

### 3.5.2.1 - חשיבות

(1) יהי  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  קמורה,  $A \in \mathbb{R}^{n \times n}$ ,  $b \in \mathbb{R}^n$  ונגדיר  $g(x) = f(Ax+b)$

נניח כי  $g$  קמורה. יהי  $\lambda \in [0,1]$  ונניח  $w_1, w_2 \in \mathbb{R}^n$

$$\begin{aligned} g(\lambda w_1 + (1-\lambda)w_2) &\stackrel{\text{הנחה}}{=} f(A(\lambda w_1 + (1-\lambda)w_2) + b) = \\ &= f(\lambda(Aw_1 + b) + (1-\lambda)(Aw_2 + b)) = \\ &= f(\lambda(Aw_1 + b) + (1-\lambda)(Aw_2 + b)) \leq \\ &\stackrel{\text{קמוריות}}{\leq} \lambda \cdot f(Aw_1 + b) + (1-\lambda) f(Aw_2 + b) = \\ &\stackrel{\text{הנחה}}{=} \lambda \cdot g(w_1) + (1-\lambda) g(w_2) \end{aligned}$$

לכן  $g$  קמורה. המעבר מהקמוריות של  $f$  לקמוריות של  $g$ .

(2)  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  פונקציה קמורה.  $f_1, \dots, f_m$  פונקציות קמורות.  $g(x) = \max_i f_i(x)$

נניח כי  $g$  קמורה. יהי  $\lambda \in [0,1]$  ונניח  $w_1, w_2 \in \mathbb{R}^d$

נניח  $i$  קובע את האינדקס של  $f_i$  שמתאים ל- $g$  בנקודה  $x$ . כלומר  $f_i(x) = g(x)$

נניח  $j$  קובע את האינדקס של  $f_j$  שמתאים ל- $g$  בנקודה  $\lambda w_1 + (1-\lambda)w_2$ . כלומר  $f_j(\lambda w_1 + (1-\lambda)w_2) = g(\lambda w_1 + (1-\lambda)w_2)$

$\Rightarrow \max_i f_i(\lambda w_1 + (1-\lambda)w_2) \leq \max_i \lambda f_i(w_1) + \max_i (1-\lambda) f_i(w_2)$

$\Rightarrow g(\lambda w_1 + (1-\lambda)w_2) \leq \lambda g(w_1) + (1-\lambda) g(w_2)$

לכן  $g$  היא קמורה.

(3)  $e_{\log}(z) = \log_2(1+e^{-z})$   $e_{\log}: \mathbb{R} \rightarrow \mathbb{R}$  נניח כי היא קמורה.

נניח  $x \in \mathbb{R}$  ונניח  $g''(x) > 0$  (כלומר  $g$  קמורה). נניח  $g(x) = e_{\log}(x)$

$$e'_{\log}(z) = \frac{-e^{-z}}{\ln 2 (1+e^{-z})}$$

$$e''_{\log}(z) = \frac{e^{-z}(\ln 2(1+e^{-z})) + e^{-z}(-\ln 2 e^{-z})}{(\ln 2(1+e^{-z}))^2} =$$

$$= \frac{e^{-z}(\ln 2 + \ln 2 e^{-z} - \ln 2 e^{-z})}{(\ln 2(1+e^{-z}))^2} = \frac{\frac{1}{\ln 2} e^{-z}}{(\ln 2(1+e^{-z}))^2} > 0$$

לכן  $e_{\log}$  היא קמורה.



# המשפט

(1) יהי  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$  ו-  $f(\omega) = \log_2(\gamma \omega \cdot x)$  עבור  $x \in \mathbb{R}^2$

$$f(\omega) = \log_2(1 + e^{-\gamma \omega \cdot x}) = \sum_{i=1}^n \log_2(1 + e^{-\gamma \omega \cdot x_i})$$

נגדיר כעת  $t_i = \log_2(1 + e^{-\gamma \omega \cdot x_i})$  עבור  $i=1, \dots, n$

כאשר  $f(\omega) = \sum_{i=1}^n t_i(\omega)$ . הרעיון כי נשתמש בקטרים ונבין  $f$  קטרים.

