\|_2^2$$

$$= \sum_{i=1}^N \left(\cancel{\langle a_s, x_i \rangle} + \langle \sum_{j=1}^N c_j x_j, x_i \rangle - y_i \right)^2 \\ + \lambda \left(\|a_s\|_2^2 + 2 \sum_{j=1}^N c_j \cancel{\langle a_s, x_j \rangle} + \sum_{j=1}^N \sum_{i=1}^N c_i c_j \langle x_i, x_j \rangle \right)$$

$\cancel{\langle a_s, x_i \rangle} = 0$
for $i=1, \dots, N$

$$= \sum_{i=1}^N \left(\sum_{j=1}^N c_j \langle x_j, x_i \rangle - y_i \right)^2 + \lambda \sum_{j=1}^N \sum_{i=1}^N c_i c_j \langle x_i, x_j \rangle$$

$$+ \lambda \|a_s\|_2^2$$

Define $K = [\langle x_i, x_j \rangle]_{i,j=1}^N$

$$\sum_{i=1}^N \left(\sum_{j=1}^N K_{ji} c_j - y_i \right)^2 + \lambda \sum_{i=1}^N \sum_{j=1}^N c_i c_j K_{ij} + \lambda \|a_s\|_2^2$$

$$= \sum_{i=1}^N \left((K^T c)_i - y_i \right)^2 + \lambda c^T K c + \lambda \|a_s\|_2^2$$

$$= \|K^T c - y\|_2^2 + \lambda c^T K c + \lambda \|a_s\|_2^2$$

$\underbrace{\quad}_{\text{denote as}} \quad \underbrace{\quad}_{\text{denote as}}$

$F(c) \quad \quad \quad G(a_s)$

Then,

the original optimization $\Leftrightarrow \begin{cases} \min_{c \in \mathbb{R}^n} F(c) + \lambda G(a_s) \\ a_s \in \mathbb{R}^n \end{cases}$

S.t. $\langle a_s, x_i \rangle = 0$ for $i=1, \dots, N$.



$$\Leftrightarrow \min_{c \in \mathbb{R}^n} F(c) \quad ①$$

and

$$\begin{cases} \min_{a_s \in \mathbb{R}^n} \lambda G(a_s) \\ \text{s.t. } \langle a_s, x_i \rangle = 0 \text{ for } i=1, \dots, N. \end{cases} \quad ②$$

To solve ② : Because $\lambda > 0$,

(in next page)

$$\begin{aligned} \textcircled{2} \iff & \left\{ \begin{array}{l} \min_{\alpha_s \in \mathbb{R}^n} \|\alpha_s\|_2^2 \\ \text{s.t. } \langle \alpha_s, x_i \rangle = 0, \quad i=1, \dots, N \end{array} \right. \\ \iff & \alpha_s^* = 0. \end{aligned}$$

Let α^* be a solution of the original optimization,
then

$$\underline{\alpha^*} = \alpha_s^* + \sum_{i=1}^N c_i^* x_i$$

$$= \sum_{i=1}^N c_i^* x_i,$$

where $c^* = \begin{bmatrix} c_1^* \\ \vdots \\ c_N^* \end{bmatrix}$ is the solution of ①.

proved.

(b) In part (a) we have already obtained the re-expression in the process of proof, as (\star) above:

$$\min_{\alpha \in \mathbb{R}^n} \sum_{i=1}^N (\langle \alpha, x_i \rangle - y_i)^2 + \lambda \|\alpha\|_2^2$$

$$\begin{aligned} \Leftrightarrow \min_{c \in \mathbb{R}^N} & \sum_{i=1}^N (\|K^T c - y\|_2^2 + \lambda c^T K c) + \lambda \|\alpha_s\|_2^2 (\star) \\ \alpha_s \in \mathbb{R}^n \\ \langle \alpha_s, x_i \rangle = 0, \quad i=1, \dots, N \end{aligned}$$

$$\text{where } K = [\langle x_i, x_j \rangle]_{i,j=1}^N$$

As we have proved in (a) : the solution

To solve :

$$\min_{\alpha_s \in \mathbb{R}^n} \lambda (\|\alpha_s\|_2^2) \text{ we have } \checkmark \quad \alpha_s^* = 0 \\ \langle \alpha_s, x_i \rangle = 0, i=1, \dots, N \quad , \text{ such that } \lambda \|\alpha_s^*\|_2^2 = 0 \quad (*-2)$$

