

1.1 Convexity

$$1. f(\omega) = \frac{1}{2} (\omega - y)^T V (\omega - y) + \frac{\lambda}{2} \|\omega\|^2$$

$$= \frac{1}{2} (\omega^T V X \omega - y^T V X \omega - \omega^T X^T y + y^T y) + \frac{\lambda}{2} \omega^T I \omega$$

$$\therefore \nabla f = X^T V X + \lambda I$$

- for $\omega \in \mathbb{R}^d$, $\omega^T V X \omega = \|V^{1/2} X \omega\|_2^2 \geq 0$, $\omega^T \lambda I \omega = \lambda \|\omega\|_2^2 \geq 0$

$$\therefore \nabla^2 f \cdot \omega > 0 \therefore \nabla^2 f = X^T V X + \lambda I \text{ is positive semidefinite} \therefore f \text{ is convex}$$

$$2. f(\omega) = -y^T X \omega + l^T V, \quad v_i = \exp(\omega^T x_i)$$

$$= -y^T X \omega + \sum_{i=1}^n \exp(\omega^T x_i)$$

$$\nabla(-y^T X \omega) = 0, \text{ Let } g(\omega) = \sum_{i=1}^n \exp(\omega^T x_i)$$

$$\frac{\partial g}{\partial \omega_j} = \sum_{i=1}^n \exp(\omega^T x_i) \cdot x_j^i$$

$$\therefore \frac{\partial g}{\partial \omega_i \partial \omega_j} = \sum_{i=1}^n \exp(\omega^T x_i) \cdot x_j^i \cdot x_k^i \quad \therefore \nabla^2 f = \nabla^2 g = X^T \begin{pmatrix} e^{\omega^T x_1} & 0 & \dots & 0 \\ 0 & e^{\omega^T x_2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \dots & e^{\omega^T x_n} \end{pmatrix} \cdot X$$

- let V be the diagonal matrix of size $n \times n$, with $V_{ii} = e^{\omega^T x_i}$. Similar to question 1, since each $V_{ii} = e^{\omega^T x_i} > 0$, $\nabla^2 f = X^T V X$ is positive semidefinite $\therefore f$ is convex

$$3. f(\omega) = \max_{j \in \{1, \dots, d\}} L_j |\omega_j|^{l+1} = \alpha l + \alpha(l+1) \geq \alpha = (l + V(l+1))$$

- for $v, w \in \mathbb{R}^d$, we examine $f((1-t)v + tw)$ for $t \in [0, 1]$

$$- v = (v_1, \dots, v_d), w = (w_1, \dots, w_d) \quad \alpha \geq (wt + v(l+1))^l \quad (1, \alpha) \geq V$$

$\forall j \in \{1, \dots, d\}$,

$$L_j |(1-t)v_j + tw_j|$$

$$\leq L_j [(1-t)|v_j| + t|w_j|] \quad (\text{triangle inequality} + \text{property of norm})$$

$$= L_j (1-t)|v_j| + t \cdot L_j |w_j|$$

$$\therefore \max_{j \in \{1, \dots, d\}} L_j |(1-t)v_j + tw_j| \leq \max_{j \in \{1, \dots, d\}} L_j (1-t)|v_j| + t \cdot L_j |w_j|$$

$$\leq (1-t) \max_{j \in \{1, \dots, d\}} L_j |v_j| + t \cdot \max_{k \in \{1, \dots, d\}} L_k |w_k|$$

$$\therefore f((1-t)v + tw) \leq (1-t)f(v) + tf(w) \quad \forall t \in [0, 1] \therefore f \text{ is convex}$$

4. $f(w) = \|xw - y\|_p + \lambda \|Aw\|_q$
- function $\|\cdot\|_p$ (norm) are convex for $1 \leq p \leq \infty$. $xw - y$, Aw are affine maps
 - convex function combine with affine map is convex: $\|xw - y\|_p$ & $\|Aw\|_q$ are convex
 - multiplication by positive constant keeps convexity: $\lambda \|Aw\|_q$ is convex
 - sum of convex functions are convex $\therefore f(w) = \|xw - y\|_p + \lambda \|Aw\|_q$ is convex

