

Final Review - Proofs

Sunday, 6 December 2020 2:18 AM

$$\sup_{\substack{x \in \mathbb{R}^n \\ x \neq 0}} \frac{\|Ax\|_p}{\|x\|_p} \Rightarrow \sup_{\substack{x \in \mathbb{R}^n \\ \|x\|_p=1}} \|Ax\|_p.$$

1. Form $\sup_{\substack{x \in \mathbb{R}^n \\ x \neq 0}} \frac{\|Ax\|_p}{\|x\|_p}$ + to $\sup_{\substack{x \in \mathbb{R}^n \\ \|x\|_p=1}} \|Ax\|_p$.

We have: $\sup_{\substack{x \in \mathbb{R}^n \\ x \neq 0}} \frac{\|Ax\|_p}{\|x\|_p}$,

not that $\|\frac{x}{\|x\|_p}\|_p = 1$, i.e. $\frac{x}{\|x\|_p}$ has one unit length,

we express $\|A\|_p$ in terms of supremum over all vectors of unit length.

$$\|A\|_p = \sup_{\|x\|_p=1} \|Ax\|_p.$$

2. Prove operator matrix norm is indeed a norm:

$$\|A\|_{p \rightarrow p}$$

(P2) First consider $A = 0$, the all-zero matrix. Clearly,

$$\|0\|_{a,b} = \sup_{\|x\|_b \leq 1} \|0x\|_a = 0 \quad (11)$$

Next suppose that $\|A\|_{a,b} = 0$ for some A . Then,

$$\sup_{\|x\|_b \leq 1} \|Ax\|_a = 0 \quad (12)$$

$$\Rightarrow \|Av\|_a = 0 \quad \text{for all } \|v\|_b \leq 1, \text{ by definition} \quad (13)$$

$$\Rightarrow Av = 0 \quad \text{from (P2), for all } \|v\|_b \leq 1 \quad (14)$$

$$\Rightarrow A = 0 \quad (15)$$

The last step follows since, if A has even a single non-zero entry, it will always be possible to find an $v : \|v\|_b \leq 1$, such that $Av \neq 0$.

(P3) Trivial, since

$$\sup_{\|x\|_b \leq 1} \|tAx\|_a = \sup_{\|x\|_b \leq 1} |t| \|Ax\|_a = |t| \sup_{\|x\|_b \leq 1} \|Ax\|_a \quad (16)$$

from (P3) for vector norm.

(P4) Consider the rhs

$$\|A\|_{a,b} + \|B\|_{a,b} = \sup_{\|x\|_b \leq 1} \|Ax\|_a + \sup_{\|y\|_b \leq 1} \|By\|_a \quad \text{by definition} \quad (17)$$

$$\geq \|Av\|_a + \|Bu\|_a \quad \forall \|v\|_b \leq 1, \|u\|_b \leq 1 \quad (18)$$

This statement holds for all u, v such that $\|v\|_b \leq 1, \|u\|_b \leq 1$. Therefore it also holds when $u = v$ with $\|u\|_b \leq 1$, i.e.,

$$\|A\|_{a,b} + \|B\|_{a,b} \geq \|Au\|_a + \|Bu\|_a \quad \forall \|u\|_b \leq 1 \quad (19)$$

$$\geq \|Au + Bu\|_a \quad \forall \|u\|_b \leq 1 \text{ from (P4)} \quad (20)$$

which implies that

$$\|A\|_{a,b} + \|B\|_{a,b} \geq \sup_{\|u\|_b \leq 1} \|Au + Bu\|_a \quad (21)$$

$$= \|A + B\|_{a,b} \quad (22)$$

$$\|A + B\|_{a,b} = \min_{\|u\|_b \leq 1} \|A + Bu\|_a$$
(22)

3. Proof the following computation of l-inf norm, l1 and l2-norm

Example: $a = \begin{pmatrix} 1 \\ 1 \\ 3 \end{pmatrix}$, $\|a\|_\infty = 1$, $\|a\|_2 = \sqrt{\frac{\pi}{6}}$, $\|a\|_1 = +\infty$, $\boxed{\text{if } 0}$
 $\therefore a \notin l_\infty, a \in l_2, a \notin l_1$

Proof:

Then $\|a\|_\infty = \sup_i |a_i| = 1$ 呢個係從 $\sin(x)$ taylor series 推的
 $\|a\|_2 = \left(\sum_{i=1}^{\infty} (1/i)^2 \right)^{1/2} = \left(\frac{\pi^2}{6} \right)^{1/2} = \frac{\pi}{\sqrt{6}}$
 $\|a\|_1 = \left(\sum_{i=1}^{\infty} \frac{1}{i} \right) = +\infty$

4. Proof of k-mean algorithms:

Step 1: Fix the representatives z_1, \dots, z_k , find the best partitions

G_1, \dots, G_k , i.e., solve

$$\min_{G_1, \dots, G_k} \sum_{j=1}^k \left(\sum_{i \in G_j} \|x_i - z_j\|_2^2 \right) \quad \dots \quad (1)$$

$$(1) \Leftrightarrow \min_{C_1, \dots, C_N} \sum_{j=1}^k \sum_{i \in C_j} \|x_i - z_j\|_2^2 \quad \text{rewrite}$$

finding the partition G_1, \dots, G_k is equivalent to

finding C_1, C_2, \dots, C_N . So (1) becomes

$$\min_{C_1, C_2, \dots, C_N} \underbrace{\|x_1 - z_1\|_2^2}_{\text{depends on } C_1 \text{ only}} + \underbrace{\|x_2 - z_2\|_2^2}_{\text{depends on } C_2 \text{ only}} + \dots + \underbrace{\|x_N - z_N\|_2^2}_{\text{depends on } C_N \text{ only}}$$

↓

$$\min_{C_i} \|x_i - z_{C_i}\|_2^2 \quad i=1, 2, \dots, N.$$

Since $C_i \in \{1, 2, \dots, k\}$, to get C_i , we only need to compare

$$\|x_i - z_1\|_2^2, \|x_i - z_2\|_2^2, \dots, \|x_i - z_k\|_2^2$$

and choose the minimum from it. i.e.,

$$C_i = \arg \min_{j \in \{1, \dots, k\}} \|x_i - z_j\|_2^2 \quad , \quad i=1, 2, \dots, N.$$

In other words,

x_i is assigned to the group whose representative vector is the closest to x_i .