Thus

$$(*) \stackrel{(*)-2}{\Leftrightarrow} \left| \min_{C \in \mathbb{R}^N} \sum_{i=1}^N \left(\|K^T c - y\|_2^2 + \lambda c^T K c \right) \right| \quad \checkmark \\ \text{where } K = [\langle x_i, x_j \rangle]_{i,j=1}^N.$$

is the final re-expression we want to find.

(in next page)

4. (a) Proof:

① First we get the gradient of $y = \mathbf{x}^T \mathbf{A} \mathbf{x}$, where $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^n$, $\mathbf{A} = [a_{ij}] \in \mathbb{R}^{n \times n}$.

$$\nabla y = \frac{\partial y}{\partial \mathbf{x}} = \left[\frac{\partial y}{\partial x_1}, \frac{\partial y}{\partial x_2}, \dots, \frac{\partial y}{\partial x_n} \right]^T,$$

$$y = \mathbf{x}^T \mathbf{A} \mathbf{x} = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j.$$

$$\text{then } \frac{\partial y}{\partial x_k} = \frac{\partial}{\partial x_k} \left(\sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j \right)$$

$$= \sum_{i \neq k} \left(\frac{\partial}{\partial x_k} \left(\sum_{j=1}^n a_{ij} x_i x_j \right) \right) + \sum_{j=1}^n a_{kj} x_k x_j$$

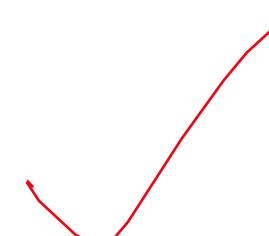
$$= \sum_{i \neq k} \left(\underbrace{\frac{\partial}{\partial x_k} \left(\sum_{j \neq k} a_{ij} x_i x_j \right)}_{=0 \text{ since } i \neq k, j \neq k} + \frac{\partial}{\partial x_k} a_{ik} x_i x_k \right) + \sum_{j \neq k} \frac{\partial}{\partial k} a_{kj} x_k x_j + \frac{\partial}{\partial k} a_{kk} x_k^2$$

$$= \underbrace{\sum_{i \neq k} \frac{\partial}{\partial x_k} a_{ik} x_i x_k}_{= \sum_{i=1}^n a_{ik} x_i} + \underbrace{\sum_{j \neq k} \frac{\partial}{\partial k} a_{kj} x_k x_j}_{= \sum_{j=1}^n a_{kj} x_j} + \frac{\partial}{\partial k} a_{kk} x_k^2$$

$$= \sum_{i=1}^n a_{ik} x_i + \sum_{j=1}^n a_{kj} x_j$$

$$= (\mathbf{A}\mathbf{x})_k + (\mathbf{A}^T \mathbf{x})_k$$

\downarrow k-th entry of $\mathbf{A}\mathbf{x}$



Since \mathbf{A} is symmetric, $\mathbf{A} = \mathbf{A}^T$

$$\text{thus } \frac{\partial y}{\partial x_k} = (2\mathbf{A}\mathbf{x})_k$$

$$\text{thus } \boxed{\Delta y = 2\mathbf{A}\mathbf{x}}$$

② Then we prove: if $\mathbf{A}^{n \times n}$ is symmetric, then $\boxed{d^T \mathbf{A} e = e^T \mathbf{A} d}$ for $\forall \mathbf{e}, \mathbf{d} \in \mathbb{R}^n$

As $d^T \mathbf{A} e$ is a scalar, $d^T \mathbf{A} e = (d^T \mathbf{A} e)^T$

$$\begin{aligned} &= e^T \mathbf{A}^T d \\ &\quad (\text{as } \mathbf{A} \text{ is symmetric}) \\ &= e^T \mathbf{A} d. \text{ proved.} \end{aligned}$$

③ Then we prove: $f(\mathbf{x})$ is differentiable and convex.

