

1 [Q1]

Last week, we talked about the gradient and the hessian. Gradient is the direction which the function descends fastest. And we also talked about convergence analysis of convex optimization. Especially, we talked about gradient descent and analyzed 2 special cases: M -smooth and strong convexity.

2 [Q2]

done

3 [Q3]

3.1 (a)

Assume $f(\mathbf{x}), f(\mathbf{y})$ are convex, then we have,

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

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

Combine the 2 inequalities together and we can get,

$$\begin{aligned} f(\mathbf{x}) + f(\mathbf{y}) &\geq f(\mathbf{y}) + f(\mathbf{x}) + \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{y}) \rangle + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle \\ - \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{y}) \rangle - \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle &\geq 0 \\ \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{y}) \rangle + \langle \mathbf{y} - \mathbf{x}, -\nabla f(\mathbf{x}) \rangle &\geq 0 \\ \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{y}) - \nabla f(\mathbf{x}) \rangle &\geq 0 \end{aligned}$$

3.2 (b)

Since ∇f is monotone, let's define

$$\phi(\alpha) = f(\mathbf{x} + \alpha(\mathbf{y} - \mathbf{x}))$$

Then, we have $\phi'(\alpha) \geq \phi'(0)$ for all $\alpha > 0$ and $\alpha \in \text{dom}\phi$ where,

$$\phi'(\alpha) = \nabla f(\mathbf{x} + \alpha(\mathbf{y} - \mathbf{x}))^T (\mathbf{y} - \mathbf{x})$$

Hence,

$$\begin{aligned} f(\mathbf{y}) = \phi(1) &= \phi(0) + \int_0^1 \phi'(\alpha) d\alpha \\ &\geq \phi(0) + \phi'(0) \\ &= f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) \end{aligned}$$

3.3 (c)

If $x > x^*$, $f'(x) > 0$ and if $x < x^*$, $f'(x) < 0$.

4 [Q4]

4.1 (a)

- In first part, I showed that (ii) \rightarrow (i). since, $g(\mathbf{x})$ is convex in $\mathbf{dom} f$, we can have,

$$\begin{aligned}
 g(\mathbf{y}) &\geq g(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle \\
 f(\mathbf{y}) - \frac{m}{2} \|\mathbf{y}\|_2^2 &\geq f(\mathbf{x}) - \frac{m}{2} \|\mathbf{x}\|_2^2 + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) - m\mathbf{x} \rangle \\
 f(\mathbf{y}) &\geq f(\mathbf{x}) + \frac{m}{2} (\|\mathbf{y} - \mathbf{x}\|_2^2) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle - m \langle \mathbf{y} - \mathbf{x}, \mathbf{x} \rangle \\
 f(\mathbf{y}) &\geq f(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle + \frac{m}{2} \langle \mathbf{y} - \mathbf{x}, \mathbf{y} - \mathbf{x} \rangle \\
 f(\mathbf{y}) &\geq f(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle + \frac{m}{2} \|\mathbf{y} - \mathbf{x}\|_2^2
 \end{aligned}$$

- In second part, I showed that (i) \rightarrow (ii).

$$\begin{aligned}
 f(\mathbf{y}) &\geq f(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle + \frac{m}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 \\
 f(\mathbf{y}) &\geq f(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle + \frac{m}{2} \langle \mathbf{y} - \mathbf{x}, \mathbf{y} - \mathbf{x} \rangle \\
 f(\mathbf{y}) &\geq f(\mathbf{x}) + \frac{m}{2} (\|\mathbf{y} - \mathbf{x}\|_2^2) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle - m \langle \mathbf{y} - \mathbf{x}, \mathbf{x} \rangle \\
 f(\mathbf{y}) - \frac{m}{2} \|\mathbf{y}\|_2^2 &\geq f(\mathbf{x}) - \frac{m}{2} \|\mathbf{x}\|_2^2 + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) - m\mathbf{x} \rangle
 \end{aligned}$$

If we define $g(\mathbf{x}) = f(\mathbf{x}) - \frac{m}{2} \|\mathbf{x}\|_2^2$, then we can have,

$$g(\mathbf{y}) \geq g(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle$$

4.2 (b)

- In first part, I showed that (ii) \rightarrow (iii).

