

תרגיל בית 3 – מבוא ללמידה חישובית

מגיש: נתן בלור

Theory Questions

שאלה 1. Max of Convex Functions

נתונות m פונקציות קמורות, $f_1(x), \dots, f_m(x)$, כאשר $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$. נגדיר פונק' חדשה – $g(x) = \max_i f_i(x)$.
סעיף א. נדרש להוכיח כי $g(x)$ קמורה. g מוגדרת מעל \mathbb{R}^d שהינה קמורה. נדרש להוכיח כי לכל $w_1, w_2 \in \mathbb{R}^d$ ולכל $\lambda \in [0,1]$, מתקיים –

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

ראשית, נתון כי $f_1(x), \dots, f_m(x)$ הינן קמורות, ולכן לכל $i \in \{1, \dots, m\}$ מתקיים אי-השוויון לכל $w_1, w_2 \in \mathbb{R}^d$ ולכל $\lambda \in [0,1]$ –

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

ולפיכך ניתן לקבל את אי-השוויון הבא –

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

נזכיר כי $g(x)$ קמורה. יהי $w_1, w_2 \in \mathbb{R}^d$ ויהי $\lambda \in [0,1]$ מתקיים –

$$g(\lambda w_1 + (1 - \lambda)w_2) \stackrel{\text{def}}{=} \max_i f_i(\lambda w_1 + (1 - \lambda)w_2) \stackrel{\text{by (1)}}{\leq} \max_i \lambda f_i(w_1) + (1 - \lambda)f_i(w_2) \leq$$

$$\stackrel{(2)}{\leq} \max_i \lambda f_i(w_1) + \max_i (1 - \lambda)f_i(w_2) = \lambda \cdot \underbrace{\max_i f_i(w_1)}_{=g(w_1)} + (1 - \lambda) \cdot \underbrace{\max_i f_i(w_2)}_{=g(w_2)} = \lambda g(w_1) + (1 - \lambda)g(w_2)$$

נכונות מעבר (2) נובע מתכונות של פונק' המקסימום (\max). הראנו בעצם **שלכל** $w_1, w_2 \in \mathbb{R}^d$ **ולכל** $\lambda \in [0,1]$ מתקיים –
 $g(\lambda w_1 + (1 - \lambda)w_2) \leq \lambda g(w_1) + (1 - \lambda)g(w_2)$

נדרש, על מנת להוכיח קמירות של $g(x)$.

סעיף ב. נדרש להוכיח כי $\nabla f_i(x)$ עבור $f_i(x) = \max\{f_1(x), \dots, f_m(x)\}$ הינו סאב-גרדיינט (sub gradient) של g בנקודה x .
 נניח כי f_i דיפרנציאבילית לכל $i \in \{1, \dots, m\}$. ראינו בהרצאה את טענה 5.3.3 לפיה, עבור פונק' דיפרנציאבילית וקמורה f בנקודה כלשהי w , מתקיים $\partial f(w) = \{\nabla f(w)\}$ ובמילים אחרות, ה- sub gradient היחיד שלה בנקודה w הינו הגרדיאנט בנקודה. יהי $x \in \mathbb{R}^d$. נסמן $f_i(x) = \max\{f_1(x), \dots, f_m(x)\}$. נדרש להוכיח כי $z = \nabla f_i(x)$ הינו sub gradient של g בנקודה x , ולכן נרצה להוכיח שלכל $u \in \mathbb{R}^d$ מתקיים –

$$g(u) \geq g(x) + \langle u - x, \nabla f_i(x) \rangle$$

ובאופן שקול, נרצה להוכיח שלכל $u \in \mathbb{R}^d$ מתקיים –

$$g(u) - g(x) \geq \langle u - x, \nabla f_i(x) \rangle$$

ואכן, כיוון ש- $f_i(x)$ קמורה ודיפרנציאבילית בכל נקודה (ובפרט בנקודה $x \in \mathbb{R}^d$), נקבל לפי משפט 5.3.3 מההרצאה, שלכל $u \in \mathbb{R}^d$ מתקיים –