③-1

③-2

③-1 $\nabla f(x) \xrightarrow{\text{as we have proved in ①}} 2Ax + 2b$, so $f(x)$ is differentiable at $\forall x \in \mathbb{R}^n$.

③-2 let $g(x) = x^T Ax$ and $h(x) = 2b^T x + C$
then $f(x) = g(x) + h(x)$.

[for $h(x)$:

$$h(y) = 2b^T y + C$$

$$= 2b^T(x + y - x) + C$$

$$= 2b^T x + 2b^T(y - x) + C$$

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

so, $h(x)$ is (differentiable and) convex.

Standard
 $\langle \cdot, \cdot \rangle$ is induced norm on \mathbb{R}^n .

[for $g(x)$:

$$g(y) = y^T A y$$

$$= (x + (y - x))^T A (x + (y - x))$$

$$= (x^T + (y - x)^T) A (x + (y - x))$$

$$= x^T A x + \underbrace{x^T A (y - x)}_{\text{equal to each other as we have proved in ②.}} + \underbrace{(y - x)^T A x}_{\text{equal to each other as we have proved in ②.}} + (y - x)^T A (y - x)$$

equal to each other as we have proved in ②.

$$= x^T A x + \underbrace{2x^T A (y - x)}_{\geq 0 \text{ as } A \text{ is spsd.}} + (y - x)^T A (y - x)$$

$$= x^T A x + \langle 2Ax, y - x \rangle + \underbrace{(y - x)^T A (y - x)}_{\geq 0 \text{ as } A \text{ is spsd.}}$$

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

so, $g(x)$ is (differentiable and) convex.

[Altogether, for $f(x)$:

$$\left\{ \begin{array}{l} h(y) = h(x) + \langle \nabla h(x), y - x \rangle \quad (*-1) \\ g(y) \geq g(x) + \langle \nabla g(x), y - x \rangle \quad (*-2) \end{array} \right.$$

(*-1) + (*-2), we have:

(in next page)

$$\underline{h(y) + g(y)} \geq \underline{h(x) + g(x)} + \underbrace{\langle \nabla h(x) + \nabla g(x), y-x \rangle}_{\checkmark}$$

i.e. $\underline{f(y)} \geq \underline{f(x)} + \underbrace{\langle \nabla f(x), y-x \rangle}_{\checkmark}$

i.e. $f(x)$ is (differentiable and) convex.

④ According to the theorem we discussed in the class,
since $f(x)$ is differentiable and convex,

the global minimizer $x^{(*)} = \arg \min_{x \in \mathbb{R}^n} f(x)$



$$\nabla f(x^{(*)}) = 0 \quad \checkmark$$

$$\nabla f(x) = 2Ax + 2b = 0 \Leftrightarrow Ax = -b$$

i.e. $x = x^{(*)} \Leftrightarrow Ax = -b$. proved. \blacksquare

(in next page)

4(b).

Notations: $\text{Col}(A) = \{Ay \mid \forall y \in \mathbb{R}^n\}$ (column space of $A_{n \times n}$)

$\text{Nul}(A) = \{x \mid Ax = 0\}$

$\text{Nul}(A)^\perp = \{y \mid y \perp \text{Nul}(A)\} \text{ i.e. for } \forall x \in \text{Nul}(A), \langle x, y \rangle = 0$

Part 0: Prove: $\text{Nul}(A^T)^\perp = \text{Col}(A)$ (★)

Proof: $\forall x \in \mathbb{R}^n$, $\forall w \in \text{Nul}(A^T)$,
 $\underbrace{\langle Ax, w \rangle}_{\in \text{Col}(A)} = 0$ since $w \in \text{Nul}(A^T)$

thus $\text{Nul}(A^T)^\perp = \text{Col}(A)$

Part 1: Prove: f is bounded below over $\mathbb{R}^n \Rightarrow b \in \{Ay : y \in \mathbb{R}^n\}$.
(i.e.) $b \in \text{Col}(A)$

Proof: We prove it by proving its converse-negative proposition.