Since $g(\mathbf{x})$ is convex, we can have,

$$f(\mathbf{y}) - \frac{m}{2} \|\mathbf{y}\|_2^2 \geq f(\mathbf{x}) - \frac{m}{2} \|\mathbf{x}\|_2^2 + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) - m\mathbf{x} \rangle$$

Simiarly, we can have,

$$f(\mathbf{x}) - \frac{m}{2} \|\mathbf{x}\|_2^2 \geq f(\mathbf{y}) - \frac{m}{2} \|\mathbf{y}\|_2^2 + \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{y}) - m\mathbf{y} \rangle$$

Combine the inequalities together, we can have that,

$$\begin{aligned}
 \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) - m\mathbf{x} \rangle + \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{y}) - m\mathbf{y} \rangle &\leq 0 \\
 \langle \mathbf{x} - \mathbf{y}, (\nabla f(\mathbf{x}) - m\mathbf{x}) - (\nabla f(\mathbf{y}) - m\mathbf{y}) \rangle &\geq 0 \\
 \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle &\geq \langle \mathbf{x} - \mathbf{y}, m(\mathbf{x} - \mathbf{y}) \rangle \\
 \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle &\geq m \|\mathbf{x} - \mathbf{y}\|_2^2
 \end{aligned}$$

- In second part, I showed that (iii) \rightarrow (ii).

$$\begin{aligned}
 \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle &\geq m \|\mathbf{x} - \mathbf{y}\|_2^2 \\
 \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle &\geq \langle \mathbf{x} - \mathbf{y}, m(\mathbf{x} - \mathbf{y}) \rangle \\
 \langle \mathbf{x} - \mathbf{y}, (\nabla f(\mathbf{x}) - m\mathbf{x}) - (\nabla f(\mathbf{y}) - m\mathbf{y}) \rangle &\geq 0 \\
 \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) - m\mathbf{x} \rangle + \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{y}) - m\mathbf{y} \rangle &\leq 0
 \end{aligned}$$

Hence, we can have that,

$$f(\mathbf{y}) - \frac{m}{2}\|\mathbf{y}\|_2^2 + f(\mathbf{x}) - \frac{m}{2}\|\mathbf{x}\|_2^2 \geq f(\mathbf{x}) - \frac{m}{2}\|\mathbf{x}\|_2^2 + f(\mathbf{y}) - \frac{m}{2}\|\mathbf{y}\|_2^2 + \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{y}) - m\mathbf{y} \rangle$$

Apparently, it indicates 2 inequalities, which follows,

$$f(\mathbf{y}) - \frac{m}{2}\|\mathbf{y}\|_2^2 \geq f(\mathbf{x}) - \frac{m}{2}\|\mathbf{x}\|_2^2 + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) - m\mathbf{x} \rangle$$

And the second inequality just swaps x with y . And $g(\mathbf{x})$ is convex.

5 [Q5]

5.1 (a)

According to Cauchy-Schwartz inequality, we have,

$$\begin{aligned} \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle &\leq \|\mathbf{x} - \mathbf{y}\|_2 * \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \\ &\leq \|\mathbf{x} - \mathbf{y}\|_2 * M\|\mathbf{x} - \mathbf{y}\|_2 \\ &\leq M\|\mathbf{x} - \mathbf{y}\|_2^2 \end{aligned}$$

5.2 (b)

$$\langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle \leq M\|\mathbf{x} - \mathbf{y}\|_2^2 \quad (1)$$

$$\langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) - M\mathbf{x} + M\mathbf{y} \rangle \leq 0 \quad (2)$$

If we define $g(\mathbf{x}) = \frac{M}{2}\|\mathbf{x}\|_2^2 - f(\mathbf{x})$, then we can have,

$$\nabla g(\mathbf{x}) = M\mathbf{x} - \nabla f(\mathbf{x})$$

And from (2) we can achieve,

$$\langle \mathbf{x} - \mathbf{y}, \nabla g(\mathbf{y}) - \nabla g(\mathbf{x}) \rangle \leq 0$$

$$\langle \mathbf{y} - \mathbf{x}, \nabla g(\mathbf{y}) - \nabla g(\mathbf{x}) \rangle \geq 0$$