$$f_i(u) \geq f_i(x) + \langle u - x, \nabla f_i(x) \rangle$$

ובאופן שקול, לכל $u \in \mathbb{R}^d$ מתקיים –

$$(1) \quad f_i(u) - f_i(x) \geq \langle u - x, \nabla f_i(x) \rangle$$

כמו כן, נזכור כי – $f_i(x) = \max_j f_j(x)$, וכן לכל $u \in \mathbb{R}^d$ מתקיים $g(u) = \max_j f_j(u) \geq f_i(u)$. מכאן לכל $u \in \mathbb{R}^d$ מתקיים –

$$g(u) - g(x) = \underbrace{\max_j f_j(u)}_{\geq f_i(u)} - \underbrace{\max_j f_j(x)}_{=f_i(x)} \geq f_i(u) - f_i(x) \stackrel{\text{by (1)}}{\geq} \langle u - x, \nabla f_i(x) \rangle$$

ולכן, לפי ההגדרה מתקיים ש- $z = \nabla f_i(x)$ הינו sub gradient של g בנקודה x , כאשר $f_i(x) = \max_j f_j(x)$, וזה כאמור לכל $x \in \mathbb{R}^d$.
 ■

שאלה 2. ℓ^2 penalty.

נתונה הבעיה הבאה –

$$\min_{w, b, \xi} \frac{1}{2} w^T w + \frac{C}{2} \sum_{i=1}^m \xi_i^2$$

$$s. t. \quad y_i(w^T x_i + b) \geq 1 - \xi_i \quad \forall i \in \{1, \dots, m\}$$

סעיף א. נוכיח כי אילוצים מהצורה $\xi_i \geq 0$ לא ישנו את הבעיה. ובמילים אחרות, נוכיח כי כל פיתרון אופטימלי לבעיה בהכרח יקיים $\xi_i \geq 0$ לכל i . **פיתרון.** יהי w^*, b^*, ξ^* פיתרון אופטימלי כלשהו לבעיה. כלומר מתקיים –

$$y_i((w^*)^T x_i + b^*) \geq 1 - \xi_i^* \quad \forall i \in \{1, \dots, m\}$$

וכן לכל w, b, ξ , מתקיים –

$$\frac{1}{2} (w^*)^T (w^*) + \frac{C}{2} \sum_{i=1}^m (\xi_i^*)^2 \leq \frac{1}{2} w^T w + \frac{C}{2} \sum_{i=1}^m \xi_i^2$$

שכן הנחנו כי w^*, b^*, ξ^* פיתרון אופטימלי. נוכיח כי מתקיים – $\xi_i^* \geq 0$ לכל i . נניח בשלילה כי קיים i עבורו $\xi_k^* < 0$. עבור k זה מתקיים –

$$(1) \quad y_k((w^*)^T x_k + b^*) \geq 1 - \xi_k^* \geq 1$$

נגדיר את הפיתרון הבא – $w' = w^*, b' = b^*$ וכן נגדיר ξ' באופן הבא –

$$\xi'_j = \begin{cases} \xi_j^* & \text{if } j \neq k \\ 0 & \text{if } j = k \end{cases}$$

נראה כי הפיתרון w', b', ξ' מקיים את האילוצים של הבעיה, ואכן לכל $j \neq k$, האילוץ מתקיים שכן $\xi'_j = \xi_j^*$, $w' = w^*, b' = b^*$ ולכן –

$$y_j((w')^T x_j + b') = y_j((w^*)^T x_j + b^*) \geq 1 - \xi_j^* = 1 - \xi'_j$$

ועבור $j = k$, מתקיים –

$$y_k((w')^T x_k + b') = y_k((w^*)^T x_k + b^*) \geq 1 = 1 - \underbrace{\xi_k^*}_{=0} = 1 - \xi'_k$$

ואכן מתקיימים כל האילוצים הנדרשים. כעת, מתקיים – $\xi_k^* < 0 = \xi'_k$ ולכן $(\xi'_k)^2 = (\xi_k^*)^2 > 0$. לכן נבחין מכך שההבדל היחיד בין הפתרונות הינו ξ_k^* , שמתקיים מכך ש – $w' = w^*, b' = b^*$ וכן –

$$\begin{aligned} \underbrace{\frac{1}{2} (w^*)^T (w^*) + \frac{C}{2} \sum_{i=1}^m (\xi_i^*)^2}_{\text{result for } w^*, b^*, \xi^*} &= \frac{1}{2} (w^*)^T (w^*) + \frac{C}{2} \sum_{i=1, i \neq k}^m (\xi_i^*)^2 + \frac{C}{2} (\xi_k^*)^2 > \frac{1}{2} (w')^T (w') + \frac{C}{2} \sum_{i=1, i \neq k}^m (\xi'_i)^2 + \frac{C}{2} (\xi'_k)^2 \\ &= \underbrace{\frac{1}{2} (w')^T (w') + \frac{C}{2} \sum_{i=1}^m (\xi'_i)^2}_{\text{result for } w', b', \xi'} \end{aligned}$$