We assume $b \notin \text{Col}(A)$, according to (★).

we have $\underbrace{b \notin \text{Nul}(A^T)^\perp}$

i.e. for $\forall x \in \text{Nul}(A^T)$, $\langle b, x \rangle \neq 0$

since $\exists b \notin \text{Col}(A)$ we can infer that $\dim(A) < n$.

as we have assumed

$\dim(A) < n \Rightarrow \dim(\text{Col}(A)) < n \Leftrightarrow \dim(\text{Nul}(A^T)^\perp) < n$

$\Rightarrow \dim(\text{Nul}(A^T)) = n - \dim(\text{Nul}(A^T)^\perp) > 0$.

so there must be non-zero vector in $\text{Nul}(A^T)$

thus we can find $x_1 \in \text{Nul}(A^T)$, $x_1 \neq 0$

such that $\langle b, x_1 \rangle = b^T x_1 \neq 0$.

Since $x_1 \in \text{Nul}(A^T)$, we have $A^T x_1 = A x_1 = 0$

Then $f(-t \cdot \text{sgn}(b^T x_1) \cdot x_1)$ denote as y

substitute x with $x = -t \cdot \text{sgn}(b^T x_1) \cdot x_1$, and then we finally get:

$$\begin{aligned} &= -t \cdot \text{sgn}(b^T x_1) \cdot y A x_1 + 2b^T y + c \\ &= -2t \cdot \text{sgn}(b^T x_1) \cdot b^T x_1 + c \quad \text{denote as } F \end{aligned}$$

where $\text{sgn}(x) = \begin{cases} 1, & x > 0 \\ -1, & x < 0 \end{cases}$

$F \rightarrow -\infty$ as $t \rightarrow +\infty$, so we proved:

$b \notin \{Ay : y \in \mathbb{R}^n\} \Rightarrow f$ is never bounded below over \mathbb{R}^n

so we can conclude that:

f is bounded below over $\mathbb{R}^n \Rightarrow b \in \{Ay : y \in \mathbb{R}^n\}$ ✓

Part 2: Prove: $b \in \{Ay : y \in \mathbb{R}^n\} \Rightarrow f$ is bounded below over \mathbb{R}^n .

Proof: Since $b \in \{Ay : y \in \mathbb{R}^n\}$

$\exists y \in \mathbb{R}^n$ such that $b = Ay$ (y is constant)

$$\begin{aligned} \text{then } f(x) &= x^T A x + 2b^T x + c = x^T A x + 2x^T A y + c \\ &\quad \left(\begin{array}{l} x^T A y = y^T A x \\ \text{since } A \text{ is symmetric} \end{array} \right) = x^T A x + x^T A y + y^T A x + c \\ &= x^T A x + x^T A y + y^T A x + y^T A y - y^T A y + c \\ &= (x+y)^T A (x+y) - y^T A y + c \\ &\quad \text{---} \geq 0 \text{ since } A \text{ is spsd.} \\ &\geq -y^T A y + c \quad \checkmark \end{aligned}$$

thus, $-y^T A y + c$ is the lower boundary of f .

i.e. we ^{have} proved:

$b \in \{Ay : y \in \mathbb{R}^n\} \Rightarrow f$ is bounded below over \mathbb{R}^n .

Altogether,

$b \in \{Ay : y \in \mathbb{R}^n\} \Leftrightarrow f$ is bounded below over \mathbb{R}^n .

(in next page)

5(a).

let $g(x) = \log(x)$, $x > 0$

$$h(x) = \sum_{i=1}^m \exp(a_i^T x + b_i), \quad x \in \mathbb{R}^n.$$

$$g'(x) = \frac{1}{x}.$$

$$\nabla h(x) = \sum_{i=1}^m \nabla(\exp(a_i^T x + b_i))$$

$$= \sum_{i=1}^m [\exp(a_i^T x + b_i) \cdot \nabla(a_i^T x + b_i)]$$

$$= \sum_{i=1}^m [\exp(a_i^T x + b_i) \cdot a_i]$$

$$\text{thus } \nabla f(x) = \nabla(g(h(x)))$$

$$= g'(h(x)) \cdot \nabla h(x)$$

$$= \frac{\sum_{i=1}^m [\exp(a_i^T x + b_i) \cdot a_i]}{\sum_{i=1}^m \exp(a_i^T x + b_i)}$$



(in next page)

5(b) Yes, it will.