5.3 (c)

Since $g(\mathbf{x})$ is convex, we can have,

$$g(\mathbf{y}) \geq g(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \nabla g(\mathbf{x}) \rangle \quad (3)$$

$$\frac{M}{2}\|\mathbf{y}\|_2^2 - f(\mathbf{y}) \geq \frac{M}{2}\|\mathbf{x}\|_2^2 - f(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, M\mathbf{x} - \nabla f(\mathbf{x}) \rangle \quad (4)$$

$$f(\mathbf{y}) \leq \frac{M}{2}\|\mathbf{y}\|_2^2 - \frac{M}{2}\|\mathbf{x}\|_2^2 + f(\mathbf{x}) - \langle \mathbf{y} - \mathbf{x}, M\mathbf{x} \rangle + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle \quad (5)$$

$$f(\mathbf{y}) \leq \frac{M}{2} \langle \mathbf{y} - \mathbf{x}, \mathbf{y} + \mathbf{x} \rangle - M \langle \mathbf{y} - \mathbf{x}, \mathbf{x} \rangle + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle + f(\mathbf{x}) \quad (6)$$

$$f(\mathbf{y}) \leq \frac{M}{2}\|\mathbf{y} - \mathbf{x}\|_2^2 + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle + f(\mathbf{x}) \quad (7)$$

5.4 (d)

since $f(x)$ is convex, we can have,

$$f(\mathbf{x}) + \langle \mathbf{z} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle \leq f(\mathbf{z}) \quad (8)$$

Besides, (iv) tells us,

$$f(\mathbf{z}) \leq f(\mathbf{y}) + \langle \mathbf{z} - \mathbf{y}, \nabla f(\mathbf{y}) \rangle + \frac{M}{2} \|\mathbf{z} - \mathbf{y}\|_2^2 \quad (9)$$

If we combine the 2 inequalities (8)(9) together, then we can have,

$$f(\mathbf{x}) - f(\mathbf{y}) \leq \langle \mathbf{x} - \mathbf{z}, \nabla f(\mathbf{x}) \rangle + \langle \mathbf{z} - \mathbf{y}, \nabla f(\mathbf{y}) \rangle + \frac{M}{2} \|\mathbf{z} - \mathbf{y}\|_2^2 \quad (10)$$

If we take the gradient with respect to \mathbf{z} , then we can have,

$$\mathbf{z} = \frac{\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})}{M} + \mathbf{y} \quad (11)$$

Plug this back into (10) expression and we can get an upper bound for $f(\mathbf{x}) - f(\mathbf{y})$,

$$\begin{aligned} (*) &\leq \left\langle \mathbf{x} - \frac{\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})}{M} - \mathbf{y}, \nabla f(\mathbf{x}) \right\rangle + \left\langle \frac{\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})}{M} + \mathbf{y}, \nabla f(\mathbf{y}) \right\rangle \\ &\quad + \left\| \frac{\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})}{M} \right\|_2^2 \\ &\leq \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) \rangle - \frac{1}{M} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2 + \frac{1}{2M} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2 \\ &\leq \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) \rangle - \frac{1}{2M} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2 \end{aligned}$$

Hence, we can have,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle + \frac{1}{2M} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2$$

5.5 (e)

(v) tells us,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle + \frac{1}{2M} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2 \quad (12)$$

Simiarly, we can have,

$$f(\mathbf{x}) \geq f(\mathbf{y}) + \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{y}) \rangle + \frac{1}{2M} \|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|_2^2 \quad (13)$$

Combine the 2 inequalities (12) (13) together and we can have,

$$\langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle \geq \frac{1}{M} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2$$

5.6 (f)

$$\begin{aligned} \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle &\geq \frac{1}{M} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2 \\ M \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle &\geq \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2 \\ M \|\mathbf{x} - \mathbf{y}\|_2 &\geq \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \end{aligned}$$

