

Homework 3

1.1

- 2.1) We prove by induction:

- Base case, $K=1$ we get that $\theta_1 x_1 \in C$
Since $\theta_1 = 1$ and $x_1 \in C$. ✓
- Assume that $\sum_{i=1}^K \theta_i x_i \in C$ with $\theta_i \geq 0$,
 $\sum_{i=1}^K \theta_i = 1$.
- Take $\sum_{i=1}^{K+1} \theta_i x_i$ with $\sum_{i=1}^{K+1} \theta_i = 1$, $\theta_i \geq 0$.

$$\text{Then } \sum_{i=1}^{K+1} \theta_i x_i = \sum_{i=1}^K \theta_i x_i + \theta_{K+1} x_{K+1}$$

$$(1) \quad = (1 - \theta_{K+1}) \sum_{i=1}^K \frac{\theta_i}{(1 - \theta_{K+1})} x_i + \theta_{K+1} x_{K+1}$$

$$\text{Now } \sum_{i=1}^K \frac{\theta_i}{1 - \theta_{K+1}} = 1 \quad \text{since } \sum_{i=1}^K \theta_i = 1 - \theta_{K+1}$$

and thus by induction hypothesis $\sum_{i=1}^K \frac{\theta_i}{1 - \theta_{K+1}} x_i \in C$.

$$\text{Say } \sum_{i=1}^K \frac{\theta_i}{1 + \theta_{K+1}} = x^* \in C.$$

Then (1) reduces to $(1 - \theta_{K+1}) x^* + \theta_{K+1} x_{K+1}$ which is in C by definition of convexity.

$$\text{Thus } \sum_{i=1}^{K+1} \theta_i x_i \in C \quad \square$$

2.5)

The distance between the two hyperplanes will be perpendicular to both hyperplanes.

Then take a point x_1 in the first hyperplane: $a^T x_1 = b_1$

Now from x_1 , we go in a direction perpendicular until we hit x_2 in the second hyperplane. So $x_2 = x_1 + ta$ for some $t \in \mathbb{R}$.

Plugging into the second hyperplane we get $a^T(x_1 + ta) = b_2$,

$$a^T x_1 + a^T a t = b_2, \quad t = \frac{b_2 - a^T x_1}{a^T a} = \frac{b_2 - b_1}{a^T a}$$

So $x_2 = x_1 + \frac{b_2 - b_1}{a^T a} a$ and thus the distance is

$$\|x_2 - x_1\| = \left\| \frac{b_2 - b_1}{a^T a} a \right\| = \frac{|b_2 - b_1|}{\|a\|^2} \|a\| = \frac{|b_2 - b_1|}{\|a\|} \quad \square$$

2.12) (a) This is the intersection of two (convex) half-spaces:

$$a^T x \leq \beta \quad \text{and} \quad a^T x \geq \alpha$$

And thus is convex.

(b) This is the intersection of $2n$ half-spaces:

$$x_i \leq \beta_i \quad \text{and} \quad x_i \geq \alpha_i \quad \text{for every } i.$$

And thus is convex.

(c) This is the intersection of two halfspaces:

$$a_1^T x \leq b_1 \quad \text{and} \quad a_2^T x \leq b_2$$

And thus is convex.

(d) First note that the set of points closer to a point than another point is a half-space:

- Geometrically take the perpendicular bisector of the two points and it's one side of it

- Algebraically note that:

$$\|x - x_0\| \leq \|x - x_1\|$$

$$\|x - x_0\|^2 \leq \|x - x_1\|^2$$

$$(x-x_0)^T(x-x_0) \leq (x-x_1)^T(x-x_1)$$

$$\cancel{x^T x} - 2x_0^T x + x_0^T x_0 \leq \cancel{x^T x} - 2x_1^T x + x_1^T x_1,$$

$$2(x_1 - x_0)^T x \leq x_1^T x_1 - x_0^T x_0$$

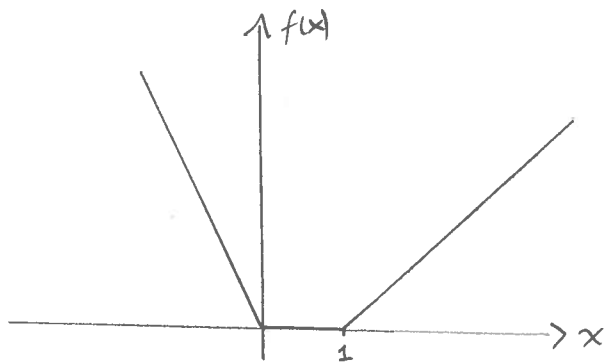
is a half-space (with $a = 2(x_1 - x_0)$, $b = x_1^T x_1 - x_0^T x_0$)