Proof: [general idea: First we prove $f(x)$ is convex, then we further conclude the convergence in GD process.]

[note: $f(x) = \log(g(x))$, $x \in \mathbb{R}^n$, $g(x): \mathbb{R}^n \rightarrow \mathbb{R}$, $g(x) > 0$ for $\forall x \in \mathbb{R}^n$. If $f(x)$ is convex, we call $g(x)$ is log-convex.]

Part 1: Prove: for $\forall g(x): \mathbb{R}^n \rightarrow \mathbb{R}$ which satisfies $g(x) > 0$ for $\forall x \in \mathbb{R}^n$, we have:
 $g(x)$ is log-convex $\iff g(\theta x + (1-\theta)y) \leq [g(x)]^\theta [g(y)]^{1-\theta}$
 $\forall \theta \in [0,1]$, $\forall x, y \in \mathbb{R}^n$.

proof: for $\forall g(x): \mathbb{R}^n \rightarrow \mathbb{R}$ and $g(x) > 0$ for $\forall x \in \mathbb{R}^n$,
 $g(x)$ is log-convex

$\iff f(x) = \log(g(x))$ is convex

by the definition of convexity
 $\iff f(\theta x + (1-\theta)y) \leq \theta f(x) + (1-\theta)f(y)$, $\forall x, y \in \mathbb{R}^n, \forall \theta \in [0,1]$

substitute f with $\log(g)$

$\iff \log(g(\theta x + (1-\theta)y)) \leq \theta \log(g(x)) + (1-\theta) \log(g(y)),$ (*)
 $\forall x, y \in \mathbb{R}^n, \forall \theta \in [0,1]$

take the logarithm on both sides (backward)

$\iff g(\theta x + (1-\theta)y) \leq \exp\left(\theta \log(g(x)) + (1-\theta) \log(g(y))\right),$
 $\forall x, y \in \mathbb{R}^n, \forall \theta \in [0,1]$

take the exponential on both sides (forward)

$\iff g(\theta x + (1-\theta)y) \leq \exp(\theta \log(g(x))) \cdot \exp((1-\theta) \log(g(y))),$
 $\forall x, y \in \mathbb{R}^n, \forall \theta \in [0,1]$

$\iff g(\theta x + (1-\theta)y) \leq [g(x)]^\theta \cdot [g(y)]^{1-\theta}, \forall x, y \in \mathbb{R}^n, \forall \theta \in [0,1].$

proved.

Part 2: Prove: $h(x) = \exp(a^T x + b_i)$, $a, x \in \mathbb{R}^n$, $b_i \in \mathbb{R}$,
is log-convex

Proof: $h(x) > 0$ for $\forall x \in \mathbb{R}^n$ since exp functions are always positive.

$$h(\theta x + (1-\theta)y) = \exp(a^T(\theta x + (1-\theta)y) + b_i)$$

$$= \exp(\underline{\theta a^T x} + \underline{(1-\theta)a^T y} + b_i)$$

(in next page)

$$\begin{aligned}
&= \exp(\underline{\theta(\alpha^T x + b_i)} + \underline{(1-\theta)(\alpha^T y + b_i)}) \\
&= \exp(\underline{\theta(\alpha^T x + b_i)}) \cdot \exp(\underline{(1-\theta)(\alpha^T y + b_i)}) \\
&= [\exp(\underline{\alpha^T x + b_i})]^\theta \cdot [\exp(\underline{\alpha^T y + b_i})]^{1-\theta} \\
&= [\underline{h(x)}]^\theta \cdot [\underline{h(y)}]^{1-\theta} \leq [\underline{h(x)}]^\theta \cdot [\underline{h(y)}]^{1-\theta},
\end{aligned}$$

According to Part 1 ↗

for $\forall x, y \in \mathbb{R}^n, \forall \theta \in [0, 1]$

So $h(x)$ is log-convex. proved.