יהי  $\lambda \in [0, 1]$  ו-  $x, y \in \mathbb{R}^2$  נקבע כי

~~הערה~~

$$f(\omega) = \sum_{i=1}^n t_i(\omega) \leq n \cdot \max_i t_i(\omega) \leq n \cdot \max_i t_i(\lambda x + (1-\lambda)y) \leq$$

$$\leq n \cdot \max_i \{ \lambda t_i(x) + (1-\lambda) t_i(y) \} = \lambda \max_i t_i(x) + (1-\lambda) \max_i t_i(y)$$

וכן  $f$  קטרי ביחס  $\omega$ .



# סמינר חישובי

(2)

כדי הרמז תחילה ~~ההצגה~~ hinge loss פונקציה  $\ell_{20}$  היא loss 0-1.

$$\min_{w \in W} \ell(w) = \sum_{i=1}^n \ell_{20}(\text{sign}(g(x_i, w)), y_i) = \sum_{i=1}^n \ell_{20}(y_i \cdot g(x_i, w))$$

ההצגה של loss 0-1 מהרה צורה:

$$\ell_{20} = \begin{cases} 1 & r \leq 0 \\ 0 & r \geq 0 \end{cases}$$

אם  $r \geq 0$  אז  $\ell_{20}(r) = 0$  אחרת  $\ell_{20}(r) = 1$

$$\ell_{\text{hinge}}(r) = \max\{0, 1-r\}$$

$$\ell_{\text{hinge}}(r) = 0 \geq 0 = \ell_{20}(r)$$

אם  $r \geq 1$  נקט:

$$\ell_{\text{hinge}}(r) = 1-r > 0 = \ell_{20}(r)$$

אם  $0 \leq r < 1$  נקט:

$$\ell_{\text{hinge}}(r) = 1+r > 1 \geq \ell_{20}(r)$$

אם  $r < 0$  נקט:

$$\ell_{\text{hinge}}(r) \geq \ell_{20}(r) \quad \forall r \in \mathbb{R}$$

נניח כי ניתן להבניה  $y_i \cdot w^* \cdot x_i \geq 0$   $\forall i \in [n]$

$$\text{sign}(w^* \cdot x_i) = y_i \quad \forall w^* \in \mathbb{R}^d$$

נרצה להראות  $w^*_{\text{hinge}}$  הוא נוסף אולי  $w^*$

אם  $\ell_{20}$  הוא loss 0-1

כיוון שאולי  $w^*_{\text{hinge}}$  ~~ההצגה~~  $\ell_{\text{hinge}}$  פונקציה  $\ell_{\text{hinge}}$  אולי  $w^*_{\text{hinge}}$

$$\sum_{i=1}^n \ell_{\text{hinge}}(y_i \cdot w^*_{\text{hinge}} \cdot x_i) = 0$$

וביון  $\ell_{\text{hinge}}(r) \geq \ell_{20}(r)$  נקט  $\ell_{\text{hinge}}(r) = 0$   $\forall i \in [n]$

$$\ell_{20}(y_i \cdot w^*_{\text{hinge}} \cdot x_i) = 0 \Rightarrow y_i \cdot w^*_{\text{hinge}} \cdot x_i \geq 1 > 0$$

אם  $\ell_{20}(r) = 0$  אז  $r \geq 0$

$$y_i = 1 \Rightarrow y_i \cdot w^*_{\text{hinge}} \cdot x_i \geq 1 \Rightarrow w^*_{\text{hinge}} \cdot x_i \geq 1 \Rightarrow \text{sign}(w^*_{\text{hinge}} \cdot x_i) = 1$$

$$y_i = -1 \Rightarrow y_i \cdot w^*_{\text{hinge}} \cdot x_i \geq 1 \Rightarrow w^*_{\text{hinge}} \cdot x_i \leq -1 \Rightarrow \text{sign}(w^*_{\text{hinge}} \cdot x_i) = -1$$

$(\text{sign}(w^*_{\text{hinge}} \cdot x_i) = y_i \quad \forall i)$  הוא המושך  $w^*_{\text{hinge}}$







proof

$$\leq \frac{1}{T} \sum_{t=1}^T \frac{1}{\eta_t} (\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2) + \frac{1}{T} \sum_{t=1}^T \frac{1}{2\eta_t} (\|x_t - y_{t+1}\|^2) \quad (3)$$