Step 2: Fix the groups G_1, \dots, G_K , find the best representatives

$$\bar{z}_1, \dots, \bar{z}_K, \text{ i.e., solve} \\ \min_{\bar{z}_1, \dots, \bar{z}_K} \sum_{j=1}^K \left(\sum_{i \in G_j} \|x_i - \bar{z}_j\|_2^2 \right) \quad \dots \quad (2)$$

$$(2) \Leftrightarrow \min_{\bar{z}_1, \dots, \bar{z}_K} \sum_{j=1}^K \sum_{i \in G_j} \|x_i - \bar{z}_j\|_2^2 \quad \text{Swap of summation order} \\ \Leftrightarrow \min_{\bar{z}_1, \dots, \bar{z}_K} \left(\sum_{i \in G_1} \|x_i - \bar{z}_1\|_2^2 + \sum_{i \in G_2} \|x_i - \bar{z}_2\|_2^2 + \dots + \sum_{i \in G_K} \|x_i - \bar{z}_K\|_2^2 \right) \\ \Leftrightarrow \min_{\bar{z}_j} \sum_{i \in G_j} \|x_i - \bar{z}_j\|_2^2, \quad j = 1, \dots, K \quad \text{independent}$$

Taking derivative w.r.t. \bar{z}_{jl} and setting it to 0, we obtain that

the solution \bar{z}_{jl} satisfies

$$2 \sum_{i \in G_j} (\bar{z}_{jl} - x_{il}) = 0 \\ \Rightarrow \bar{z}_{jl} = \left(\sum_{i \in G_j} x_{il} \right) / |G_j| \quad (|G_j| \text{ is the number of elements in } G_j) \\ l = 1, 2, \dots, n.$$

In vector form,

$$\begin{pmatrix} \bar{z}_{j1} \\ \bar{z}_{j2} \\ \vdots \\ \bar{z}_{jn} \end{pmatrix} = \frac{1}{|G_j|} \cdot \sum_{i \in G_j} \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{in} \end{pmatrix} \Leftrightarrow \bar{z}_j = \frac{1}{|G_j|} \left(\sum_{i \in G_j} x_i \right)$$

In other words,

\bar{z}_j is the mean of all vectors in G_j .

5.

Lemma: If $x^{(k)} \rightarrow x \in V$, then $\{x^{(k)}\}$ is a cauchy sequence.

Proof. $x^{(k)} \rightarrow x$ implies: $\forall \frac{\varepsilon}{2} > 0, \exists K$ st. $\forall k \geq K \quad \|x^{(k)} - x\| \leq \frac{\varepsilon}{2}$.

Therefore, $\|x^{(k)} - x^{(l)}\| \leq \|x^{(k)} - x\| + \|x^{(l)} - x\| \leq \varepsilon$. $\forall k, l \geq K$. \blacksquare

6. If $\{x^{(k)}\}$ is a cauchy sequence, $x^{(k)} \rightarrow x \in V$ may not hold.

Consider this example:

Let $a^{(k)} = \begin{pmatrix} 1 \\ y_2 \\ \vdots \\ y_k \\ 0 \\ \vdots \end{pmatrix}$ $a = \begin{pmatrix} 1 \\ y_2 \\ y_3 \\ \vdots \\ y_k \\ \vdots \end{pmatrix}$

$$\|a^{(k)} - a\|_\infty = \left\| \begin{pmatrix} 0 \\ \vdots \\ y_k \\ \vdots \\ 0 \end{pmatrix} \right\|_\infty = \frac{1}{k} < \varepsilon \quad \text{if } k, l \geq \frac{1}{\varepsilon} + 1 \quad (\text{assume } l \geq k)$$

$\Rightarrow \{a^{(k)}\}$ is Cauchy sequence.

$$\|a^{(k)}\|_1 = \sum_{i=1}^k f_i) < +\infty \quad \text{as } k < +\infty$$

$$\|a\|_1 = \sum_{i=1}^k f_i + \sum_{i=1}^\infty f_i = +\infty$$

$$\lim_{k \rightarrow \infty} \|a^{(k)} - a\|_\infty = \lim_{k \rightarrow \infty} \left\| \begin{pmatrix} 0 \\ \vdots \\ y_{k+1} \\ \vdots \\ 0 \end{pmatrix} \right\|_\infty = \lim_{k \rightarrow \infty} \frac{1}{k+1} = 0$$

$a^{(k)} \rightarrow a$, $a^k \in V$, but $a \notin V$

\therefore It is proved that if $\{x^{(k)}\}$ is Cauchy sequence,

$$x^{(k)} \rightarrow x \in V$$

7.

Note that the constants C_1, C_2 depend on V .

Example: Consider \mathbb{R}^n and $\|\cdot\|_1, \|\cdot\|_2, \|\cdot\|_\infty$

- $\|\cdot\|_1$ and $\|\cdot\|_2$ are equivalent because

$$\|a\|_2 \leq \|a\|_1 \leq \sqrt{n} \|a\|_2$$

- $\|\cdot\|_2$ and $\|\cdot\|_\infty$ are equivalent

$$\|a\|_\infty \leq \|a\|_2 \leq \sqrt{n} \|a\|_\infty$$

- $\|\cdot\|_1$ and $\|\cdot\|_\infty$ are equivalent

$$\|a\|_\infty \leq \|a\|_1 \leq n \|a\|_\infty$$

TBC....

8. Proof of Cauchy Schwartz Inequality

Proof • if $y=0$. Then obviously

$$\langle x, y \rangle^2 = \langle x, 0 \rangle^2 = 0 \leq 0 = \langle x, x \rangle \langle y, y \rangle$$

- It remains to prove the inequality with $y \neq 0$.

Let $\lambda \in \mathbb{R}$ be an arbitrary number

$$\begin{aligned} 0 \leq \langle x + \lambda y, x + \lambda y \rangle &= \langle x, x \rangle + \lambda \langle y, x \rangle + \lambda \langle x, y \rangle + \lambda^2 \langle y, y \rangle \\ &= \langle x, x \rangle + 2\lambda \langle x, y \rangle + \lambda^2 \langle y, y \rangle \end{aligned}$$

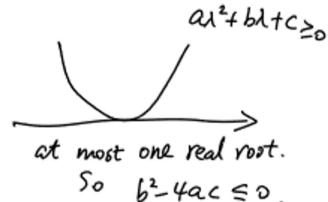
Thus, $\lambda^2 \langle y, y \rangle + 2\lambda \langle x, y \rangle + \langle x, x \rangle \geq 0$. $\forall \lambda \in \mathbb{R}$.

Since $y \neq 0$, $\langle y, y \rangle > 0$. $\Rightarrow f(\lambda) \geq 0$

So, $f(\lambda)$ a quadratic function of λ that takes non-negative values only.

