

### מבוא ללמידה חישובית – מטלה 3

1) נחשב את ההסיאן של  $g(x) = f(Ax + b)$

נשתמש בכלל השרשרת:

$$\nabla g(x) = \nabla f(Ax + b)A = A^T \nabla f(Ax + b)$$

ושוב לי כלל השרשרת:

$$H_g(x) = A^T H_f(Ax + b)A$$

נתון ש  $f$  קמורה לכן  $H_f(Ax + b)$  היא PSD וסימטרית לכן קיימת  $Q$  הפיכה כך ש:

$$H_f(Ax + b) = Q^T D Q$$

$D$  אלכסונית לכן מתקיים

$$x^T A^T H_f(Ax + b) A x = x^T A^T Q^T D Q A x = \sum_{i=1}^n (Q A x)_i^2 D_{i,i}$$

$D_{i,i} \geq 0$  לכל  $i$  כי ההסיאן הוא PSD לכן

$$\forall x \in \mathbb{R}^n \sum_{i=1}^n (Q A x)_i^2 D_{i,i} \geq 0 \Rightarrow g(x) \text{ is convex}$$

(2)

יהיה  $x, y \in \mathbb{R}^d$  ו  $\lambda \in [0,1]$  נסמן  $t = \operatorname{argmax}_{i \in [m]} f_i(\lambda x + (1 - \lambda)y)$

(a)

$$\begin{aligned} \max_i f_i(\lambda x + (1 - \lambda)y) &= f_t(\lambda x + (1 - \lambda)y) \leq \lambda f_t(x) + (1 - \lambda)f_t(y) \leq \lambda \max_i f_i(x) \\ &+ (1 - \lambda) \max_i f_i(y) \Rightarrow \max_i f_i(x) \text{ is convex} \Rightarrow g(x) \text{ is convex} \end{aligned}$$

(b)

שנסמן  $t = \operatorname{argmax}_{i \in [m]} f_i(x)$  נתון ש  $f_i$  דיפרנצבילית לכן רציפות לכל  $i$  לכן קיימת סביבה קמורה (כדור עם רדיוס קטן מספיק)  $U \subset \mathbb{R}^d$  כך ש:

$$\forall x \in U \quad g(x) = f_t(x)$$

לכן מתקיים (מקמירות):

$$\forall w, u \in U \quad g(u) = f_t(u) \geq f_t(w) + \langle u - w, \nabla f_t(u) \rangle$$

לכן  $\nabla f_t$  הוא סאב-גדיאנט של  $g$

(3)

(a) נראה ש  $\ell$  כפונקציה של המשקולות של קמורה:

יהיה  $w_1, \dots, w_L$  לכן קיים  $\hat{y}$  כך  $w_{\hat{y}} x$  מקסימלי

$$\ell(w_1, \dots, w_L, x, y) = \max_{\hat{y} \in [L]} w_{\hat{y}} x - w_y x + \Delta_{zo}(\hat{y}, y)$$

נשים לב ש  $\Delta_{zo}(\hat{y}, y)$  הוא קבוע ולכן קמור (כי פונקציה קבועה היא קמורה)

$\max_{\hat{y} \in [L]} w_{\hat{y}} x$  קמורה משאלה 2,  $w_y x$  קמורה (פונקציה לינארית היא קמורה) ולכן  $\ell(w_1, \dots, w_L, x, y)$  קמורה כפונקציה של  $w_1, \dots, w_L$  (b)

$$\begin{aligned} \max_{\hat{y} \in [L]} w_{\hat{y}} x - w_y x &\geq 0 \Rightarrow \\ \ell(w_1, \dots, w_L, x, y) &= \max_{\hat{y} \in [L]} w_{\hat{y}} x - w_y x + \Delta_{zo}(\hat{y}, y) \geq \Delta_{zo}(\hat{y}, y) = \Delta_{zo}(f(x; w_1, \dots, w_L), y) \end{aligned} \quad (c)$$

$$\begin{aligned} \min_{w_1, \dots, w_L} \sum_i \ell(w_1, \dots, w_L, x, y) &= \sum_i \ell(w_1^{opt}, \dots, w_L^{opt}, x_i, y_i) \Rightarrow 0 \\ &= \sum_i \ell(w_1^*, \dots, w_L^*, x_i, y_i) \geq \sum_i \ell(w_1^{opt}, \dots, w_L^{opt}, x_i, y_i) \geq 0 \end{aligned}$$

לכן  $\sum_i \ell(w_1^{opt}, \dots, w_L^{opt}, x_i, y_i) = 0$  ובסעיף b ראינו ש

$$\ell(w_1, \dots, w_L, x, y) \geq \Delta_{zo}(f(x; w_1, \dots, w_L), y) \in \{0, 1\} \Rightarrow \Delta_{zo}(f(x; w_1^*, \dots, w_L^*), y) = 0 \quad 4$$