6 [Q6]

6.1 (a)

Since we have $\{x_k\}_{k=0}^{\infty}$ which satisfies $|x_{k+1} - x^*| \leq \beta|x_k - x^*|$, we can get,

$$\begin{aligned} |x_1 - x^*| &\leq \beta|x_0 - x^*| \\ |x_2 - x^*| &\leq \beta|x_1 - x^*| \\ &\vdots \\ |x_k - x^*| &\leq \beta|x_{k-1} - x^*| \end{aligned}$$

If we combine these inequalities together we can have,

$$|x_k - x^*| \leq \beta^k|x_0 - x^*| \quad (14)$$

Besides, we know that,

$$k \geq \frac{\log\left(\frac{|x_0 - x^*|}{\epsilon}\right)}{\log\left(\frac{1}{\beta}\right)}$$

we can achieve in addition,

$$\frac{1}{\beta^k} \geq \frac{|x_0 - x^*|}{\epsilon} \quad (15)$$

$$\epsilon \geq \beta^k|x_0 - x^*| \quad (16)$$

Put (16) back to (14) we can achieve,

$$|x_k - x^*| \leq \epsilon$$

6.2 (b)

$$\lim_{k \rightarrow \infty} \frac{x_{k+1} - x^*}{x_k - x^*} = \lim_{k \rightarrow \infty} \frac{k}{k+1} \quad (17)$$

$$= \beta = 1 \quad (18)$$

Hence, it's sublinear.

Besides, we can get

$$\begin{aligned} |x_1 - x^*| &= \frac{1}{2}|x_0 - x^*| \\ |x_2 - x^*| &= \frac{2}{3}|x_1 - x^*| \\ &\vdots \\ |x_k - x^*| &= \frac{k}{k+1}|x_{k-1} - x^*| \end{aligned}$$

Combine these inequalities together, we can have,

$$|x_k - x^*| \leq \frac{1}{k+1} |x_0 - x^*|$$

To make sure $|x_k - x^*| \leq \epsilon$, we can get,

$$\frac{1}{k+1} |x_0 - x^*| \geq \epsilon \quad (19)$$

$$k \leq \frac{|x_0 - x^*|}{\epsilon} - 1 \quad (20)$$

since k is integer, I think it needs to round up.

$$k \leq \lceil \frac{|x_0 - x^*|}{\epsilon} - 1 \rceil$$

6.3 (c)

For $q = 1$, since $x_k = \frac{1}{2^k}$, x^* must be 0 when $k \rightarrow \infty$, thus we can have,

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - x^*|}{|x_k - x^*|} = \lim_{k \rightarrow \infty} \frac{\frac{1}{2^{k+1}}}{\frac{1}{2^k}} \quad (21)$$

$$= \frac{1}{2} \quad (22)$$

Similarly, for $q = 2$, $z_k = \frac{1}{2^{2k}}$, x^* must be 0 when $k \rightarrow \infty$, and we can have,

$$\begin{aligned} \lim_{k \rightarrow \infty} \frac{|z_{k+1} - z^*|}{|z_k - z^*|^2} &= \frac{\frac{1}{2^{2(k+1)}}}{\left(\frac{1}{2^{2k}}\right)^2} \\ &= 1 \end{aligned}$$

7 [Q7]

7.1 (a)

Since the monotone gradient property of convex functions, the function is monotonically increasing or decreasing, bisection in the optimal solution, the gradient of convex function is 0. Only when x_k is close enough to optimal solution x^* can the iteration stop. Else, the gradient is either greater than 0 or less than 0, choosing the midpoint can shorten interval. And finally, it's guaranteed that the the optimal solution or quasi-optimal solution will be found.