: p.s.d. n. d. s.p.w

$$\sum_{t=1}^T (\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2) =$$

$$= \|x_1 - x^*\|^2 - \underbrace{\|x_2 - x^*\|^2 + \|x_2 - x^*\|^2}_{\leq 0} - \dots - \underbrace{\|x_T - x^*\|^2 + \|x_T - x^*\|^2}_{\leq 0} - \|x_{T+1} - x^*\|^2$$

p.s.d. n. d. s.p.w

$$\|x_1 - x^*\|^2 - \underbrace{\|x_{T+1} - x^*\|^2}_{\geq 0} \leq \|x_1 - x^*\|^2 = \|x^*\|^2$$

$$\|x_t - x_{t+1}\|^2 \leq \|\nabla f(x_t) \cdot \eta\|^2 \leq \eta^2 \cdot \|\nabla f(x_t)\|^2 \quad \text{p.s.d. n. d. s.p.w} \quad (4)$$

$$\|x^*\| < B \quad , \quad \|\nabla f(x_t)\| < G \quad \text{p.s.d. n. d. s.p.w}$$

$$\eta_0 = \frac{\varepsilon}{G^2} \quad , \quad T = \frac{\varepsilon^2}{B \cdot G^2} \quad : \text{p.s.d. n. d. s.p.w}$$

$$f(\bar{x}) - f(x^*) \leq \frac{B \cdot T}{2\eta} + \frac{\eta \cdot G^2}{2} = \varepsilon \quad \text{p.s.d. n. d. s.p.w}$$



# משפט חשיבות

ה'  $f: \mathbb{R}^n \rightarrow \mathbb{R}$   $\beta$ -smooth  $\Rightarrow$   $f$  חסומה מלמעלה על כל קטע סגור.

(4)

נניח  $\eta \in \left(\frac{\beta}{2}, 1\right)$ .

$$\begin{aligned} f(x_{t+1}) &\leq f(x_t) + \nabla f(x_t)^\top (x_{t+1} - x_t) + \frac{\beta}{2} \|x_{t+1} - x_t\|^2 = \\ &\stackrel{\text{משפט חשיבות}}{=} f(x_t) + \langle \nabla f(x_t), (x_t - \eta \nabla f(x_t)) - x_t \rangle + \frac{\beta}{2} \|\eta \nabla f(x_t)\|^2 = \\ &= f(x_t) + \langle \nabla f(x_t), -\eta \nabla f(x_t) \rangle + \frac{\beta \eta^2}{2} \|\nabla f(x_t)\|^2 = \\ &= f(x_t) - \eta \|\nabla f(x_t)\|^2 + \frac{\beta \eta^2}{2} \|\nabla f(x_t)\|^2 = \\ &= f(x_t) + \left(1 - \eta - \frac{\beta \eta^2}{2}\right) \|\nabla f(x_t)\|^2 \Rightarrow \end{aligned}$$

$$\Rightarrow \|\nabla f(x_t)\|^2 \leq \frac{f(x_t) - f(x_{t+1})}{1 - \eta - \frac{\beta \eta^2}{2}}$$

$$\Rightarrow \sum_{t=1}^m \|\nabla f(x_t)\|^2 \leq \sum_{t=1}^m \frac{f(x_t) - f(x_{t+1})}{1 - \eta - \frac{\beta \eta^2}{2}} =$$

$$\stackrel{\text{סדר גורמים}}{=} \frac{1}{1 - \eta - \frac{\beta \eta^2}{2}} \left( \underbrace{f(x_1)}_{\text{קבוע}} - \underbrace{f(x_{m+1})}_{\geq 0} \right)$$

כאשר  $m \rightarrow \infty$  ו-  $f(x_{m+1}) \geq 0$  ו-  $f$  חסומה מלמעלה על כל קטע סגור.