There is at most one root of the quadratic function.

$$\begin{aligned} \text{So, } (2\langle x, y \rangle)^2 - 4\langle y, y \rangle \langle x, x \rangle &\leq 0 \\ \Rightarrow \langle x, y \rangle^2 &\leq \langle x, x \rangle \langle y, y \rangle \end{aligned}$$



9. Conditions for $\langle x, y \rangle^2 = \langle x, x \rangle \langle y, y \rangle$: $x = \alpha * y$ or $y = \alpha * x$

- Next when the equality holds true, i.e., $\langle x, y \rangle^2 = \langle x, x \rangle \langle y, y \rangle$,

- if $y=0$, then obviously $y=\alpha x$, where $\alpha=0$.

- if $y \neq 0$, there is exactly one real root of $f(\lambda)$.

$$\exists \text{ a unique } \lambda \in \mathbb{R}, \lambda^2 \langle y, y \rangle + 2\lambda \langle x, y \rangle + \langle x, x \rangle = 0$$

$$\Rightarrow \langle x + \lambda y, x + \lambda y \rangle = 0 \Rightarrow x + \lambda y = 0 \Rightarrow x = \alpha y \text{ with } \alpha = -\lambda.$$

Finally, we show that if $x = \alpha y$ or $y = \alpha x$, then $\langle x, y \rangle^2 = \langle x, x \rangle \langle y, y \rangle$.

10. $\|x\| = \sqrt{\langle x, x \rangle}$ defines a norm

Proof. ① $\|x\| = (\langle x, x \rangle)^{1/2} \geq 0$ and $\|x\| = (\langle x, x \rangle)^{1/2} = 0 \Leftrightarrow x = 0$.

$$\text{② } \|\alpha x\| = (\langle \alpha x, \alpha x \rangle)^{1/2} = (\alpha^2 \langle x, x \rangle)^{1/2} = |\alpha| \|x\|$$

$$\begin{aligned} \text{③ } \|x+y\|_2^2 &= \langle x+y, x+y \rangle = \langle x, x \rangle + \langle x, y \rangle + \langle y, x \rangle + \langle y, y \rangle \\ &= \|x\|^2 + \|y\|^2 + 2\langle x, y \rangle \\ &\leq \|x\|^2 + \|y\|^2 + 2\|x\| \|y\| \quad \left(\begin{array}{l} \text{Note that Cauchy-Schwarz} \\ \text{becomes} \\ |\langle x, y \rangle| \leq \|x\| \|y\| \end{array} \right) \\ &= (\|x\| + \|y\|)^2 \quad \text{⊗} \end{aligned}$$

11. L2 space is a hilbert space

<https://math.stackexchange.com/questions/1259364/show-that-l2-is-a-hilbert-space>

12. Rewrite k-mean using kernel trick to map it to hilbert space

第一步，消除 $\phi(x_j)$.

$$\text{简化成: } \begin{cases} \textcircled{1} & c_i = \arg \min_{j=1 \dots k} \|\phi(x_i) - \frac{1}{|G_j|} \left(\sum_{x \in G_j} \phi(x_e) \right)\| \\ \textcircled{2} & G_j = \{i \mid c_i = j\}, \quad j=1 \dots k. \end{cases}$$

$$c_i = \arg \min_{j=1 \dots k} \|\phi(x_i) - \frac{1}{|G_j|} \left(\sum_{x \in G_j} \phi(x_e) \right)\|^2$$

$\stackrel{\text{第k步}}{\Rightarrow}$

took square term, 无影响, same minimizer.



$$c_i = \arg \min_{j=1 \dots k} \langle \phi(x_i) - \frac{1}{|G_j|} \left(\sum_{x \in G_j} \phi(x_e) \right), \phi(x_i) - \frac{1}{|G_j|} \sum_{x \in G_j} \phi(x_e) \rangle$$



$$c_i = \arg \min_{j=1 \dots k} \langle \phi(x_i), \phi(x_i) \rangle - 2 \langle \phi(x_i), \frac{1}{|G_j|} \sum_{x \in G_j} \phi(x_e) \rangle + \langle \frac{1}{|G_j|} \sum_{x \in G_j} \phi(x_e), \frac{1}{|G_j|} \sum_{x \in G_j} \phi(x_e) \rangle$$



$$c_i = \arg \min_{j=1 \dots k} \langle \phi(x_i), \phi(x_i) \rangle - \frac{2}{|G_j|} \sum_{x \in G_j} \langle \phi(x_i), \phi(x_e) \rangle + \frac{1}{|G_j|^2} \sum_{x_1 \in G_j} \sum_{x_2 \in G_j} \langle \phi(x_{e1}), \phi(x_{e2}) \rangle$$

13. Prove requirements for kernel function

Prove

$K(x,y)$ 一些必要條件:

$$\text{inner product property}$$

$$\textcircled{1} \quad K(x,y) = \langle \phi(x), \phi(y) \rangle \stackrel{?}{=} \langle \phi(y), \phi(x) \rangle = K(y,x)$$

So, $K(x,y) = K(y,x) \quad \forall x, y \in \mathbb{R}^n \Rightarrow$ show K is symmetric

$$\textcircled{2} \quad \text{Let } y_1, y_2, \dots, y_m \text{ 係 } \mathbb{R}^n \text{ space 中 } m \text{ 支 vector, for any } c = \begin{pmatrix} c_1 \\ \vdots \\ c_m \end{pmatrix} \in \mathbb{R}^m$$

$$0 \leq \underbrace{\left\langle \sum_{i=1}^m c_i \phi(y_i), \sum_{i=1}^m c_i \phi(y_i) \right\rangle}_{\text{in H space 由,}} \stackrel{\uparrow}{\substack{\text{By inner product property}}} \sum_{i=1}^m \sum_{j=1}^m c_i c_j \langle \phi(y_i), \phi(y_j) \rangle$$

\because c 只係實數

$$= \sum_{i=1}^m \sum_{j=1}^m c_i c_j K(y_i, y_j)$$

$$= \sum_{i=1}^m \sum_{j=1}^m c_i c_j K_{ij}$$

$$= C^T K C$$

$$\Rightarrow C^T K C \geq 0$$

$\therefore K$ is SPSD

$$\Rightarrow \text{where } K = \begin{bmatrix} K(y_1, y_1) & \cdots & K(y_1, y_m) \\ K(y_2, y_1) & \ddots & \vdots \\ \vdots & \cdots & K(y_m, y_m) \end{bmatrix} \in \mathbb{R}^{m \times m}$$

關係呢個 matrix 裏 \Rightarrow

住 所以 K 是 SPSD 係其中一個必要條件

14. Any linear function in \mathbb{R}^n can be written as $\langle a, x \rangle$

To see this, let e_1, e_2, \dots, e_n , where $e_i = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \end{pmatrix}$ (ith component), be a basis of \mathbb{R}^n ,

So that any $x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^n$ is written as $x = x_1 e_1 + x_2 e_2 + \dots + x_n e_n$.