7.2 (b)

In first step, we have,

$$|x_0 - x^*| \leq \frac{x_u - x_l}{2}$$

If we use bisection algorithm, we compute and have $x_1, x_{u1}, x_{l1}, x_2, x_{u2}, x_{l2}, \dots$, then,

$$|x_k - x^*| \leq \frac{x_{uk} - x_{lk}}{2}$$

Since the interval is halved at each step, we have

$$x_{uk} - x_{lk} = \frac{x_{u(k-1)} - x_{l(k-1)}}{2} = \dots = \frac{x_u - x_l}{2^{k+1}}$$

Hence, we achieve

$$|x_k - x^*| \leq \frac{x_u - x_l}{2^{k+1}}$$

Thus, the conclusion here is the 'maximum error' possible between x_k and x^* decreases as k increases, which means x_k converges to x^* .

7.3 (c)

Apparently it converges linearly. The convergence rate β is

$$\beta = \lim_{k \rightarrow \infty} \frac{|x_k - x^*|}{|x_{k-1} - x^*|} \quad (23)$$

$$= \lim_{k \rightarrow \infty} \frac{\frac{x_{uk} - x_{lk}}{2}}{\frac{x_{u(k-1)} - x_{l(k-1)}}{2}} \quad (24)$$

$$= \lim_{k \rightarrow \infty} \frac{x_{uk} - x_{lk}}{x_{u(k-1)} - x_{l(k-1)}} \quad (25)$$

Since the interval is halved at each step, we have

$$\beta = \lim_{k \rightarrow \infty} \frac{x_{uk} - x_{lk}}{x_{u(k-1)} - x_{l(k-1)}} = \frac{1}{2}$$

7.4 (d)

First we need to find the bound. Note that in convexity, any local minimum is also a global minimum. Hence we can find the bound x_l, x_r which satisfies $f'(x_l) > 0$ and $f'(x_r) < 0$. Hence, we can make sure the bound is $[x_l, x_r]$. Then we can treat it as an constrained optimization problem and use the method above to solve it.

8 [Q8]

8.1 (a)

if $g(\mathbf{a}^T \mathbf{x} + b) \geq \frac{1}{2}$, we can have

$$\frac{1}{1 + e^{-(\mathbf{a}^T \mathbf{x} + b)}} \geq \frac{1}{2} \quad (26)$$

$$e^{-(\mathbf{a}^T \mathbf{x} + b)} \leq 1 \quad (27)$$

$$\mathbf{a}^T \mathbf{x} + b \geq 0 \quad (28)$$

And vice versa, namely we can deduce $g(\mathbf{a}^T \mathbf{x} + b) \geq \frac{1}{2}$ from $\mathbf{a}^T \mathbf{x} + b \geq 0$. Thus, we can conclude,

$$g(\mathbf{a}^T \mathbf{x} + b) \geq \frac{1}{2} \iff \mathbf{a}^T \mathbf{x} + b \geq 0$$

Similarly, we can have

$$g(\mathbf{a}^T \mathbf{x} + b) \leq \frac{1}{2} \iff \mathbf{a}^T \mathbf{x} + b \leq 0$$

In conclusion, $h(x)$ can equivalently be written as

$$h(x) = \begin{cases} 1 & \text{if } \mathbf{a}^T \mathbf{x} + b \geq 0 \\ 0 & \text{if } \mathbf{a}^T \mathbf{x} + b < 0 \end{cases}$$

8.2 (b)