Thus the set of points closer to a given point than a set is an infinite intersection of half-spaces and thus is convex:

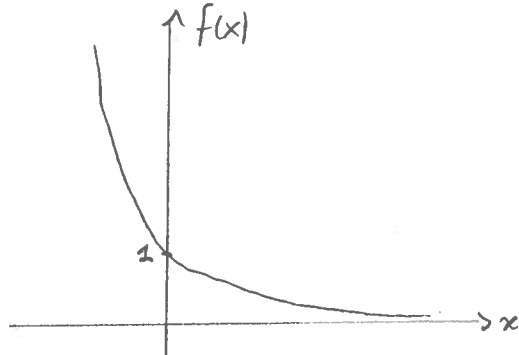
$$\bigcap_{y \in S} \{x \mid \|x - x_0\| \leq \|x - y\|\}$$

2.1

- (i) Take $f(x) = \max(-x, 0, x-1)$, then as seen in the graph below f has a unique minimum value at zero but multiple minimizers. Any $x \in [0, 1]$ is a minimizer:



- (ii) Take $f(x) = e^{-x}$ for $x \in \mathbb{R}$. Then f is convex since $f''(x) = e^{-x} > 0$ for all x , but f has no minimum value since we can make $f(x) = e^{-x}$ as close to zero as desired but never zero:



2.2

(i) For every γ , we have that:

$$\frac{\partial L(hw(x), \gamma)}{\partial x} = L'(hw(x), \gamma) \cdot w$$

$$\frac{\partial^2 L(hw(x), \gamma)}{\partial x^2} = L''(hw(x), \gamma) \cdot w^2$$

Using the property:

Also note that since $hw(x) = w \cdot x$ is an affine mapping and $L(h, \gamma)$ is convex, then $L(hw(x), \gamma)$ is convex by properties of convexity

Since $L(h, \gamma)$ is convex then $L''(hw(x), \gamma) \geq 0 \forall x, \gamma$ and thus $L''(hw(x), \gamma) \cdot w^2 \geq 0 \forall x, \gamma$ so $L(hw(x), \gamma)$ is convex.

(ii) $\hat{L}(w) = \frac{1}{M} \sum_{i=1}^M L(hw(x_i), \gamma_i) = \sum_{i=1}^M \frac{1}{M} L(hw(x_i), \gamma_i)$ is a weighted sum of convex functions and thus is convex.

3.1

To reduce the optimality gap by a factor of 10 we need:

$$c^k \leq \frac{1}{10} \Rightarrow \left(1 - \frac{1}{M}\right)^k \leq \frac{1}{10}$$

(for $C_f = 3$)

$$\left(\frac{2}{3}\right)^k \leq \frac{1}{10} \quad (\text{since the condition number } \frac{M}{M-1} = 3$$

$$k \geq \log_{\frac{2}{3}}\left(\frac{1}{10}\right) \quad \text{so } \frac{M}{M-1} = \frac{3}{2})$$

$$k \geq 5.6$$

So we need 6 iterations.

If the condition number is 100 instead, we get:

$$\left(1 - \frac{1}{100}\right)^k \leq \frac{1}{10}$$

$$\left(\frac{99}{100}\right)^k \leq \frac{1}{10}$$

$$k \geq \log_{\frac{99}{100}}\left(\frac{1}{10}\right)$$

$$k \geq 229.1$$

So we need 230 iterations.

4.1

$$\textcircled{1} \hat{L}(w) = \frac{1}{n} \sum_{i=1}^n (w - \gamma_i)^2 / 2$$

$$\hat{L}'(w) = \frac{1}{n} \sum_{i=1}^n (w - \gamma_i)$$

setting $\hat{L}'(w)$ to zero we get: $\sum_{i=1}^n (w^* - \gamma_i) = 0$

$$nw^* = \sum_{i=1}^n \gamma_i$$

$$w^* = \frac{1}{n} \sum_{i=1}^n \gamma_i = \bar{\gamma}$$

$$\textcircled{2} \text{ From above } \hat{L}'(w) = \frac{1}{n} \sum_{i=1}^n (w - \gamma_i)$$

$$= w - \frac{1}{n} \sum_{i=1}^n \gamma_i$$

$$= w - \bar{\gamma}$$

$\textcircled{3}$ With learning rate $h=1$, initial $w=w_0$ we get that:

$$w_1 = w_0 - \hat{L}'(w_0) = w_0 - (w_0 - \bar{\gamma}) = \bar{\gamma}$$

so $w_1 = \bar{\gamma} = w^*$ and GD converges in one step.