5. $f(w) = \sum_{i=1}^n \max\{0, (w^T x_i - y_i) - \varepsilon\} + \frac{\lambda}{2} \|w\|_2^2$
- norm squared is convex, and multiplication of positive constant is convex: $\frac{\lambda}{2} \|w\|_2^2$ is convex
 - $w^T x_i - y_i$ is affine, and absolute value is convex. Composition of convex function with linear affine map is convex $\therefore (w^T x_i - y_i) - \varepsilon$ is convex
 - subtracting a constant remains convex, and 0 is convex, maximum of 2 fractions is convex $\therefore \max\{0, (w^T x_i - y_i) - \varepsilon\}$ is convex
 - sum of convex function is convex $\therefore f(w) = \sum_{i=1}^n \max\{0, (w^T x_i - y_i) - \varepsilon\} + \frac{\lambda}{2} \|w\|_2^2$ is convex

6. note that $\{w \in \mathbb{R}^d \mid Aw \leq b\} \subset C$ is a convex set. we use the convention that $\infty - \infty = \infty$.
- we show that f is convex by cases. let $v, w \in \mathbb{R}^d$
 - Case 1: $v, w \in C$. since $Aw \leq b$ $\Rightarrow Av \leq b$ $\Rightarrow Aw \leq b$ $\Rightarrow Av \leq b$ $\Rightarrow Aw \leq b$
 - $\therefore \forall t \in [0, 1], (1-t)v + tw \in C$
 - $\therefore f((1-t)v + tw) = \infty \leq -(1-t)\cdot 0 + t \cdot 0 = (1-t)f(v) + t f(w)$
- Case 2: one of v, w is in C^c . Suppose, wlog, $v \in C^c$, $w \in C$.
 $\forall t \in (0, 1) \quad f((1-t)v + tw) \leq \infty = (1-t)f(v) + t f(w) \quad (\because f_w = \infty, t > 0 \Rightarrow t f_w = \infty)$

- Case 3: $v, w \in C^c$
- $\forall t \in [0, 1], f((1-t)v + tw) \leq \infty = (1-t)f(v) + t f(w) + 1 (f(v) = f(w) = \infty)$
 - $\therefore f$ is convex.

$$(w^T v + t(v^T w))_{\min} \geq (w^T v + t(v^T w))_{\min}$$

$$(\min(w^T v + t(v^T w)))_{\min} \geq (\min(w^T v + t(v^T w)))_{\min}$$

$$\min(w^T v + t(v^T w)) \geq \min(w^T v + t(v^T w))$$

1.2 Convergence of Gradient Descent

$$\min_{k \in \{1, \dots, t\}} \|\nabla f(\omega^k)\|^2 \leq \frac{2L(f(\omega^0) - f^*)}{t^{\frac{4}{3}}} \leq \epsilon$$

$$\therefore \frac{2L(f(\omega^0) - f^*)}{\epsilon} \leq t^{\frac{4}{3}} \quad \therefore \frac{(2L(f(\omega^0) - f^*))^{\frac{3}{4}}}{\epsilon^{\frac{3}{4}}} \leq t, \quad \boxed{t = O\left(\frac{1}{\epsilon^{\frac{3}{4}}}\right) \text{ iterations}}$$

2. Note that by the algorithm, $L^0 \leq L^1 \leq \dots \leq L^k$

- we claim that $L^k \leq 2L$.

- suppose, for a contradiction, that $L^k > 2L$. Then \exists a minimum $t \in \{1, \dots, k\}$, $L^t > 2L$

$$\therefore L^{t-1} = \frac{L^t}{2^i} \text{ for some } i \in \mathbb{N}, \text{ and } L^{t-1} \leq 2L.$$

There are two cases:

Case 1: $i = 1$. Then $L^{t-1} = \frac{L^t}{2} > L$. But f is L -Lipschitz $\Rightarrow L^{t-1}$ -Lipschitz,

$$\forall w \in \mathbb{R}^d, f(w - \frac{1}{L^{t-1}} \nabla f(w)) \leq f(w) - \frac{1}{2L^{t-1}} \|\nabla f(w)\|^2 \quad (\star)$$

\therefore there is no need to double L^{t-1} for $w^t \therefore L^t = L^{t-1}$, contradiction.

Case 2: $i \geq 2$.

- if $L \leq L^{t-1} \leq 2L$, then similar as above, there is no need to double L^{t-1} for $w^t \therefore L^t = L^{t-1}$, contradiction.