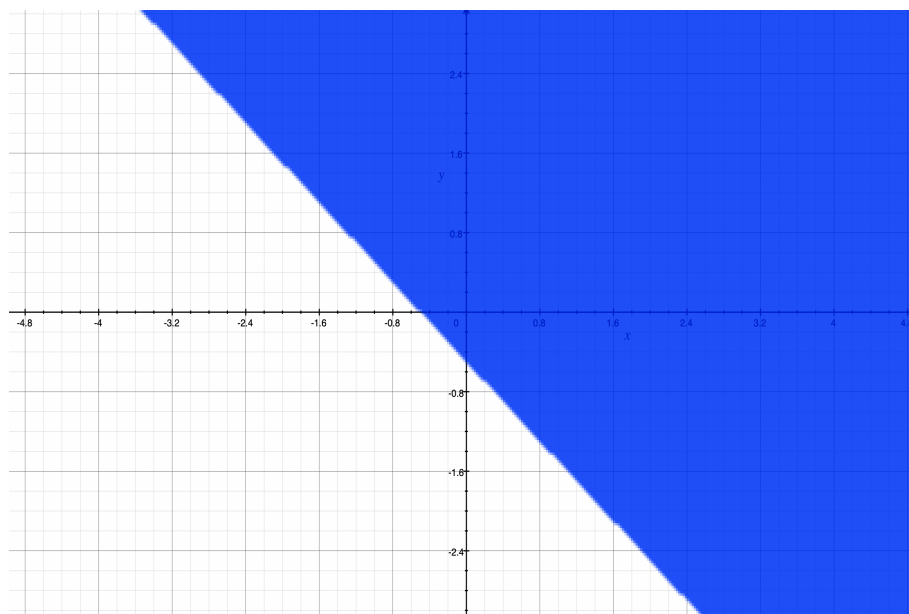


Figure 1: Illustration of Set

Hence, the blue region is where the set lies

8.3 (c)

$$\sum_{n=1}^N y_n \log(g(\mathbf{a}^T \mathbf{x} + b)) + (1 - y_n) \log(1 - g(\mathbf{a}^T \mathbf{x} + b)) \quad (29)$$

$$= \sum_{n=1}^N y_n \log\left(\frac{1}{1 + e^{-(\mathbf{a}^T \mathbf{x} + b)}}\right) + (1 - y_n) \log\left(1 - \frac{1}{1 + e^{-(\mathbf{a}^T \mathbf{x} + b)}}\right) \quad (30)$$

$$= \sum_{n=1}^N -y_n \log\left(1 + e^{-(\mathbf{a}^T \mathbf{x} + b)}\right) + (1 - y_n) \log\left(\frac{e^{-(\mathbf{a}^T \mathbf{x} + b)}}{1 + e^{-(\mathbf{a}^T \mathbf{x} + b)}}\right) \quad (31)$$

$$= \sum_{n=1}^N -y_n \log\left(1 + e^{-(\mathbf{a}^T \mathbf{x} + b)}\right) + \log\left(\frac{e^{-(\mathbf{a}^T \mathbf{x} + b)}}{1 + e^{-(\mathbf{a}^T \mathbf{x} + b)}}\right) - y_n \log\left(\frac{e^{-(\mathbf{a}^T \mathbf{x} + b)}}{1 + e^{-(\mathbf{a}^T \mathbf{x} + b)}}\right) \quad (32)$$

$$= \sum_{n=1}^N -y_n \log\left[\left(1 + e^{-(\mathbf{a}^T \mathbf{x} + b)}\right) \left(\frac{e^{-(\mathbf{a}^T \mathbf{x} + b)}}{1 + e^{-(\mathbf{a}^T \mathbf{x} + b)}}\right)\right] + \log\left(\frac{1}{1 + e^{\mathbf{a}^T \mathbf{x} + b}}\right) \quad (33)$$

$$= \sum_{n=1}^N -y_n \log\left(e^{-(\mathbf{a}^T \mathbf{x} + b)}\right) - \log\left(1 + e^{\mathbf{a}^T \mathbf{x} + b}\right) \quad (34)$$

$$= \sum_{n=1}^N y_n (\mathbf{a}^T \mathbf{x} + b) - \log\left(1 + e^{\mathbf{a}^T \mathbf{x} + b}\right) \quad (35)$$

Hence, the origin problem is equivalent to

$$\underset{\mathbf{a}, b}{\text{maximize}} \sum_{n=1}^N y_n (\mathbf{a}^T \mathbf{x} + b) - \log\left(1 + e^{\mathbf{a}^T \mathbf{x} + b}\right)$$