וזו סתירה כאמור, שכן קיבלנו כי הפיתרון עבור w^*, b^*, ξ^* אינו פיתרון אופטימלי. על כן, נסיק את הנדרש – מתקיים – $\xi_i^* \geq 0$ לכל i , ולכן אילוצים מהצורה $\xi_i \geq 0$ לא ישנו את הבעיה.

סעיף ב. הלגרנג'יאן של הבעיה.

פיתרון. נבחין כי ניתן לכתוב את האילוצים באופן הבא –

$$1 - \xi_i - y_i(w^T x_i + b) \leq 0 \quad \forall i \in \{1, \dots, m\}$$

ולכן הלגרנג'יאן במקרה זה הינו –

$$\mathcal{L}(w, b, \xi, \alpha) = \frac{1}{2} w^T w + \frac{C}{2} \sum_{i=1}^m \xi_i^2 + \sum_{i=1}^m \alpha_i (1 - \xi_i - y_i(w^T x_i + b))$$

סעיף ג. נביא את הלגרנג'יאן למינימום ביחס ל- w, b, ξ , ע"י גזירה ביחס למשתנים הללו והשוואת הנגזרות החלקיות ל-אפס.

פיתרון. מתקיים –

$$\frac{\partial \mathcal{L}(w, b, \xi, \alpha)}{\partial w_k} = w_k - \sum_{i=1}^m \alpha_i y_i x_i^{(k)}$$

כאשר $x_j^{(k)}$ מייצג את האינדקס ה- k של הוקטור x_i . ע"י השיוויון $\frac{\partial \mathcal{L}(w, b, \xi, \alpha)}{\partial w_k} = 0$, נקבל – $w_k = \sum_{i=1}^m \alpha_i y_i x_i^{(k)}$, ונסמן שיוויון זה ב- $(*)$. כמו כן, מתקיים –

$$\frac{\partial \mathcal{L}(w, b, \xi, \alpha)}{\partial b} = - \sum_{i=1}^m \alpha_i y_i$$

ע"י השיוויון $\frac{\partial \mathcal{L}(w, b, \xi, \alpha)}{\partial b} = 0$, נקבל – $\sum_{i=1}^m \alpha_i y_i = 0$, ונסמן שיוויון זה ב- $(**)$. כמו כן, מתקיים –

$$\frac{\partial \mathcal{L}(w, b, \xi, \alpha)}{\partial \xi_k} = C \xi_k - \alpha_k$$

ע"י השיוויון $\frac{\partial \mathcal{L}(w, b, \xi, \alpha)}{\partial \xi_k} = 0$, נקבל – $C \xi_k = \alpha_k$, לכל k , ונסמן שיוויונות אלה ב- $(***)$. נבחין כי מתקיים –

$$\sum_{i=1}^m \alpha_i y_i x_i = \alpha_1 y_1 x_1 + \dots + \alpha_m y_m x_m = \begin{pmatrix} \alpha_1 y_1 x_1^{(1)} \\ \vdots \\ \alpha_1 y_1 x_1^{(n)} \end{pmatrix} + \dots + \begin{pmatrix} \alpha_m y_m x_m^{(1)} \\ \vdots \\ \alpha_m y_m x_m^{(n)} \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^m \alpha_i y_i x_i^{(1)} \\ \vdots \\ \sum_{i=1}^m \alpha_i y_i x_i^{(n)} \end{pmatrix} \stackrel{(*)}{=} \begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix} = w$$