- if $L^{t-1} < L$, then we may need to double L^{t-1} multiply times (to get L^t for the inequality to hold for w^t). But $\frac{L^{t-1}}{2^{i-1}} L^{t-1} = \frac{L^t}{2} > L$, and f is $2^{i-1} L^{t-1}$ -Lipschitz

\therefore the algorithm will produce an $L^t \leq 2^{i-1} L^{t-1}$, but $L^t = 2^i \cdot L^{t-1}$, contradiction.

$$- \forall w \in \mathbb{R}^d, \exists \mu > 0, \frac{1}{2} \|\nabla f(w)\|^2 \geq \mu (f(w) - f^*(w^*))$$

$$\forall k \in \mathbb{N}: -f(w^k) \leq f(w^{k-1}) - \frac{1}{2L^{k-1}} \|\nabla f(w^{k-1})\|^2 \leq f(w^{k-1}) - \frac{\mu}{L^{k-1}} (f(w^{k-1}) - f(w^*))$$

$$\therefore f(w^k) \leq (1 - \frac{\mu}{L^{k-1}}) f(w^{k-1}) + \frac{\mu}{L^{k-1}} f(w^*) \Rightarrow f(w^k) - f(w^*) \leq (1 - \frac{\mu}{L^{k-1}})(f(w^{k-1}) - f(w^*))$$

$$\therefore f(w^k) - f(w^*) \leq \left[\prod_{j=0}^{k-1} \left(1 - \frac{\mu}{L^j}\right) \right] (f(w^0) - f(w^*))$$

$$\therefore \forall j \in \{0, 1, \dots, k-1\}, L^j \leq L^k \leq 2L \quad \therefore \frac{\mu}{L^j} \geq \frac{\mu}{2L}, \quad \therefore 1 - \frac{\mu}{L^j} \leq 1 - \frac{\mu}{2L}.$$

$$\therefore f(w^k) - f(w^*) \leq \left[\prod_{j=0}^{k-1} \left(1 - \frac{\mu}{2L}\right) \right] (f(w^0) - f(w^*))$$

$$= \left(1 - \frac{\mu}{2L}\right)^k (f(w^0) - f(w^*)).$$

$$3. L^k = \frac{1}{2} L^{k-1}.$$

- Let $C^k = \{w \in \mathbb{R}^d \mid f(w) - \leq f(w^{k-1})\}$, then if f restricted to the set C^k is Lipschitz continuous with a smaller L^k (\Leftrightarrow composed f to L where $f: \mathbb{R}^d \rightarrow \mathbb{R}$), then decreasing L^k will increase the step size.

- This is effective when the step gradient descent algorithm is approaching the global minimum of the function, and we want to increase the step size for the algorithm to converge to the global minimum faster.

$$4. f(v) = f(w) + \nabla f(w)^T(v-w) + \frac{1}{2} (v-w)^T \nabla^2 f(w) (v-w) \text{ for some } u \text{ between } v \text{ and } w$$

- by strong convexity, $d^T \nabla^2 f(u) d \geq \mu \|d\|^2$

$$\therefore f(v) \geq f(w) + \nabla f(w)^T(v-w) + \frac{\mu}{2} \|v-w\|^2$$

- let $w = w^* \therefore \nabla f(w^*) = 0$

$$\therefore f(v) \geq f(w^*) + \frac{\mu}{2} \|v-w^*\|^2$$

- in class we have shown that $(f(w^k) - f(w^*)) \leq (1 - \frac{\mu}{L})^k (f(w^0) - f(w^*))$,

$$\rho = 1 - \frac{\mu}{L}$$

- from (*), we have that $f(w^k) - f(w^*) \geq \frac{\mu}{2} \|w^k - w^*\|^2$

$$\therefore \frac{\mu}{2} \|w^k - w^*\|^2 \leq (1 - \frac{\mu}{L})^k (f(w^0) - f(w^*))$$

$$\|w^k - w^*\| \leq \left(\frac{2}{\mu}\right)^{\frac{1}{2}} \left(1 - \frac{\mu}{L}\right)^{\frac{k}{2}} (f(w^0) - f(w^*))^{\frac{1}{2}}$$

- thus $\|w^k - w^*\|$ has a convergence rate of $O((1 - \frac{\mu}{L})^{\frac{k}{2}}) = O(\rho^{\frac{k}{2}})$.