הערה: כדי נסמן את מספר הצעדים ב $\mathbb{T}$  כדי ששיחלוף ומספר הצעדים לא יסמנו באותו סימן בנוסף נניח כי  $w \in \mathbb{R}^d$

$$f(w) = w^T A w \Rightarrow \nabla f(w) = (A + A^T)w = 2Aw$$

לכן

$$w_{t+1} = w_t - 2\eta A w_t = (1 - 2\eta A)w_t \Rightarrow w_{\mathbb{T}} = (1 - 2\eta A)^{\mathbb{T}-1} w_1$$

נרצה לבחור  $\eta$  כך ש

$$f(w_{\mathbb{T}}) = f((1 - 2\eta A)^{\mathbb{T}-1} w_1) = ((1 - 2\eta A)^{\mathbb{T}-1} w_1)^T A ((1 - 2\eta A)^{\mathbb{T}-1} w_1) \leq \epsilon$$

נסמן  $\lambda_{max}$  את הערך העצמי המקסימלי של  $A$  ו $\lambda_{min}$  את הערך העצמי המינימלי של  $A$  בנוסף בגלל ש PCD מתקיים:

$$0 > \lambda_{min} > \lambda_{max}, \forall v \in \mathbb{R}^d \lambda_{min} \|v\| \leq v^T A v \leq \lambda_{max} \|v\|$$

$$\begin{aligned} ((1 - 2\eta A)^{\mathbb{T}-1} w_1)^T A ((1 - 2\eta A)^{\mathbb{T}-1} w_1) &\leq \lambda_{max} \|(1 - 2\eta A)^{\mathbb{T}-1} w_1\| \\ &= \lambda_{max} ((1 - 2\eta A)^{\mathbb{T}-1} w_1)^T ((1 - 2\eta A)^{\mathbb{T}-1} w_1) \leq \lambda_{max} (1 - 2\eta \lambda_{min})^{2\mathbb{T}-2} \|w_1\| \end{aligned}$$

$$\eta = \frac{1 - \frac{\epsilon}{\lambda_{max} \|w_1\|}}{2\lambda_{min}} \quad \text{נבחר}$$

ונקבל

$$f(w_{\mathbb{T}}) \leq \epsilon$$

לכן

$$(1 - 2\eta\lambda_{\min})^{2\mathbb{T}-2} \leq \frac{\epsilon}{\lambda_{\max}\|w_1\|} \Rightarrow 2\mathbb{T} - 2 \leq \frac{\log\left(\frac{\epsilon}{\lambda_{\max}\|w_1\|}\right)}{\log((1 - 2\eta\lambda_{\min}))} = \frac{\log\left(\frac{\lambda_{\max}\|w_1\|}{\epsilon}\right)}{-\log((1 - 2\eta\lambda_{\min}))}$$

כאשר  $1 - 2\eta\lambda_{\min}$  חיובי קטן מאחד עבור אפסילון קטן מספיק לכן

$$\mathbb{T} \leq \frac{\log\left(\frac{\lambda_{\max}\|w_1\|}{\epsilon}\right)}{-2\log((1 - 2\eta\lambda_{\min}))} + 2 \Rightarrow \mathbb{T} = O\left(\log\left(\frac{1}{\epsilon}\right)\right)$$

(5)

(1)

$$\text{יהי } 0 < \eta < \frac{2}{\beta} \text{ (} (*) \text{) נשים לב } 1 - \frac{\beta}{2}\eta > 0$$

$\ell$  היא בטא-חלקה אז עבור  $x_t$  ו  $x_{t+1} = x_t - \eta \nabla \ell(x_t)$  מתקיים:

$$\begin{aligned} \ell(x_{t+1}) &\leq \ell(x_t) + \nabla \ell(x_t)^T (x_{t+1} - x_t) + \frac{\beta}{2} \|x_t - x_{t+1}\|^2 \\ &= \ell(x_t) + \nabla \ell(x_t)^T (-\eta \nabla \ell(x_t)) + \frac{\beta}{2} \|\eta \nabla \ell(x_t)\|^2 \\ &= \ell(x_t) - \eta \|\nabla \ell(x_t)\|^2 + \frac{\beta}{2} \eta^2 \|\nabla \ell(x_t)\|^2 \\ &= \ell(x_t) - \|\nabla \ell(x_t)\|^2 \eta \left(1 - \frac{\beta}{2}\eta\right) \\ \Rightarrow \|\nabla \ell(x_t)\|^2 \eta \left(1 - \frac{\beta}{2}\eta\right) &\leq \ell(x_t) - \ell(x_{t+1}) \stackrel{(*)}{\Rightarrow} \|\nabla \ell(x_t)\|^2 \leq \frac{(\ell(x_t) - \ell(x_{t+1})))}{\eta \left(1 - \frac{\beta}{2}\eta\right)} \end{aligned}$$