נסתכל כעת על הביטוי $\sum_{i=1}^m \alpha_i (1 - \xi_i - y_i (w^T x_i + b))$ כגורם בלגרנג'יאן ונציב בו את האילוצים שקיבלנו כך ש – $\nabla \partial \mathcal{L}(w, b, \xi, \alpha) = 0$. ביחס למשתנים w, b, ξ . נקבל –

$$\begin{aligned} \sum_{i=1}^m \alpha_i (1 - \xi_i - y_i (w^T x_i + b)) &= \sum_{i=1}^m \alpha_i - \sum_{i=1}^m \alpha_i \xi_i - \sum_{i=1}^m \alpha_i y_i (w^T x_i + b) = \sum_{i=1}^m \alpha_i - \sum_{i=1}^m \alpha_i \xi_i - \sum_{i=1}^m \alpha_i y_i w^T x_i - \sum_{i=1}^m \alpha_i y_i b = \\ &= \sum_{i=1}^m \alpha_i - \sum_{i=1}^m \alpha_i \xi_i - w^T \underbrace{\sum_{i=1}^m \alpha_i y_i x_i}_{=w} - b \underbrace{\sum_{i=1}^m \alpha_i y_i}_{=0} \stackrel{\substack{\text{by } (***) \\ \text{and } (**)}}{=} \sum_{i=1}^m \alpha_i - \frac{1}{C} \sum_{i=1}^m \alpha_i \alpha_i - \left(\underbrace{\sum_{i=1}^m \alpha_i y_i x_i}_{=w} \right)^T \left(\underbrace{\sum_{i=1}^m \alpha_i y_i x_i}_{=w} \right) \\ &= \sum_{i=1}^m \alpha_i - \frac{1}{C} \sum_{i=1}^m \alpha_i \alpha_i - \sum_{i=1}^m \sum_{j=1}^m \alpha_i y_i x_i \alpha_j y_j x_j \end{aligned}$$

נציב בלגרנג'יאן את האילוצים –

$$\begin{aligned} \mathcal{L}(w, b, \xi, \alpha) &= \frac{1}{2} w^T w + \frac{C}{2} \sum_{i=1}^m \xi_i^2 + \sum_{i=1}^m \alpha_i (1 - \xi_i - y_i (w^T x_i + b)) = \\ &= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y_i x_i \alpha_j y_j x_j + \frac{C}{2} \sum_{i=1}^m \alpha_i^2 + \sum_{i=1}^m \alpha_i - \frac{1}{C} \sum_{i=1}^m \alpha_i \alpha_i - \sum_{i=1}^m \sum_{j=1}^m \alpha_i y_i x_i \alpha_j y_j x_j = \\ &= \boxed{-\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y_i x_i \alpha_j y_j x_j + \sum_{i=1}^m \alpha_i - \frac{1}{2} \cdot \frac{1}{C} \sum_{i=1}^m \alpha_i^2} \end{aligned}$$

סעיף ד. התוכנית הדואלית.

ראינו בהרצאה כי התוכנית הפרימאלית, שקולה לבעיית ה-min-max הבאה –

$$\min_{w,b,\xi} \max_{\alpha} \mathcal{L}(w, b, \xi, \alpha)$$

וכן ראינו כי מתקיים האי-שיוויון –

$$\min_{w,b,\xi} \max_{\alpha} \mathcal{L}(w, b, \xi, \alpha) \geq \max_{\alpha} \min_{w,b,\xi} \mathcal{L}(w, b, \xi, \alpha)$$

ואם נגדיר $g(\alpha) = \min_{w,b,\xi} \mathcal{L}(w, b, \xi, \alpha)$, ופונק' המטרה של התוכנית הדואלית תהיה מקסימיזציה של החסם התחתון לפונק' המטרה המקורית, ולכן –

$$\max_{\alpha} g(\alpha) = \max_{\alpha} \min_{w,b,\xi} \mathcal{L}(w, b, \xi, \alpha)$$

וכיוון שהפונק' $\mathcal{L}(w, b, \xi, \alpha)$ הינה קמורה (בסכום של פונק' קמורות כל אחת), נקבל כי המינימום שלה מתקבל בנקודה בה הגרדיאנט שלה לפי w, b, ξ מקיים $\nabla \mathcal{L}(w, b, \xi, \alpha) = 0$. ראינו כי במקרה זה, המינימום מקיים –

$$\mathcal{L}(w, b, \xi, \alpha) = -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y_i x_i \alpha_j y_j x_j + \sum_{i=1}^m \alpha_i - \frac{1}{2} \cdot \frac{1}{C} \sum_{i=1}^m \alpha_i^2$$