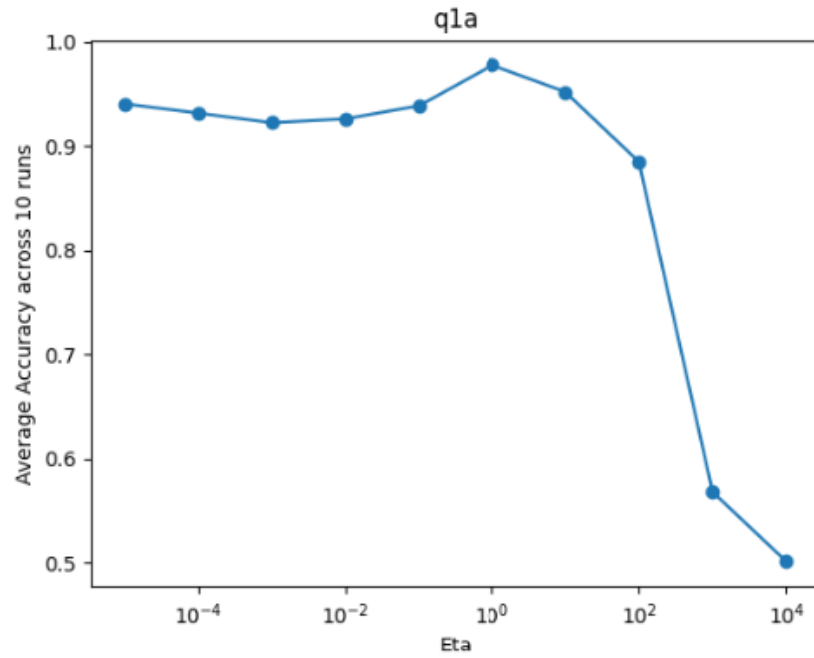
אם  $f$  חסומה מלמעלה על כל קטע סגור ו-  $f$  חסומה מלמטה על כל קטע סגור, אז  $\sum_{t=1}^{\infty} \|\nabla f(x_t)\|^2 < \infty$ .

~~אם  $f$  חסומה מלמעלה על כל קטע סגור ו-  $f$  חסומה מלמטה על כל קטע סגור, אז  $\sum_{t=1}^{\infty} \|\nabla f(x_t)\|^2 < \infty$ .~~

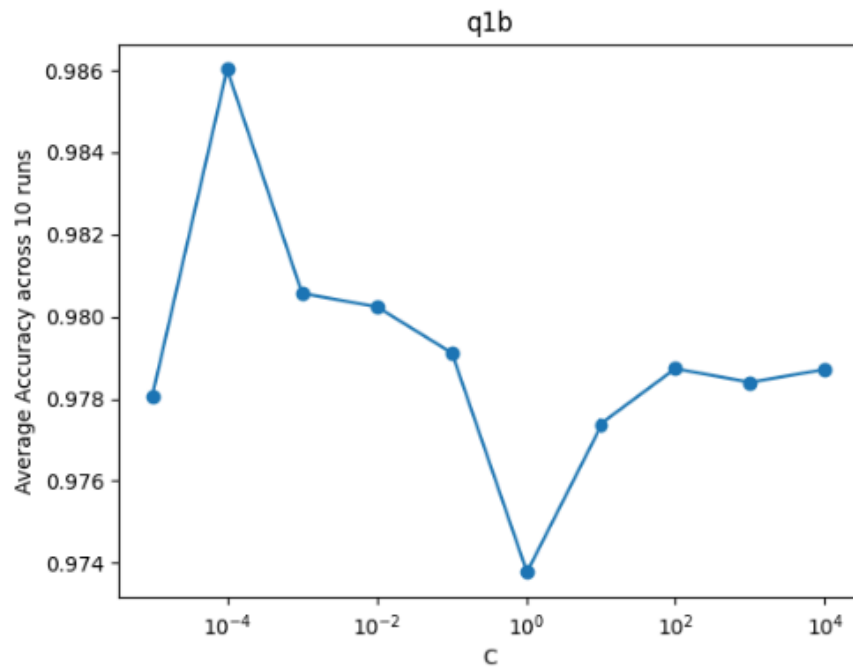
## חלק מעשי

### שאלה 1

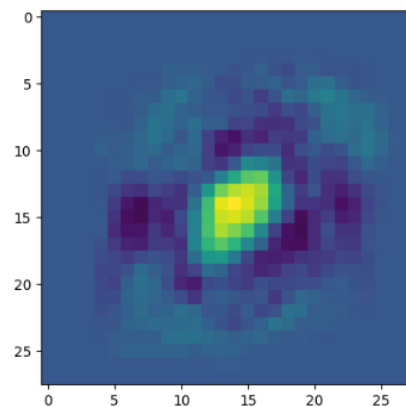
א. מקבלים כי ה  $\eta$  הכי טוב הוא  $10^0 = 1$