8.4 (d)

$$\nabla f(\boldsymbol{\theta}) = \sum_{n=1}^N \left(y_n \tilde{\mathbf{x}}_n - \frac{e^{\boldsymbol{\theta}^T \tilde{\mathbf{x}}_n}}{1 + e^{\boldsymbol{\theta}^T \tilde{\mathbf{x}}_n}} \tilde{\mathbf{x}}_n \right) \quad (36)$$

$$= \sum_{n=1}^N \tilde{\mathbf{x}}_n \left(y_n - \frac{e^{\boldsymbol{\theta}^T \tilde{\mathbf{x}}_n}}{1 + e^{\boldsymbol{\theta}^T \tilde{\mathbf{x}}_n}} \right) \quad (37)$$

$$\nabla^2 f(\boldsymbol{\theta}) = - \sum_{n=1}^N \frac{e^{-\boldsymbol{\theta}^T \tilde{\mathbf{x}}_n}}{(1 + e^{-\boldsymbol{\theta}^T \tilde{\mathbf{x}}_n})^2} \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^T$$

8.5 (e)

It's M -smooth for some finite value of M and strongly convex.

8.6 (f)

Here, I set the stopping criteria as the maximum value of θ_k **is less than 0.01**.
In details, I code it in python as below. The first part is to load data.

```
import numpy as np
from sklearn import datasets
import math

np.random.seed(2020) # Set random seed so results are repeatable
x,y = datasets.make_blobs(n_samples=100,n_features=2,centers=2,
                           cluster_std=6.0)
```

The second part is to define a function in gradient so that i can make my life easier.

```
# define logit function
def logit(x):
    return 1./(1+np.exp(-x))
```

The third part is to iteratively find the optimal value

```
# add 1 for every row of x
x_tilde = np.ones([100,3])
for i in range(len(x)):
    x_tilde[i] = np.r_[x[i], np.ones(1)]

X_train = np.matrix(x_tilde) #100*3
y_train = np.matrix(y).T #100*1

# initialize theta 3*1
theta = np.zeros(X_train.shape[1])
theta = np.matrix(theta).T

k=0
dk = np.ones((3,1))
alpha = 0.001

# iteration number
while max(abs(dk)) > 0.01:
    error = y_train - logit(X_train * theta) #100*1
    dk = X_train.T * error #3*1
    theta += alpha * dk
    k = k + 1
    print("current dk is {} and iteration is {}".format(dk.T, k))
print("theta is {}".format(theta.T))
```

I found in this condition, the it breaks the loop in **2695**. As is shown in the picture,

```
current dk is [[-0.00089327 -0.00122307  0.01033231]] and iteration is 2678
current dk is [[-0.0008915  -0.00122065  0.01031184]] and iteration is 2679
current dk is [[-0.00088973 -0.00121824  0.0102914  ]] and iteration is 2680
current dk is [[-0.00088797 -0.00121583  0.01027101]] and iteration is 2681
current dk is [[-0.00088621 -0.00121342  0.01025066]] and iteration is 2682
current dk is [[-0.00088446 -0.00121102  0.01023034]] and iteration is 2683
current dk is [[-0.0008827  -0.00120862  0.01021007]] and iteration is 2684
current dk is [[-0.00088096 -0.00120623  0.01018984]] and iteration is 2685
current dk is [[-0.00087921 -0.00120385  0.01016965]] and iteration is 2686
current dk is [[-0.00087747 -0.00120146  0.0101495  ]] and iteration is 2687
current dk is [[-0.00087573 -0.00119909  0.01012939]] and iteration is 2688
current dk is [[-0.000874  -0.00119671  0.01010932]] and iteration is 2689
current dk is [[-0.00087227 -0.00119435  0.01008928]] and iteration is 2690
current dk is [[-0.00087054 -0.00119198  0.01006929]] and iteration is 2691
current dk is [[-0.00086881 -0.00118962  0.01004934]] and iteration is 2692
current dk is [[-0.00086709 -0.00118727  0.01002943]] and iteration is 2693
current dk is [[-0.00086538 -0.00118492  0.01000956]] and iteration is 2694
current dk is [[-0.00086366 -0.00118258  0.00998973]] and iteration is 2695
theta is [[-0.28047652 -0.45702889  2.20888147]]
```

Figure 2: Illustration of Code

And the final optimal d_k , namely gradient of $f(\theta)$ is $\begin{bmatrix} -0.00086366 & -0.00118258 & 0.00998973 \end{bmatrix}$ and corresponding iteration is 2695. And the final θ is $\begin{bmatrix} -0.28047652 & -0.45702889 & 2.20888147 \end{bmatrix}$

8.7 (g)

Here, in outer loop I set the stopping criteria as the maximum absolute value of θ_k is less than 0.01.