נסמן  $q = \frac{1}{\eta(1 - \frac{\beta}{2}\eta)}$  אז מתקיים לכל  $n \in \mathbb{N}$ :

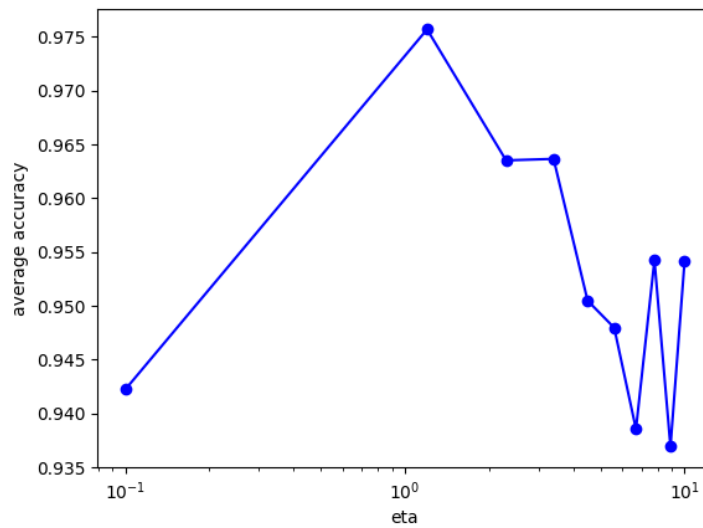
$$\sum_{t=0}^n \|\nabla \ell(x_t)\|^2 \leq q \sum_{i=0}^n (\ell(x_t) - \ell(x_{t+1})) = q(\ell(x_0) - \ell(x_n))$$

קיבלנו שסדרת הסכומים החלקיים של הטור האינסופי  $\sum_{t=0}^n \|\nabla \ell(x_t)\|^2$  חסומה כלומר הטור מתכנס ומתקיים

$$\lim_{t \rightarrow \infty} \|\nabla \ell(x_t)\|^2 = 0$$

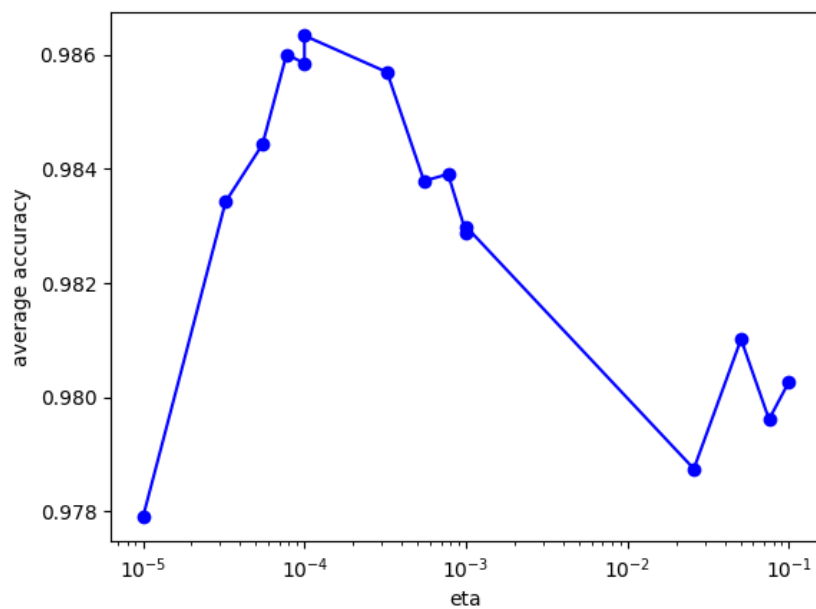
### חלק מעשי

(a) התחום המצומצם שמצאתי הוא  $[10^{-1}, 10]$  קיבלנו את הגרף



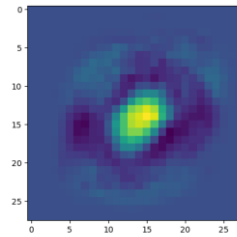
$\eta_0$  האופטימלי שמצאנו הוא 0.85 בקירוב

(b) בהינתן  $\eta_0 = 0.85$  התחום המצומצם שבו  $C$  אופטימלי הוא  $[10^{-5}, 10^{-1}]$  והגרף



וה  $C$  האופטימי הוא 0.00011

(c)