1.3 Beyond Gradient Descent

$$\begin{aligned}
 1. & \frac{1}{2} \|v - w^{k+1}\|^2 + \alpha_k r(v) \\
 &= \frac{1}{2} \|v - w^k + \alpha_k \nabla f(w^k)\|^2 + \alpha_k r(v) \\
 &= \frac{1}{2} (v^T - w^{kT} + \alpha_k \nabla f(w^k)^T)(v - w^k + \alpha_k \nabla f(w^k)) + \alpha_k r(v) \\
 &= \frac{1}{2} (v^T - w^{kT})(v - w^k) + \alpha_k \nabla f(w^k)^T (v - w^k) + \frac{1}{2} \alpha_k^2 \|\nabla f(w^k)\|^2 + \alpha_k r(v) \\
 &= \cancel{\frac{1}{2\alpha_k} (\frac{1}{2} \|v - w^k\|^2 + \nabla f(w^k)^T (v - w^k) + r(v))} + \frac{1}{2} \alpha_k^2 \|\nabla f(w^k)\|^2
 \end{aligned}$$

- since multiplying by a constant (α_k) and adding a constant ($\frac{1}{2} \alpha_k^2 \|\nabla f(w^k)\|^2$ does not depend on v) does not change the minimizer,

$$w^{k+1} \in \underset{v \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \frac{1}{2\alpha_k} \|v - w^k\|^2 + \nabla f(w^k)^T (v - w^k) + r(v) + f(w^k) \right\}$$

$$2. f(w) = \sum_{i=1}^n \sum_{c \neq y_i} [1 - w_{y_i}^T x_i + w_c^T x_i]^+ + \frac{\lambda}{2} \|w\|_F^2$$

Notation: w_e^t : t^{th} weight of for class e ($t \in \mathbb{N}, 1 \leq t \leq d$).

$$- 1_A(x) = \begin{cases} 1 & \text{if } x \in A \\ 0 & \text{if } x \notin A \end{cases}, \delta_{c,e} = \begin{cases} 1 & \text{if } c = e \\ 0 & \text{if } c \neq e \end{cases}$$

- $\theta_{i,c}^{e,t} \in [0, 1]$ used for to describe subgradient when normal gradient does not exist.

$$\begin{aligned}
 \therefore \frac{\partial f}{\partial w_e^t} &= \sum_{i=1}^n \sum_{c \neq y_i} \left[1_{\{z>0\}} (1 - w_{y_i}^T x_i + w_c^T x_i) \cdot (-x_i^T \delta_{e,y_i} + x_i^T \delta_{e,c}) \right. \\
 &\quad \left. + 1_{\{z=0\}} (1 - w_{y_i}^T x_i + w_c^T x_i) \cdot \theta_{i,c}^{e,t} \cdot (-x_i^T \delta_{e,y_i} + x_i^T \delta_{e,c}) \right] + \lambda w_e^t
 \end{aligned}$$

$$3. f(v) \leq f(u) + \nabla f(u)^T(v-u) + \frac{L_\infty}{2} \|v-u\|_\infty^2.$$

$$\begin{aligned} \therefore f(w^{k+1}) &\leq f(w^k) - \nabla f(w^k)^T \cdot \frac{\|\nabla f(w^k)\|_1}{L_\infty} \cdot \text{sign}(\nabla f(w^k)) + \frac{L_\infty}{2} \left\| \frac{\nabla f(w^k)}{L_\infty} \right\|_1 \cdot \text{sign}(\nabla f(w^k)) \\ &= f(w^k) - \frac{\|\nabla f(w^k)\|_1}{L_\infty} \cdot \nabla f(w^k)^T \cdot \text{sign}(\nabla f(w^k)) + \frac{\|\nabla f(w^k)\|_1^2}{2L_\infty} \cdot \frac{\|\text{sign}(\nabla f(w^k))\|_\infty^2}{L_\infty}, \\ &\leq f(w^k) - \frac{\|\nabla f(w^k)\|_1}{L_\infty} \cdot \|\nabla f(w^k)\|_1 + \frac{\|\nabla f(w^k)\|_1^2}{2L_\infty}. \quad \leq 1 \\ &= f(w^k) - \frac{\|\nabla f(w^k)\|_1^2}{2L_\infty}. \end{aligned}$$

$$\therefore \|\nabla f(w^k)\|_1^2 \leq 2L_\infty (f(w^k) - f(w^{k+1}))$$

$$\therefore t \cdot \min_{k=0, \dots, t-1} \|\nabla f(w^k)\|_1^2 \leq 2L_\infty (f(w^0) - f(w^t))$$

If $\|\nabla f(w^k)\|_1^2 \leq \varepsilon$ to be guaranteed, then $\frac{2L_\infty (f(w^0) - f(w^t))}{t} \leq \varepsilon$.

$$\therefore \frac{2L_\infty (f(w^0) - f(w^t))}{\varepsilon} \leq t \quad \therefore \|\nabla f(w^k)\|_1^2 \leq \varepsilon \text{ after } O\left(\frac{1}{\varepsilon}\right) \text{ iterations.}$$