Therefore, if f is a linear function, then

$$\begin{aligned} f(x) &= f(x_1 e_1 + x_2 e_2 + \dots + x_n e_n) \\ &= x_1 f(e_1) + x_2 f(e_2) + \dots + x_n f(e_n) \quad \text{by property of linear functions} \\ &= \langle a, x \rangle, \end{aligned}$$

where $a = \begin{pmatrix} f(e_1) \\ f(e_2) \\ \vdots \\ f(e_n) \end{pmatrix} \in \mathbb{R}^n$.

15. The inner product for a linear product is unique

- Furthermore, the representation of a linear function $f(x) = \langle a, x \rangle$ is unique, which means there is only one vector $a \in \mathbb{R}^n$ for which $f(x) = \langle a, x \rangle$ holds for all x .

Indeed, suppose that a is not unique, i.e., we have two vectors a, b such that $f(x) = \langle a, x \rangle$ and $f(x) = \langle b, x \rangle$ for all $x \in \mathbb{R}^n$.

Then, let $x = e_i$: $f(e_i) = \langle a, e_i \rangle = a_i$ and $f(e_i) = \langle b, e_i \rangle = b_i$

So, $a_i = b_i$, $i = 1, 2, \dots, n$.

Therefore $a = b$.

16. a fixed t_0 function cannot have $h(t)$ in $\langle f, h \rangle$

$F : C[a,b] \rightarrow \mathbb{R}$ defined by $F(f) = f(t_0)$ for a fixed $t_0 \in \mathbb{R}$.

F is linear on $C[a,b]$

If we define an inner product on $C[a,b]$

$$\text{by } \langle f, g \rangle = \int_a^b f(t)g(t)dt$$

and the induced norm is $\|f\|_n = \left(\int_a^b |f(t)|^2 dt \right)^{\frac{1}{2}}$

If we want to represent F into inner product form

$$f(t_0) = F(f) = \int_a^{t_0} f(t)h(t)dt$$

we require $h(t) = 0$ when $t \neq t_0$.

呢到意思係，我地要令當 t 不是 t_0 的時候，

But this will make make $\int_a^b f(t)h(t)dt = 0$ contribute to integral 嘅其他 term 都係 0

We can't find h .

所以 $h(t)$ 要係 0 when $t \neq t_0$

17.

$S_{a,0} = \{x \in V \mid \langle a, x \rangle = 0\}$ is a hyperplane.

Then, $\forall x, y \in S_{a,0}$ and $\alpha, \beta \in \mathbb{R}$,

$$\langle a, \alpha x + \beta y \rangle = \alpha \langle a, x \rangle + \beta \langle a, y \rangle = 0 \Rightarrow \alpha x + \beta y \in S.$$

Therefore, $S_{a,0}$ is a plane.

linear

18.

$S_{a,b} = \{x \in V \mid \langle a, x \rangle = b\}$ is also a hyperplane.

Now let's consider $a \in H$

$$S_{a,b} = \{x \mid \langle a, x \rangle = b\} \text{ for } b \in \mathbb{R} \text{ is given.}$$

Let $x_0 \in S_{a,b}$, i.e., $\langle a, x_0 \rangle = b$, be fixed.

Then $S_{a,b} = S_{a,0} + x_0$ because:

$$\begin{matrix} b-b=0 \\ \swarrow \end{matrix}$$

$$\begin{aligned} \textcircled{1} \quad \forall x \in S_{a,b} \quad \langle a, x - x_0 \rangle &= \langle a, x \rangle - \langle a, x_0 \rangle = 0, \Rightarrow x - x_0 \in S_{a,0}. \\ &\Rightarrow x \in S_{a,0} + x_0. \end{aligned}$$

any element in $S_{a,b}$ can be written in $S_{a,0} + x_0$, so it is subset

$$\textcircled{2} \quad \forall x \in S_{a,0} \quad \langle a, x + x_0 \rangle = \langle a, x \rangle + \langle a, x_0 \rangle = b \Rightarrow x + x_0 \in S_{a,b}$$

In other words,

$S_{a,b}$ is a shift of a hyperplane, still called a hyperplane.

19. Proof of this theorem in projection onto hyperplanes

19. Proof of this theorem in projection onto hyperplanes

Theorem: z is a solution of $\min_{x \in S} \|x - y\|$ if and only if $z \in S$ and $\langle z - y, z - x \rangle = 0 \quad \forall x \in S$.

Proof. ① We first prove that: If $z \in S$ is a solution of $\min_{x \in S} \|x - y\|$, then $\langle z - y, z - x \rangle = 0 \quad \forall x \in S$. arbitrary x

z 與 x 的那條直線

Since z is a solution, $z \in S$, i.e., $\langle a, z \rangle = b$.

$\forall x \in S$ and $t \in \mathbb{R}$, it is easy to see that

$$\langle a, (1+t)z - tx \rangle = (1+t)\langle a, z \rangle - t\langle a, x \rangle = b.$$

Therefore, $(1+t)z - tx \in S$.

Since z is closest to y on S , we have

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

$$\text{i.e., } t\langle z - y, z - x \rangle \geq -\frac{t^2}{2}\|z - x\|^2.$$

- If we choose $t > 0$,

$$\langle z - y, z - x \rangle \geq -\frac{t}{2}\|z - x\|^2$$

Letting $t \rightarrow 0_+$ gives $\langle z - y, z - x \rangle \geq 0$.

- If we choose $t < 0$,

$$\langle z - y, z - x \rangle \leq -\frac{t}{2}\|z - x\|^2$$

Letting $t \rightarrow 0_-$ gives $\langle z - y, z - x \rangle \leq 0$

Altogether, z satisfies $\langle z - y, z - x \rangle = 0 \quad \forall x \in S$.

- We then show that $z \in S$ satisfies $\langle z - y, z - x \rangle = 0$, then z is a solution of $\min_{x \in S} \|x - y\|$, by direct calculation.