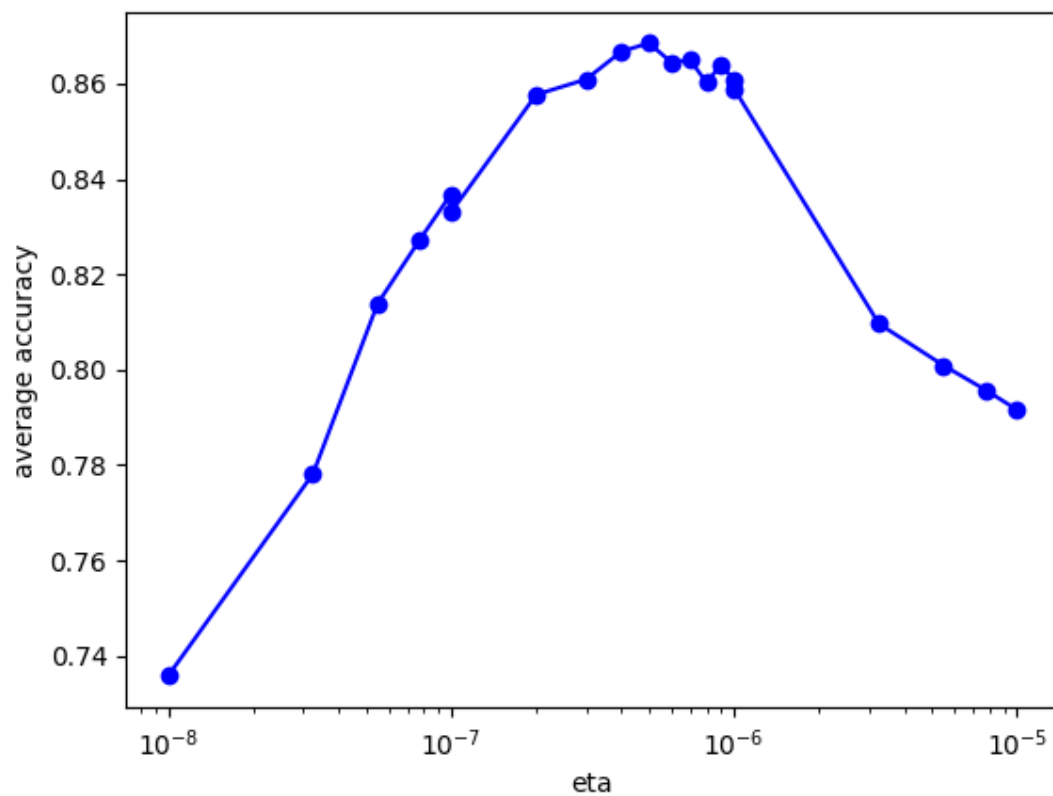


קיבלנו שרוב ה"משקל" הוא במרכז התמונה

(d) לאחר הרצה על *test set* קיבלנו דיוק של 0.993858

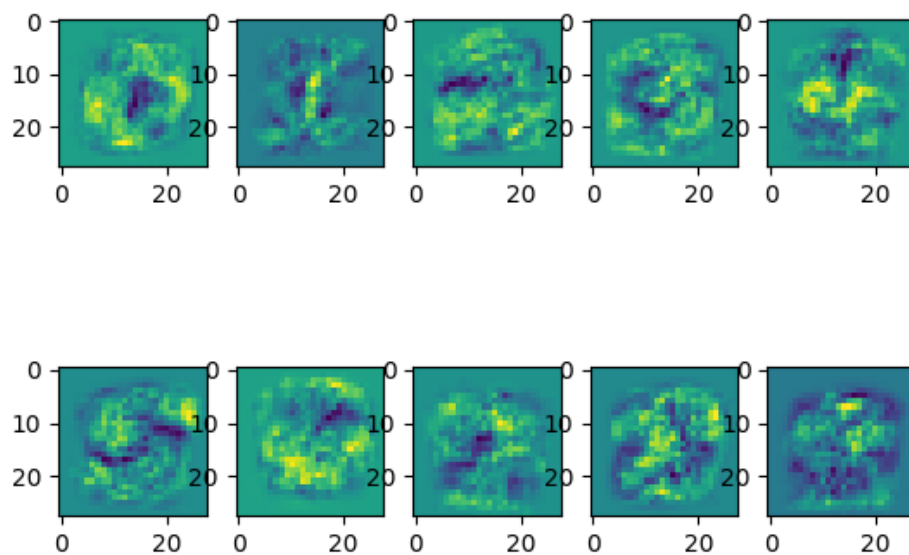
2.

(a) לאחר בדיקה כל כמה תחומים התחום הרלוונטי שבו מקבלים הטא אופטימלי הוא  $[10^{-8}, 10^{-5}]$  והגרף



והטא האופטימלית היא בערך  $5 * 10^{-7}$

(b)



בתמונות נראה שכל מסווג מקבל "משקל" גבוהה כך שבערך רואים במטושטש את המספר המתאים.

(c) לאחר הרצה על *test set* קיבלנו דיוק של 0.8785