In inner loop, I set the stopping criteria as $\alpha_l - \alpha_r > 1e - 10$. In details, I code it in python as below. The first and second part is the same as before.

```
import numpy as np
from sklearn import datasets
import math

np.random.seed(2020) # Set random seed so results are repeatable
x,y = datasets.make_blobs(n_samples=100,n_features=2,centers=2,
                           cluster_std=6.0)

# define logit function
def logit(x):
    return 1./(1+np.exp(-x))
```

And the third part is to iteratively find the optimal value. There are 2 loops. The Inner loop is to find optimal α in each iteration while the outer loop is to find the optimal solution.

```
# add 1 for every row of x
x_tilde = np.ones([100,3])
for i in range(len(x)):
    x_tilde[i] = np.r_[x[i], np.ones(1)]

X_train = np.matrix(x_tilde) #100*3
y_train = np.matrix(y).T #100*1

# initialize theta 3*1
theta = np.zeros(X_train.shape[1])
theta = np.matrix(theta).T

k=0
dk = np.ones((3,1))

alpha_left = 0
alpha_right = 1
# iteration number
while max(abs(dk)) > 0.01:

    m = 0
    alpha_left = 0
    alpha_right = 1
    h_prime = 1
    while abs(alpha_left - alpha_right) > 1e-10:
        alpha = (alpha_left + alpha_right) / 2
        h_prime = (X_train.T * (- y_train + logit(X_train * (theta +
                                                                alpha * dk)))) .T * dk

        if h_prime > 0:
            alpha_right = alpha
        else:
            alpha_left = alpha
        m += 1
    print("optimal alpha is {} and iteration is {}".format(alpha,
                                                            m))

    error = y_train - logit(X_train * theta) #100*1
    dk = X_train.T * error #3*1
    theta += alpha * dk
```

```

k = k + 1
print("current dk is {} and iteration is {}".format(dk.T, k))
print("theta is {} and total iteration is {}".format(theta.T, m*k))

```

The final results are shown in picture below,

```

optimal alpha is 7.2479248046875e-05 and iteration is 18
optimal alpha is 7.43865966796875e-05 and iteration is 19
optimal alpha is 7.534027099609375e-05 and iteration is 20
optimal alpha is 7.486343383789062e-05 and iteration is 21
optimal alpha is 7.510185241699219e-05 and iteration is 22
optimal alpha is 7.49826431274414e-05 and iteration is 23
optimal alpha is 7.492303848266602e-05 and iteration is 24
optimal alpha is 7.495284080505371e-05 and iteration is 25
optimal alpha is 7.496774196624756e-05 and iteration is 26
optimal alpha is 7.497519254684448e-05 and iteration is 27
optimal alpha is 7.497146725654602e-05 and iteration is 28
optimal alpha is 7.496960461139679e-05 and iteration is 29
optimal alpha is 7.49705359339714e-05 and iteration is 30
optimal alpha is 7.49700702726841e-05 and iteration is 31
optimal alpha is 7.497030310332775e-05 and iteration is 32
optimal alpha is 7.497018668800592e-05 and iteration is 33
optimal alpha is 7.497012848034501e-05 and iteration is 34
current dk is [[-0.00987779 -0.00068754  0.00972636]] and iteration is 233
theta is [[-0.28043516 -0.45700234  2.20864798]] and total iteration is 7922

```

Figure 3: Illustration of Code

When using this strategy to find the “optimal” step size at each iteration, 34 iterations are needed. And total number of iterations taken by the combined bisection searches is 7922. And θ is $\begin{bmatrix} -0.28043516 & -0.45700234 & 2.20864798 \end{bmatrix}$. I don’t think it’s worth the additional costs since it has more iterations.