Since $\langle z - y, z - x \rangle = 0 \quad \forall x \in S$,

$$\|x - y\|^2 = \|(z - x) - (z - y)\|^2$$

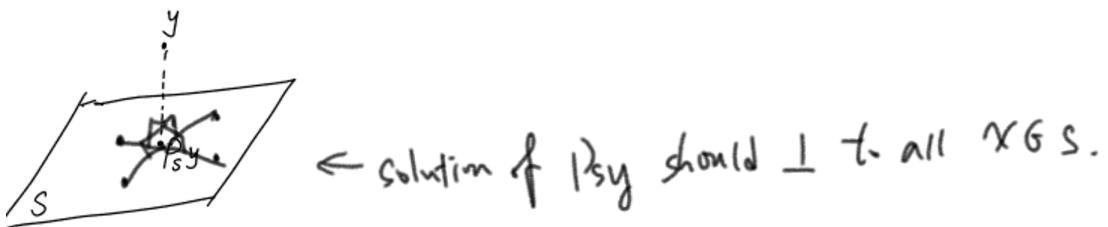
$\langle z - y, z - x \rangle$ 只能是 0, 因为 t 是任意的

since $\|x - z\| \geq \|x - y\|$

$$\begin{aligned}\|x - y\|^2 &= \|(z-x) - (z-y)\|^2 \\ &= \|z-x\|^2 + \|z-y\|^2 - 2\langle z-x, z-y \rangle \\ &= \|z-y\|^2 + \|z-x\|^2 \geq \|z-y\|^2 \quad \forall x \in S.\end{aligned}$$

This, together with $z \in S$, implies z minimizes $\|x-y\|^2$ in $x \in S$.

□



20. P_Sy should be the following formula:

Theorem: The solution of $\min_{x \in S} \|x - y\|^2$ exists and unique, which is given by $y - \left(\frac{\langle a, y \rangle - b}{\|a\|^2} \right) a$

proof. denote $z = y - \left(\frac{\langle a, y \rangle - b}{\|a\|^2} \right) a$.

$$\begin{aligned} \textcircled{1} \quad \langle a, z \rangle &= \langle a, y \rangle - \left(\frac{\langle a, y \rangle - b}{\|a\|^2} \right) \langle a, a \rangle \\ &= \langle a, y \rangle - (\langle a, y \rangle - b) = b, \quad \text{so } z \in S. \end{aligned}$$

\textcircled{2} $\forall x \in S$,

$$\begin{aligned} \langle z - y, z - x \rangle &= - \frac{\langle a, y \rangle - b}{\|a\|^2} \langle a, z - x \rangle \\ &= - \frac{\langle a, y \rangle - b}{\|a\|^2} (\langle a, z \rangle - \langle a, x \rangle) = 0 \\ &\quad (\text{because } \langle a, z \rangle = \langle a, x \rangle = b). \end{aligned}$$

By the previous theorem, z is a solution of $\min_{x \in S} \|x - y\|$.

It remains to show the uniqueness.

Suppose we have two solutions z_1 and z_2 . Then,

z_1 is a solution, $\Rightarrow \langle z_1 - y, z_1 - z_2 \rangle = 0$

z_2 is a solution, $\Rightarrow \langle z_2 - y, z_2 - z_1 \rangle = 0$

Taking difference leads to $\langle z_1 - z_2, z_1 - z_2 \rangle = 0$

$$\Rightarrow \|z_1 - z_2\|^2 = 0 \Rightarrow \underline{z_1 = z_2} \quad \text{⊗}.$$

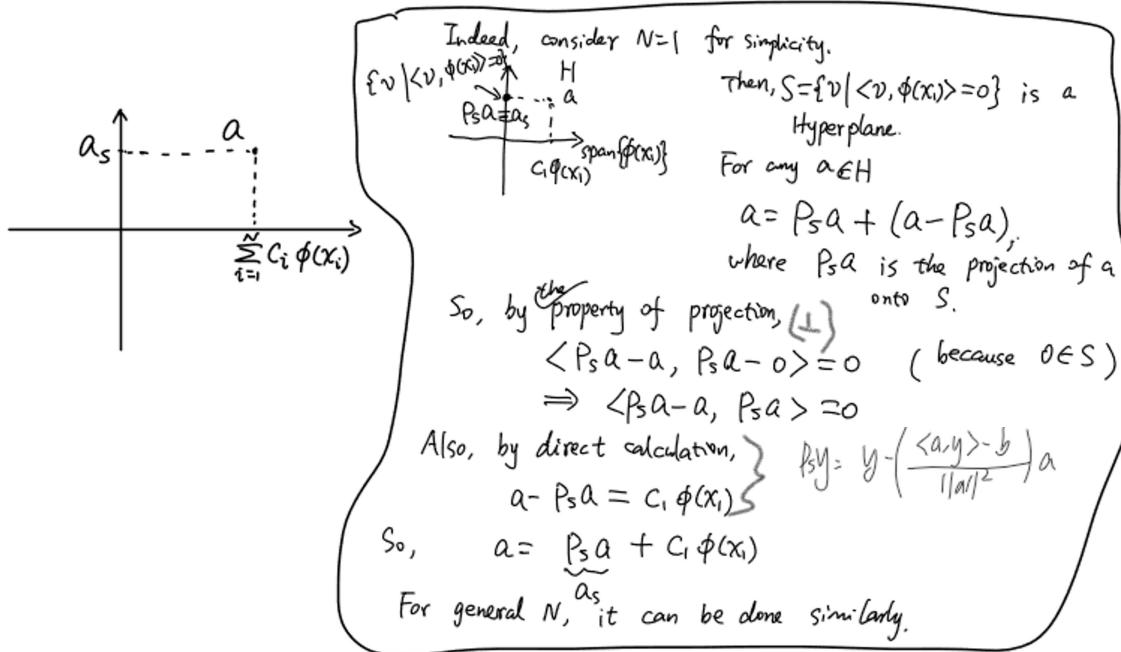
wadiu

21. Proof of representer theorem

Proof. For any $a \in H$, we claim that a can be decomposed as

$$a = a_s + \sum_{i=1}^N c_i \phi(x_i)$$

where $C = [c_1 \dots c_N] \in \mathbb{R}^N$ and $\langle a_s, \phi(x_i) \rangle = 0$ for $i=1, 2, \dots, N$.