Part 3: Prove: if $a > 0, b > 0, \alpha, \beta \in [0, 1]$ and $\alpha + \beta = 1$
then $a^\alpha b^\beta \leq \alpha a + \beta b$.

Proof: First we prove: If $A > 0, B > 0, \frac{1}{p} + \frac{1}{q} = 1, p > 1, q > 1$

$$\text{then } AB \leq \frac{A^p}{p} + \frac{B^q}{q}. \quad \boxed{(1)}$$

th-1

$$\text{Let } t = \frac{A^p}{B^q}, \varphi = \frac{1}{p}$$

$$\begin{aligned}
\text{then } \boxed{(1)} \iff & \left(\frac{A^p}{B^q} \cdot B^q \right)^{\frac{1}{p}} \cdot (B^q)^{1-\frac{1}{p}} \leq \frac{1}{p} \left(\frac{A^p}{B^q} \cdot B^q \right) + \left(1 - \frac{1}{p} \right) \cdot B^q \\
\iff & (t \cdot B^q)^{\varphi} \cdot (B^q)^{1-\varphi} \leq \varphi(t \cdot B^q) + (1-\varphi)B^q \\
\iff & t^\varphi \cdot B^q \leq (\varphi t + 1 - \varphi) B^q \\
\iff & t^\varphi \leq \varphi t + 1 - \varphi \quad \boxed{(2)}
\end{aligned}$$

$$\text{Let } d(t) = t^\varphi - \varphi t, \text{ thus } d'(t) = \varphi t^{\varphi-1} - \varphi = \varphi(t^{\varphi-1} - 1).$$

Since $p > 1$, then $0 < \varphi < 1$, i.e., in $d'(t)$ the exponent of t is negative
so $d'(t)$ is monotonically decreasing in $t \in (0, +\infty)$,

And since $d'(t)$ is differentiable at $t > 0$, ($d''(t) = \varphi(\varphi-1)t^{\varphi-2}$), we have:

$\begin{cases} \text{when } t \in (0, 1), d'(t) > d'(1) = 0, \text{ thus } d(t) \text{ is monotonically increasing.} \\ \text{when } t > 1, d'(t) < d'(1) = 0, \text{ thus } d(t) \text{ is monotonically decreasing.} \end{cases}$

$$\text{Therefore, } t^\varphi - \varphi t = d(t) \leq d(1) = 1 - \varphi \iff \boxed{(2)} \iff \boxed{(1)}$$

Now, we've proved th-1.

We reformulate th-1 by some transformation:

$$\text{Let } \alpha = \frac{1}{p}, \beta = \frac{1}{q}, A^p = a, B^q = b,$$

(thus $A > 0, B > 0, \alpha, \beta \in (0, 1)$ and $\alpha + \beta = 1$)

$$\text{then } AB \leq \frac{A^p}{p} + \frac{B^q}{q} \iff a^\alpha b^\beta \leq \alpha a + \beta b \quad (3)$$

There are still conditions with boundary values.

$$\text{when } \alpha=1, \beta=0, a^\alpha b^\beta = a = \alpha a + \beta b$$

$$\text{when } \alpha=0, \beta=1, a^\alpha b^\beta = b = \alpha a + \beta b$$

Altogether, we have proved:

When $a > 0, b > 0, \alpha, \beta \in [0, 1]$ and $\alpha + \beta = 1$
 then $a^\alpha b^\beta \leq \alpha a + \beta b$.