ב. מקבלים כי ה  $C$  הכי טוב הוא  $10^{-4} = 0.0001$



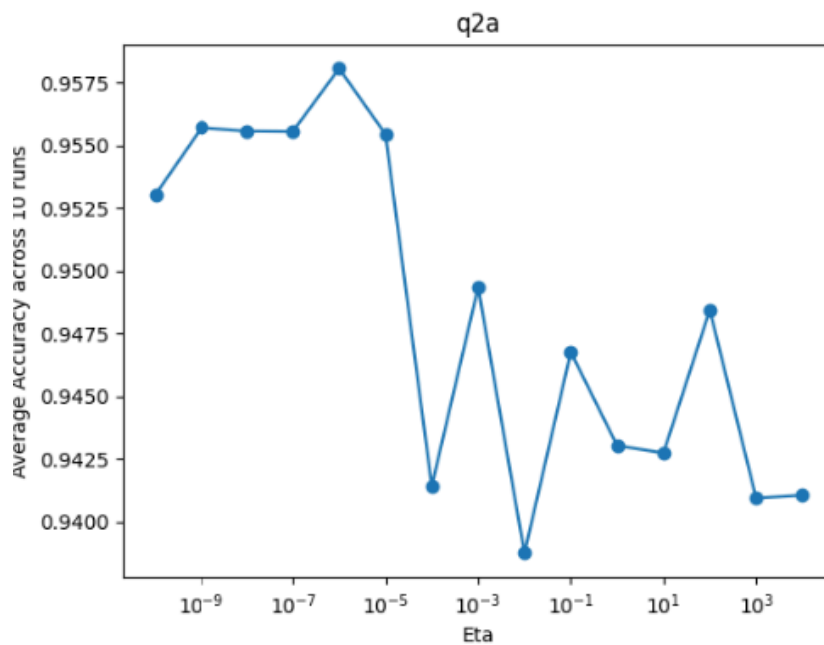
ג. התמונה מנרמלת את וקטור  $w$ , האזורים הכהים יותר מציינים סיכוי גבוה יותר לקבל 0 ואילו הבהירים יותר (במרכז) מציינים סיכוי גבוה יותר לקבל 8, ככל שהצבע בהיר המשקל גדול יותר. לכן ההבדל העיקר ביניהם הוא באמצע.



ד. מקבלים כי הדיוק הטוב ביותר הוא 0.9912998976458547.

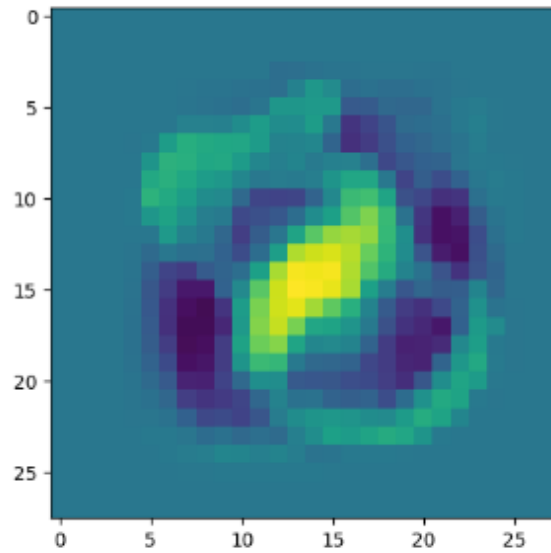
## שאלה 2

א. מקבלים כי ה  $\eta$  הכי טוב הוא  $10^{-6}$





ב. מקבלים כי הדיוק הטוב ביותר הוא 0.9616171954964176.



ג. מתחילים עם  $w=0$  לכן בהתחלה הנורמה היא 0.  $\eta$  כיוון שהיא מנורמלת לפי  $t$  היא קטנה עם הזמן ולכן קצב הגידול של  $w$  קטן עם הזמן. לכן בהתחלה אנחנו רחוקים מהפ האופטימלי אבל עם התקדמות באיטרציות של ה SGD ובגלל הרנדומליות נצפה להתקרב למינימום.