解釋：

a 在 H space 上面

然後將 a 拆成兩個 component

一個係 linear combination of $\phi(x)$

一個係與其對應的 orthogonal vector

然後我 claim that $a =$ 兩個 component 總和

22. prove how to change objective function to inner product only in kernel regression problem

Therefore,

$$\begin{aligned}
 & \frac{1}{2} \sum_{i=1}^N (\langle a, \phi(x_i) \rangle - y_i)^2 + \lambda \|a\|_H^2 \\
 & = \frac{1}{2} \sum_{i=1}^N \left(\sum_{j=1}^N c_j \phi(x_j) + a_s \phi(x_i) - y_i \right)^2 + \lambda \left\| \sum_{i=1}^N c_i \phi(x_i) + a_s \right\|_H^2 \\
 & = \frac{1}{2} \sum_{i=1}^N \left(\sum_{j=1}^N c_j \langle \phi(x_j), \phi(x_i) \rangle - y_i \right)^2 + \\
 & \quad \lambda \left(\sum_{i=1}^N c_i \phi(x_i), \sum_{j=1}^N c_j \phi(x_j) \right) + 2 \langle a_s, \sum_{i=1}^N c_i \phi(x_i) \rangle + \langle a_s, a_s \rangle \\
 & = \frac{1}{2} \sum_{i=1}^N \left(\sum_{j=1}^N c_j K(x_i, x_j) - y_i \right)^2 + \lambda \sum_{i=1}^N \sum_{j=1}^N c_i c_j K(x_i, x_j) + \lambda \|a_s\|_H^2 \\
 & = \frac{1}{2} \|Kc - y\|_2^2 + \lambda c^T K c + \lambda \|a_s\|_H^2
 \end{aligned}$$

where $K = \begin{bmatrix} K(x_1, x_1) & K(x_1, x_2) & \dots & K(x_1, x_N) \\ \vdots & & & \\ K(x_N, x_1) & K(x_N, x_2) & \dots & K(x_N, x_N) \end{bmatrix} \in \mathbb{R}^{N \times N}$ $c = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_N \end{bmatrix} \in \mathbb{R}^N$.

Let $F_1(c) = \frac{1}{2} \|Kc - y\|_2^2 + \lambda c^T K c$ — depends on $c \in \mathbb{R}^N$ only.
 $F_2(a_s) = \lambda \|a_s\|_H^2$ — depends on $a_s \in H$ only.

Then, the minimization is the same as

$$\begin{aligned}
 & \min_{\substack{c \in \mathbb{R}^N \\ a_s \in H}} F_1(c) + F_2(a_s) \\
 & \langle a_s, \phi(x_i) \rangle = 0, i=1, \dots, N. \\
 & \bigcap a = a_s + \sum_{i=1}^N c_i \phi(x_i) \\
 & \min_{c \in \mathbb{R}^N} F_1(c) \quad \text{and} \quad \min_{\substack{a_s \in H \\ \langle a_s, \phi(x_i) \rangle = 0 \\ i=1, \dots, N}} F_2(a_s)
 \end{aligned}$$

Obviously, because $\|a_s\|_H^2 \geq 0$, $\min_{\substack{a_s \in H \\ \langle a_s, \phi(x_i) \rangle = 0}} F_2(a_s)$ is solved by $a_s = 0$.

Thus, the solution of the original minimization is

$$a = \sum_{i=1}^N c_i \phi(x_i)$$

where $c \in \mathbb{R}^N$ is a solution of $\min_{c \in \mathbb{R}^N} F_1(c)$. \square

Pf23. SVM-1的推論過程

Since x_+ is a projection of x_- onto $S_+ = \{x | \langle a, x \rangle = 1-b\}$

$$\begin{aligned}
 x_+ &= x_- - \frac{\langle a, x_- \rangle + b - 1}{\|a\|_2^2} a \\
 &= x_- - \frac{-1-b+b-1}{\|a\|_2^2} a \quad (\text{since } x_- \in S_-) \\
 &= x_- + \frac{2}{\|a\|_2^2} a
 \end{aligned}$$

因為 x 係 S -到, $\langle a, x \rangle$ 自然符合

Thus, $\|x_+ - x_-\|_2 = \left\| \frac{2}{\|a\|_2^2} a \right\|_2 = \frac{2}{\|a\|_2}$

S -定義
= -1-b

i.e., the margin = $\frac{2}{\|a\|_2}$

Pf24. SVM-2的邏輯推論

這裡的邏輯是這樣的：
 如果 $y_i(\langle a, x_i \rangle + b) - 1 \geq 0$, 那麼 $-e_i \leq 0$ 時 fulfill $e_i \geq 0$,
 $\min e_i$ 的話只能是 0. (e.g. -15)
 如果 $-e_i \leq y_i(\langle a, x_i \rangle + b) - 1 < 0$, $-e_i$ 要經過某個負數,
 同時 $e_i \geq 0$. $\Rightarrow e_i$ 只能是那負數之絕對值 (15)

Pf25: Kernel SVM 推論過程

First step of proof:

We know that: $\forall a \in H$, we can decompose it as
 $a = a_s + \sum_{i=1}^N c_i \phi(x_i)$, where $\langle a_s, \phi(x_i) \rangle = 0$
 $i=1, \dots, N$.

Then

$$\begin{aligned} F(a) &= \frac{\lambda}{2} \left(\|a_s\|_H^2 + \left\| \sum_{i=1}^N c_i \phi(x_i) \right\|_H^2 \right) \\ &\quad + \sum_{i=1}^N h(y_i \langle a_s + \sum_{j=1}^N c_j \phi(x_j), \phi(x_i) \rangle - 1) \\ &= \frac{\lambda}{2} \|a_s\|_H^2 + \left(\frac{\lambda}{2} \left\| \sum_{i=1}^N c_i \phi(x_i) \right\|_H^2 + \sum_{i=1}^N h(y_i \sum_{j=1}^N c_j \langle \phi(x_j), \phi(x_i) \rangle - 1) \right) \\ &= \frac{\lambda}{2} \|a_s\|_H^2 F_1(a_s) \Rightarrow a_s = 0 \\ &\quad + \left(\frac{\lambda}{2} \sum_{i=1}^N \sum_{j=1}^N c_i c_j \langle \phi(x_i), \phi(x_j) \rangle + \sum_{i=1}^N h(y_i \sum_{j=1}^N c_j \langle \phi(x_j), \phi(x_i) \rangle - 1) \right) \end{aligned}$$

$$c^T K c$$

$$h(y_i (Kc)_i - 1)$$

Then $\min_{a \in H} F(a)$

$$F_2(c)$$

$$\begin{array}{l} \min_{a \in H} F_1(a_s) + F_2(c) \Leftrightarrow \begin{cases} a_s = 0 \\ c \text{ is the solution of } \min_{c \in \mathbb{R}^N} F_2(c) \end{cases} \\ \text{subject to } \begin{cases} a_s \in H \\ \langle a_s, \phi(x_i) \rangle = 0 \\ i=1, \dots, N \end{cases} \end{array}$$