Part 4: Prove:

Assume $H_1(x), H_2(x) : \mathbb{R}^n \rightarrow \mathbb{R}$, $H_1(x), H_2(x) > 0$ for $\forall x \in \mathbb{R}^n$,
 If $H_1(x)$ and $H_2(x)$ are both log-convex,
 then $H(x) = H_1(x) + H_2(x)$ is also log-convex. (*-0)

Proof: Since $H_1(x), H_2(x)$ are log-convex, as we proved in Part 1,

$$\left\{ \begin{array}{l} H_1(\theta x + (1-\theta)y) \leq [H_1(x)]^\theta [H_1(y)]^{1-\theta}, \quad \forall x, y \in \mathbb{R}^n, \forall \theta \in [0, 1] \end{array} \right. \quad (4)$$

$$\left\{ \begin{array}{l} H_2(\theta x + (1-\theta)y) \leq [H_2(x)]^\theta [H_2(y)]^{1-\theta}, \quad \forall x, y \in \mathbb{R}^n, \forall \theta \in [0, 1] \end{array} \right. \quad (5)$$

(4)+(5), we have

$$\begin{aligned} H(\theta x + (1-\theta)y) &= H_1(\theta x + (1-\theta)y) + H_2(\theta x + (1-\theta)y) \\ &\leq \underbrace{[H_1(x)]^\theta [H_1(y)]^{1-\theta}}_{\text{pink underline}} + \underbrace{[H_2(x)]^\theta [H_2(y)]^{1-\theta}}_{\text{pink underline}} \end{aligned}$$

Notice that what we want to prove is: $H(x)$ is log-convex, i.e.,

$$H(\theta x + (1-\theta)y) \leq [H(x)]^\theta [H(y)]^{1-\theta} = \underbrace{[H_1(x) + H_2(x)]^\theta}_{\text{yellow underline}} \underbrace{[H_1(y) + H_2(y)]^{1-\theta}}_{\text{yellow underline}}$$

so to prove (*-0), we just need to prove (*-1) as below

$$\underbrace{[H_1(x)]^\theta [H_1(y)]^{1-\theta}}_{\text{pink underline}} + \underbrace{[H_2(x)]^\theta [H_2(y)]^{1-\theta}}_{\text{pink underline}} \leq \underbrace{[H_1(x) + H_2(x)]^\theta}_{\text{yellow underline}} \underbrace{[H_1(y) + H_2(y)]^{1-\theta}}_{\text{yellow underline}} \quad (*-1)$$

let $a = H_1(x), b = H_1(y), c = H_2(x), d = H_2(y)$

then to prove (*-1) becomes to prove (*-2) as below:

$$a^\theta b^{1-\theta} + c^\theta d^{1-\theta} \leq (a+c)^\theta (b+d)^{1-\theta} \quad (*-2)$$

divide $\frac{(a+c)^\theta (b+d)^{1-\theta}}{>0}$ on both side, we have

$$\frac{a^\theta b^{1-\theta} + c^\theta d^{1-\theta}}{(a+c)^\theta (b+d)^{1-\theta}} \leq 1 \quad (*-3)$$

Since $\begin{cases} a = H_1(x) \\ c = H_2(x) \end{cases}$, $\begin{cases} b = H_1(y) \\ d = H_2(y) \end{cases}$ → there is no relationship between the value of a and c (or: b and d), → thus, $a+c > 0$ can be arbitrary. Similarly, $b+d > 0$ can be arbitrary.

Thus, we assume that $\underline{a+c=k_1}$, $\underline{b+d=k_2}$, $k_1, k_2 > 0$ are both arbitrary.

and we set $a_1 = \frac{a}{k_1}$, $c_1 = \frac{c}{k_1}$, such that $a_1 + c_1 = 1$

Similarly we set $b_1 = \frac{b}{k_2}$, $d_1 = \frac{d}{k_2}$, such that $b_1 + d_1 = 1$

So the LHS of (*-3) can be transformed as below:

$$\begin{aligned} \frac{a^\theta b^{1-\theta} + c^\theta d^{1-\theta}}{(a+c)^\theta (b+d)^{1-\theta}} &= \frac{(k_1 a_1)^\theta (k_2 b_1)^{1-\theta} + (k_1 c_1)^\theta (k_2 d_1)^{1-\theta}}{k_1^\theta (a_1 + c_1)^\theta k_2^{1-\theta} (b_1 + d_1)^{1-\theta}} \\ &= \frac{\cancel{(k_1^\theta k_2^{1-\theta})} a_1^\theta b_1^{1-\theta} + \cancel{(k_1^\theta k_2^{1-\theta})} c_1^\theta d_1^{1-\theta}}{\cancel{(k_1^\theta k_2^{1-\theta})} (a_1 + c_1)^\theta (b_1 + d_1)^{1-\theta}} \\ &= \frac{a_1^\theta b_1^{1-\theta} + c_1^\theta d_1^{1-\theta}}{(a_1 + c_1)^\theta (b_1 + d_1)^{1-\theta}} \quad (6) \end{aligned}$$