ולכן פונק' המטרה של התוכנית הדואלית הינה –

$$\max_{\alpha} g(\alpha) = \max_{\alpha} -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y_i x_i \alpha_j y_j x_j + \sum_{i=1}^m \alpha_i - \frac{1}{2} \cdot \frac{1}{C} \sum_{i=1}^m \alpha_i^2$$

נבחין כי אספנו בדרך אילוח נוסף לפיו –

$$\sum_{i=1}^m \alpha_i y_i = 0$$

אילוח זה נובע מכך שבלגרנג'יאן קיבלנו –

$$\begin{aligned} \mathcal{L}(w, b, \xi, \alpha) &= \frac{1}{2} w^T w + \frac{C}{2} \sum_{i=1}^m \xi_i^2 + \sum_{i=1}^m \alpha_i (1 - \xi_i - y_i (w^T x_i + b)) = \\ &= \frac{1}{2} w^T w + \frac{C}{2} \sum_{i=1}^m \xi_i^2 + \sum_{i=1}^m \alpha_i - \sum_{i=1}^m \alpha_i \xi_i - \sum_{i=1}^m \alpha_i y_i w^T x_i - b \sum_{i=1}^m \alpha_i y_i \end{aligned}$$

ונבחין כי אם $\sum_{i=1}^m \alpha_i y_i \neq 0$, אזי נקבל תמיד כי $\min_{w,b,\xi} \mathcal{L}(w, b, \xi, \alpha) = -\infty$.

כמו כן, כיוון ש- α_i הינם כופלים של אילוצי אי-שיוויון, נקבל את אילוצי אי-שליליות לגביהם – $\alpha_i \geq 0$.

בסה"כ, התוכנית הדואלית הינה –

$$\boxed{\begin{aligned} \max_{\alpha} \quad & -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y_i x_i \alpha_j y_j x_j + \sum_{i=1}^m \alpha_i - \frac{1}{2} \cdot \frac{1}{C} \sum_{i=1}^m \alpha_i^2 \\ \text{s.t.} \quad & \sum_{i=1}^m \alpha_i y_i = 0 \\ & \alpha_i \geq 0 \quad \forall i \in \{1, \dots, m\} \end{aligned}}$$

שאלה 3. The Multi-Class Hinge-Loss.

נגדיר את פונק' ההפסד (multi class hinge loss function) באופן הבא –

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

סעיף א. נדרש להוכיח כי ℓ פונקציית קמורה של המשתנים w_1, \dots, w_L . ידוע כי פונק' המקסימום של פונק' קמורות כל אחת הינה גם כן פונק' קמורה, ולכן מספיק להראות שלכל $\hat{y} \in [L]$ הפונק' $w_{\hat{y}} \cdot x - w_y \cdot x + \Delta_{zo}(\hat{y}, y)$ קמורה של המשתנים w_1, \dots, w_L . כמו כן, ידוע כי סכום של פונק' קמורות הינה פונק' קמורה ולכן מספיק להראות כי כל גורם בפונק' $w_{\hat{y}} \cdot x - w_y \cdot x + \Delta_{zo}(\hat{y}, y)$ הינו פונק' קמורה של המשתנים w_1, \dots, w_L . ואכן, יהיו x, y, \hat{y} קבועים כלשהם ונוכיח קמירות. הגורם $\Delta_{zo}(\hat{y}, y)$ הינה פונק' קבועה של המשתנים w_1, \dots, w_L ולכן זו פונק' קמורה. הגורמים של $w_{\hat{y}} \cdot x$ ו- $-w_y \cdot x$ הינם פונק' ליניאריות של המשתנים w_1, \dots, w_L , ולכן מדובר בפונק' קמורות גם כן. ובסה"כ, כל הגורמים מהווים פונק' קמורות ביחס למשתנים w_1, \dots, w_L , ולכל $\hat{y} \in [L]$ הפונק' $w_{\hat{y}} \cdot x - w_y \cdot x + \Delta_{zo}(\hat{y}, y)$ הינה פונק' קמורה של המשתנים w_1, \dots, w_L . מכאן, נסיק כי הפונק' $\ell(w_1, \dots, w_L, x, y)$ הינה קמורה במקסימום של פונק' קמורות כל אחת, כנדרש.