So, a solution of $\min_{a \in H} F(a)$ is

$$a = 0 + \sum_{i=1}^N c_i \phi(x_i), \text{ where } c \text{ is the solution of } \min_{c \in \mathbb{R}^N} F_2(c).$$

26. Proving if f is continuous and coercive, there exists at least one solution of unconstrained smooth optimization

Proof. Consider the set, known as the level set,

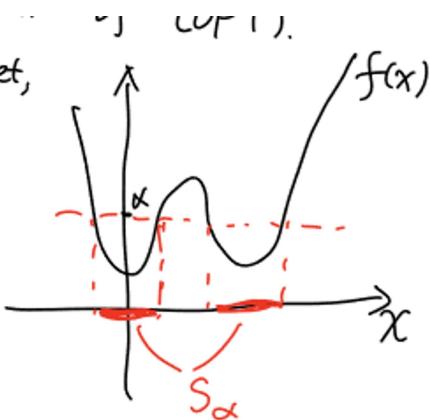
$$S_\alpha = \{x \mid f(x) \leq \alpha\}$$

S_α is closed, because of ①:

Let $x^{(n)} \rightarrow x$ and $\{x^{(n)}\} \subset S_\alpha$.

Then $f(x^{(n)}) \leq \alpha$

$$\text{① } \Rightarrow f(x) \leq \alpha \Rightarrow x \in S_\alpha$$



S_α is bounded, because of ②:

Suppose S_α is unbounded, i.e., $\exists \{x^{(n)}\} \subset S_\alpha$ s.t. $\|x^{(n)}\| \rightarrow \infty$.

Then ② $\Rightarrow f(x^{(n)}) \rightarrow +\infty$

but $f(x^{(n)}) \leq \alpha$

} contradiction.

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

We choose α s.t. S_α is non-empty. (For example, choose $\alpha = f(\omega)$)

By Weierstrass's theorem (any continuous function on a bounded and closed non-empty set must have a minimizer and maximizer)

So, $\min_{x \in S_\alpha} f(x)$ exists a solution, called x^* .

Then, $\forall x \in S_\alpha, f(x^*) \leq f(x)$.

$\forall x \notin S_\alpha, f(x^*) \leq \alpha \leq f(x)$

$x^* \in S_\alpha$



$$\left. \begin{array}{l} f(x^*) \leq f(x) \\ f(x^*) \leq \alpha \leq f(x) \end{array} \right\} \Rightarrow f(x^*) \leq f(x) \quad \forall x \in \mathbb{R}^n.$$

27. proof of first order condition in unconstrained smooth optimization problem

Proof. By definition of the gradient,

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

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

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

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

Because $x^{(*)}$ is fixed, $\|\nabla f(x^{(*)})\| \neq 0$ is a constant.

$$\lim_{t \rightarrow 0} \frac{t \|\nabla f(x^{(*)})\|^2}{t} = \|\nabla f(x^{(*)})\|^2 \neq 0$$

$$\lim_{t \rightarrow 0} \frac{o(t \|\nabla f(x^{(*)})\|)}{t} = 0$$

So, by choosing a small enough t ,

$$t \|\nabla f(x^{(*)})\|^2 > o(t \|\nabla f(x^{(*)})\|)$$

This implies $f(\tilde{x}) < f(x^{(*)})$ contradiction.

28. proof of the theorem: Convex + differentiable, grad = 0 \iff argmin

Proof. " \Rightarrow " is proved. (for any f , convex or not)

" \Leftarrow "

Since f is convex and differentiable,

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

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

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



Pf29: Theorem If a function is strictly convex, the solution is unique if it exists

Proof. Suppose there are two solutions $x^{(k)}, y^{(k)}$

$$\text{of } \min_{x \in \mathbb{R}^n} f(x) \quad (\Rightarrow f(x^{(k)}) = f(y^{(k)}))$$

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

Then $z \neq x^{(k)}, z \neq y^{(k)}$

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

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

$$= f(x^{(k)}) \Rightarrow f(x^{(k)}) > f(z)$$

$$x^{(k)}, y^{(k)} = \arg \min_{x \in \mathbb{R}^n} f(x) \Rightarrow f(x^{(k)}) \leq f(z)$$

f is strictly convex

convex but not strictly convex
It may have multiple solutions of $\min_{x \in \mathbb{R}^n} f(x)$

? contradiction \square

Pf30 It can be checked $F(x)$ is continuous, coercive, strictly convex

① F is continuous, and coercive.

— Continuity is obvious

— Coercivity:

$$\begin{aligned} F(x) &= f(x^{(k)}) + \frac{1}{\alpha_k} \left(\frac{1}{2} \|x - x^{(k)}\|_2^2 + \langle x - x^{(k)}, \alpha_k \nabla f(x^{(k)}) \rangle \right) \\ &= f(x^{(k)}) + \frac{1}{\alpha_k} \left(\frac{1}{2} \|x - x^{(k)} + \alpha_k \nabla f(x^{(k)})\|_2^2 - \frac{1}{2} \alpha_k^2 \|\nabla f(x^{(k)})\|_2^2 \right) \\ &= \frac{1}{2\alpha_k} \|x - (x^{(k)} - \alpha_k \nabla f(x^{(k)}))\|_2^2 - \frac{\alpha_k}{2} \|\nabla f(x^{(k)})\|_2^2 + f(x^{(k)}) \\ &\geq \frac{1}{2\alpha_k} \left(\|x\|_2 - \|x^{(k)} - \alpha_k \nabla f(x^{(k)})\|_2 \right)^2 - \frac{\alpha_k}{2} \|\nabla f(x^{(k)})\|_2^2 + f(x^{(k)}) \end{aligned}$$

Obviously, when $\|x\|_2 \rightarrow +\infty$, $F(x) \rightarrow +\infty$.

So, $\min_{x \in \mathbb{R}^n} F(x)$ has at least a solution.

② F is strictly convex,

because

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

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

↑
strictly convex

↑
Convex

is strictly convex in x .

So $F = \text{convex} + \text{strictly convex function}$

is strictly convex.