$\textcircled{4}$ With learning rate $h=1/2$ and initial $w=w_0$ we get that:

$$w_1 = w_0 - \frac{1}{2} \hat{L}'(w_0) = w_0 - \frac{1}{2}(w_0 - \bar{\gamma}) = \frac{1}{2}w_0 + \frac{1}{2}\bar{\gamma}$$

$$(1) \Rightarrow w_1 - \bar{\gamma} = \frac{1}{2}(w_0 - \bar{\gamma})$$

$$\text{And } w_{k+1} = w_k - \frac{1}{2} \hat{L}'(w_k) = w_k - \frac{1}{2}(w_k - \bar{\gamma}) = \frac{1}{2}w_k + \frac{1}{2}\bar{\gamma}$$

$$(2) \Rightarrow w_{k+1} - \bar{\gamma} = \frac{1}{2}(w_k - \bar{\gamma})$$

From (1) and (2) by induction we get that $(w_k - \bar{\gamma}) = (\frac{1}{2})^k (w_0 - \bar{\gamma})$

(which implies $|w_k - \bar{\gamma}| \leq (\frac{1}{2})^k |w_0 - \bar{\gamma}|$)

4.2

$$M \hat{L}'(\omega_0) = 0 - 1 - 1 - 1 - 1 - 1 = -5$$

$$6 \hat{L}'(-3) = -5 \quad (\text{note that in this case } M=6)$$

$$\hat{L}'(-3) = -\frac{5}{6}$$

$$\text{Then } \omega_1 = \omega_0 - 6 \hat{L}'(\omega_0) = -3 - 6\left(-\frac{5}{6}\right) = 2 \quad (1 \text{ iteration})$$

$$M \hat{L}'(\omega_1) = +1 + 1 + 1 + 1 + 0.5 - 1 = 3.5$$

$$6 \hat{L}'(2) = 3.5$$

$$\hat{L}'(2) = \frac{7}{12}$$

$$\text{Then } \omega_2 = \omega_1 - 6 \hat{L}'(\omega_1) = 2 - 6\left(\frac{7}{12}\right) = -\frac{3}{2} \quad (2 \text{ iterations})$$

$$M \hat{L}'(\omega_2) = +1 + 0.5 - 1 - 1 - 1 - 1 = -2.5$$

$$6 \hat{L}'\left(-\frac{3}{2}\right) = -\frac{5}{2}$$

$$\hat{L}'\left(-\frac{3}{2}\right) = -\frac{5}{12}$$

$$\text{Then } \omega_3 = \omega_2 - 6 \hat{L}'(\omega_2) = -\frac{3}{2} - 6\left(-\frac{5}{12}\right) = 1 \quad (3 \text{ iterations})$$

4.3

$$\hat{L}(\omega) = \frac{1}{M} \sum_{i=1}^M (\omega \cdot x_i - \omega^* \cdot x_i)^2 / 2$$

$$\frac{\partial \hat{L}(\omega)}{\partial \omega_j} = \frac{1}{M} \sum_{i=1}^M (\omega \cdot x_i - \omega^* \cdot x_i) (x_i)_j$$

$$= \frac{1}{M} \sum_{i=1}^M (\omega - \omega^*) \cdot x_i (x_i)_j$$

$$\text{Hence } \nabla \hat{L}(\omega) = \frac{1}{M} \sum_{i=1}^M (\omega - \omega^*) \cdot x_i x_i$$

To express this as a matrix equation: $\nabla \hat{L}(\omega) = H(\omega - \omega^*)$ we

See that the j^{th} component of $H(\omega - \omega^*)$ is the j^{th} row of H dotted with $(\omega - \omega^*)$.

Comparing this with $\nabla \hat{L}(\omega) = \frac{1}{N} \sum_{i=1}^N (\omega - \omega^*) \cdot x_i x_i$ we see that the j^{th} component of $\nabla \hat{L}(\omega)$ is $\frac{1}{N} \sum_{i=1}^N (\omega - \omega^*) \cdot x_i (x_i)_j$
 $= \frac{1}{N} (\omega - \omega^*) \cdot \sum_{i=1}^N x_i (x_i)_j$

Thus the j^{th} row of H must be $\frac{1}{N} \sum_{i=1}^N (x_i)_j x_i^T$ (*) but this is the j^{th} row of $\frac{1}{N} X^T X$ so $H = \frac{1}{N} X^T X$

From * we can also see that if the l^{th} row of H is $\frac{1}{N} \sum_{i=1}^N (x_i)_l x_i^T$ then $H_{kl} = \frac{1}{N} \sum_{i=1}^N (x_i)_l (x_i)_k = \frac{1}{N} \sum_{i=1}^N (x_i)_k (x_i)_l$

For $N=4$ and $x_i = (1, i)$ we have that:

$$H = \frac{1}{4} X^T X = \frac{1}{4} \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{pmatrix}$$

$$H = \frac{1}{4} \begin{pmatrix} 4 & 10 \\ 10 & 30 \end{pmatrix}$$

□