According to (6), we can reformulate (*-3) as:

$$\frac{a_1^\theta b_1^{1-\theta} + c_1^\theta d_1^{1-\theta}}{(a_1 + c_1)^\theta (b_1 + d_1)^{1-\theta}} \leq 1 \iff a_1^\theta b_1^{1-\theta} + c_1^\theta d_1^{1-\theta} \leq 1, \text{ where } a_1 + c_1 = b_1 + d_1 = 1$$

Thus, to prove (*-0), we just need to prove (*-4).

$$a_1, b_1, c_1, d_1 > 0, \theta \in [0, 1]$$

As we have proved in Part 3,

$$a_1^\theta + b_1^{1-\theta} \leq \theta a_1 + (1-\theta)b_1 \quad (7)$$

(in next page)

$$c_i^\theta d_i^{1-\theta} \leq \theta c_i + (1-\theta)d_i \quad (8)$$

(7) + (8), we have :

$$\begin{aligned} a_i^\theta b_i^{1-\theta} + c_i^\theta d_i^{1-\theta} &\leq \theta a_i + (1-\theta)b_i + \theta c_i + (1-\theta)d_i \\ &= \theta(a_i + c_i) + (1-\theta)(b_i + d_i) \\ &= \theta + 1-\theta \\ &= 1 \end{aligned}$$

Thus we proved (*-4),

i.e., we proved (*-o). ✓

Part 5:

Prove: $j(x) = \sum_{i=1}^m \exp(a_i^T x + b_i)$ is log-convex
 $a_i, x \in R^n, b_i \in R.$

let $h_i(x) = \exp(a_i^T x + b_i)$, $i=1, 2, \dots, m$,
then $j(x) = \sum_{i=1}^m h_i(x).$

As we have proved in Part 2,

$h_i(x)$ is log-convex, $i=1, 2, \dots, m$

Then as we have proved in Part 4,

$h_1(x) + h_2(x)$ is log-convex.

Then $(h_1(x) + h_2(x)) + h_3(x)$ is log-convex

since $h_1(x) + h_2(x)$ is convex and $h_3(x)$ is ^{log-}convex.

Thus, inductively, we can finally conclude that

$j(x) = \sum_{i=1}^m h_i(x)$ is log-convex.

Part 6: Get the conclusion.

Since $j(x)$ is log-convex.

$$f(x) = \log(j(x)) \text{ is convex}$$

$$= \log\left(\sum_{i=1}^m \exp(a_i^T x + b_i)\right)$$

(condition-1)

According to 5(a), we know $f(x)$ is differentiable. (condition-2)

Some interpretation for (condition-2):

And since $a_i^T x + b_i$, $\exp(x)$, $\log(x)$, and sum-function
 $x \in R^n$ $x \in R$ respective

are all differentiable in their definition domains.

so $f(x)$ which is composed of these differentiable function

is also differentiable. (condition-2')

by condition-1
 condition-2, we have

$$x^{(*)} = \underset{x \in R^n}{\operatorname{argmin}} f(x) \Leftrightarrow \nabla f(x^{(*)}) = 0.$$

If we use gradient descent to solve $\min_{x \in R^n} f(x)$,

and GD converges, then

we have

$$\lim_{k \rightarrow \infty} \nabla f(x^{(k)}) = 0 \quad \Longleftrightarrow \quad \lim_{k \rightarrow \infty} x^{(k)} = \underset{x \in R^n}{\operatorname{argmin}} f(x)$$

i.e. It will converge to the global minimizer.

The minimizer may not exist.

-2