Thus, $\min_{x \in \mathbb{R}^n} F(x)$ has a unique solution

$x^{(k+1)} = \arg \min_{x \in \mathbb{R}^n} F(x)$ is well-defined

Pf31: proving gradient descent(cont'd)

3. Since F is convex and differentiable

$$x^{(k+1)} = \arg \min_{x \in \mathbb{R}^n} F(x) \iff \nabla F(x^{(k+1)}) = 0$$

Since

$$\nabla F(x) = \nabla f(x^{(k)}) + \frac{1}{\alpha_k} (x - x^{(k)})$$

$$\text{So, } \nabla F(x^{(k+1)}) = 0 \iff \nabla f(x^{(k)}) + \frac{1}{\alpha_k} (x^{(k+1)} - x^{(k)}) = 0$$

$$\iff \boxed{x^{(k+1)} = x^{(k)} - \alpha_k \nabla f(x^{(k)})}$$

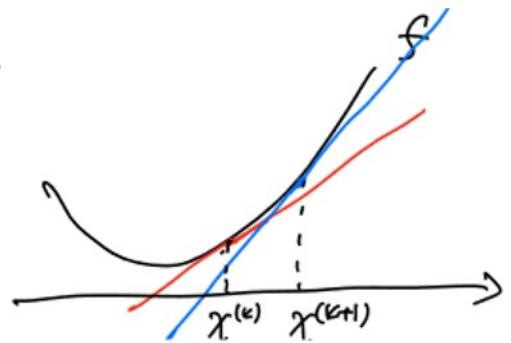
Pf32: Proving range of alpha needs to be (1, m)

For simplicity, we assume f is convex.

By convexity,

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

To give an upper bound of $f(x^{(k+1)})$



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

i.e.,

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

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

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

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

$$\stackrel{\downarrow}{\underbrace{}} \leq f(x^{(k)}) - \alpha \|\nabla f(x^{(k)})\|_2^2 + M \|x^{(k+1)} - x^{(k)}\|_2^2 \quad \alpha^2 \|\nabla f(x^{(k)})\|_2^2$$

Assume

$$\leq f(x^{(k)}) - \alpha(1-M\alpha) \|\nabla f(x^{(k)})\|_2^2$$

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

If we want $f(x^{(k+1)}) < f(x^{(k)})$, we require

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

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

Pf33. objective function for Least Squares is differentiable

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

- f is differentiable, because:

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

$$f(y) = \frac{1}{2} \|Ay - b\|_2^2 = \frac{1}{2} \|(Ax - b) + A(y - x)\|_2^2$$

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

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

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

So, if M is finite, then

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

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

Indeed, $M = \max_{\|z\|_2 \leq 1} \|Az\|_2^2 = \left(\max_{\|z\|_2 \leq 1} \|Az\|_2 \right)^2 = \|A\|_{2 \rightarrow 2}^2 \equiv \|A\|_2^2$

Because, $g(z) = \|Az\|_2 = \sqrt{\sum_i \left(\sum_j a_{ij} z_j \right)^2}$ is continuous

$\{z \mid \|z\|_2 \leq 1\}$ is closed, bounded, non-empty }

by Weierstrass thm, the maximum of $g(z)$ on $\{z \mid \|z\|_2 \leq 1\}$ must be finite.

Part 1: proving

$$\boxed{\frac{\partial F}{\partial w^{(j)}}(W, c) = \sum_{i=1}^m 2(f_{w,c}(x^{(i)}) - y_i) \left(\frac{\partial}{\partial w^{(j)}} f_{w,c}(x^{(i)}) \right)}$$

In order to implement GD, we need the gradient of F .

Since $F: \mathbb{R}^{(n+1)p} \rightarrow \mathbb{R}$ is a function on the Euclidean space,

$$\nabla F(W, c) = \begin{bmatrix} \frac{\partial F}{\partial w^{(1)}}(W, c) \\ \vdots \\ \frac{\partial F}{\partial w^{(p)}}(W, c) \\ \frac{\partial F}{\partial c}(W, c) \end{bmatrix}, \quad \text{where } \frac{\partial F}{\partial w^{(j)}} \text{ stands for}$$

the gradient of F
w.r.t. $w^{(j)}$ when the
others are fixed.

∇F

For $\frac{\partial F}{\partial w^{(j)}}$:

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

Let $g_1(t) = (t - y_i)^2$, $g_2(w^{(j)}) = f_{w,c}(x^{(i)})$

$$\begin{aligned} \text{Then } g_1: \mathbb{R} \rightarrow \mathbb{R}, \quad g_2: \mathbb{R}^n \rightarrow \mathbb{R} \quad \text{and } (g_1 \circ g_2)(w^{(j)}) &= g_1(g_2(w^{(j)})) \\ &= (f_{w,c}(x^{(i)}) - y_i)^2 \end{aligned}$$

By the chain rule:

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

So,

$$\boxed{\frac{\partial F}{\partial w^{(j)}}(W, c) = \sum_{i=1}^m 2(f_{w,c}(x^{(i)}) - y_i) \left(\frac{\partial}{\partial w^{(j)}} f_{w,c}(x^{(i)}) \right)}$$

Pf34 Part 2: Proving:

$$\boxed{\frac{\partial}{\partial w^{(j)}} f_{w,c}(x^{(i)}) = \nabla(g_4 \circ g_3)(w^{(i)}) = g'_4(g_3(w^{(i)})) \cdot \nabla g_3(w^{(i)})} \\ = c_j \cdot \sigma'(\langle w^{(i)}, x^{(i)} \rangle) \cdot x^{(i)}$$

- Since $f_{w,c}(x^{(i)}) = \langle c, \sigma(W^T x^{(i)}) \rangle$
- $= \sum_{k=1}^p c_k \sigma(\langle w^{(k)}, x^{(i)} \rangle)$
- $= c_j \sigma(\langle w^{(j)}, x^{(i)} \rangle) + \sum_{\substack{k=1 \\ k \neq j}}^p c_k \sigma(\langle w^{(k)}, x^{(i)} \rangle)$ (II)
- Let $g_3(w^{(j)}) = \langle w^{(j)}, x^{(i)} \rangle$ ($g_3: \mathbb{R}^n \rightarrow \mathbb{R}$)
- $g_4(t) = c_j \sigma(t) + A$ ($g_4: \mathbb{R} \rightarrow \mathbb{R}$)
- Also, $(g_4 \circ g_3)(w^{(i)}) = g_4(g_3(w^{(i)})) = f_{w,c}(x^{(i)})$

$$\boxed{\frac{\partial}{\partial w^{(j)}} f_{w,c}(x^{(i)}) = \nabla(g_4 \circ g_3)(w^{(i)}) = g'_4(g_3(w^{(i)})) \cdot \nabla g_3(w^{(i)})} \\ = c_j \cdot \sigma'(\langle w^{(j)}, x^{(i)} \rangle) \cdot x^{(i)}$$

Pf34. Part 3-

Altogether

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