סעיף ב. נדרש להוכיח שמתקיים $\ell(w_1, \dots, w_L, x, y) \geq \Delta_{zo}(f(x; w_1, \dots, w_L), y)$ לכל w, x, y . יהי w, x, y כלשהם. נסמן כך $i = \arg\max_{y \in [L]} w_y \cdot x =: i$. מתקיים לכן $i \in [L]$ וכן $w_i \cdot x \geq w_y \cdot x$ לכל $y \in [L]$. כיוון שלפי ההגדרה $i \in [L]$, לפי הגדרת מקסימום מתקיים –

$$(1) \quad \max_{\hat{y} \in [L]} w_{\hat{y}} \cdot x - w_y \cdot x + \Delta_{zo}(\hat{y}, y) \geq w_i \cdot x - w_y \cdot x + \Delta_{zo}(i, y)$$

שכן הגורם $w_i \cdot x - w_y \cdot x + \Delta_{zo}(i, y)$ הינו גורם בפונק' \max -ה"ל. לכן,

$$\ell(w_1, \dots, w_L, x, y) = \max_{\hat{y} \in [L]} w_{\hat{y}} \cdot x - w_y \cdot x + \Delta_{zo}(\hat{y}, y) \underset{(1)}{\geq} \underbrace{w_i \cdot x - w_y \cdot x + \Delta_{zo}(i, y)}_{\substack{\geq 0 \text{ as} \\ \forall y \in [L]. \\ w_i \cdot x \geq w_y \cdot x}} \geq \Delta_{zo}(i, y) = \Delta_{zo}(f(x; w_1, \dots, w_L), y)$$

כאשר השוויון האחרון נובע מכך שסימנו $i := \arg\max_{y \in [L]} w_y \cdot x = f(x; w_1, \dots, w_L)$

ובסה"כ הראנו בסעיף זה כי לכל w, x, y , מתקיים –

$$\ell(w_1, \dots, w_L, x, y) \geq \Delta_{zo}(f(x; w_1, \dots, w_L), y)$$

כנדרש.

סעיף ג. נתון כי קיימים w_1^*, \dots, w_L^* , שמשיגים טעות אמפירית (טעות על גבי סט האימון) של אפס. במילים אחרות, לכל $i \in [L]$, מתקיים –

$$\Delta_{zo}(f(x_i; w_1^*, \dots, w_L^*), y_i) = 0$$

נוכיח כי גם בעבור $w_1^{opt}, \dots, w_L^{opt}$, לכל $i \in [L]$, מתקיים $\Delta_{zo}(f(x_i; w_1^{opt}, \dots, w_L^{opt}), y_i) = 0$. פיתרון. נוכיח כי מתקיים שוויון ע"י כך שנוכיח כי מתקיימים שני אי-שוויונים – $0 \leq \Delta_{zo}(f(x_i; w_1^{opt}, \dots, w_L^{opt}), y_i) \leq 0$.

בכיוון הראשון, בו נוכיח כי $0 \leq \Delta_{zo}(f(x_i; w_1^{opt}, \dots, w_L^{opt}), y_i)$ הינו מידי מהגדרת הפונק' $\Delta_{zo}(\cdot, \cdot)$, שכן לכל $\hat{y} \in [L]$, מתקיים $\Delta_{zo}(y, \hat{y}) \in \{0, 1\}$, ובפרט אי-שלילי. ולכן $0 \leq \Delta_{zo}(f(x_i; w_1^{opt}, \dots, w_L^{opt}), y_i)$ כנדרש בכיוון זה.

בכיוון השני, נדרש להוכיח כי מתקיים $\Delta_{zo}(f(x_i; w_1^{opt}, \dots, w_L^{opt}), y_i) \leq 0$ לפי ההנחה, לכל $i \in [L]$, מתקיים –

$$\Delta_{zo}(f(x_i; w_1^*, \dots, w_L^*), y_i) = 0 \quad \underset{\substack{\text{def of} \\ \Delta_{zo}}}{\implies} f(x_i; w_1^*, \dots, w_L^*) = y_i \quad \underset{\substack{\text{def of} \\ f}}{\implies} \arg\max_{y \in [L]} w_y \cdot x = y_i$$

$$\implies \max_{y \in [L]} w_y \cdot x \leq w_{y_i} \cdot x \implies \max_{\substack{y \in [L] \\ y \neq y_i}} w_y \cdot x < w_{y_i} \cdot x$$

כאשר המעבר האחרון נובע מכך ש- $f(x_i; w_1^*, \dots, w_L^*) = \arg\max_{y \in [L]} w_y \cdot x$ ו- f פונק' ולכן חד-ערכית, ולכן לא ייתכנו שני ארגומנטים שהם maximizers .

לכל $i \in [L]$, נגדיר –

$$\varepsilon_i = \max_{\substack{y \in [L] \\ y \neq y_i}} w_y \cdot x - w_{y_i} \cdot x$$

ולפי מה שהוכחנו ב-(*), מתקיים לכל $i \in [L]$, $\varepsilon_i < 0$. נגדיר גם כן $\varepsilon := \min_i |\varepsilon_i| > 0$. כעת, נגדיר את הוקטורים w_1', \dots, w_L' , באופן הבא –

$$w_i' = \frac{1}{\varepsilon} \cdot w_i^*$$

כמו כן, נסתכל על הביטוי $\max_{y \in [L]} w'_y \cdot x_i - w'_{y_i} \cdot x_i + \Delta_{zo}(y, y_i) -$ ונבחין שעבור $y = y_i$, נקבל $0 = \underbrace{w'_y \cdot x_i - w'_{y_i} \cdot x_i}_{=0} + \underbrace{\Delta_{zo}(y_i, y_i)}_{=0}$ ולכן ניתן

להציג זאת באופן הבא **השקול** –

$$\max_{y \in [L]} w'_y \cdot x_i - w'_{y_i} \cdot x_i + \Delta_{zo}(y, y_i) = \max\{0, \max_{\substack{y \in [L] \\ y \neq y_i}} w'_y \cdot x_i - w'_{y_i} \cdot x_i + \Delta_{zo}(y, y_i)\} = \max\{0, \max_{\substack{y \in [L] \\ y \neq y_i}} w'_y \cdot x_i - w'_{y_i} \cdot x_i + 1\} \quad (*)$$

ובנוסף, נבחין כי במקרה זה, המקסימום נלקח על גבי איברי הקבוצה $[L]$ **השונים** מ- y_i , ולכן מתקיים $\Delta_{zo}(y, y_i) = 1$

מתקיים –

$$\begin{aligned} \sum_i \Delta_{zo}(f(x_i; w_1^{opt}, \dots, w_L^{opt}), y_i) &\stackrel{(b)}{\leq} \sum_i \ell(w_1^{opt}, \dots, w_L^{opt}, x_i, y_i) \stackrel{\substack{\leq \\ w^{opt} \\ optimizer}}{\leq} \sum_i \ell(w'_1, \dots, w'_L, x_i, y_i) = \\ &= \sum_i \max_{y \in [L]} w'_y \cdot x_i - w'_{y_i} \cdot x_i + \Delta_{zo}(y, y_i) \stackrel{(*)}{=} \sum_i \max\{0, \max_{\substack{y \in [L] \\ y \neq y_i}} w'_y \cdot x_i - w'_{y_i} \cdot x_i + 1\} = \\ &= \sum_i \max\{0, 1 + \max_{\substack{y \in [L] \\ y \neq y_i}} \frac{1}{\varepsilon} \cdot w_y^* \cdot x_i - \frac{1}{\varepsilon} \cdot w_{y_i}^* \cdot x_i\} = \sum_i \max\{0, 1 + \underbrace{\frac{1}{\varepsilon} \cdot \max_{\substack{y \in [L] \\ y \neq y_i}} w_y^* \cdot x_i - w_{y_i}^* \cdot x_i}_{=\varepsilon_i}\} = \\ &= \sum_i \max\{0, 1 + \frac{1}{\varepsilon} \cdot \varepsilon_i\} \stackrel{(**)}{=} \sum_i \max\{0, 1 - \frac{1}{\varepsilon} \cdot |\varepsilon_i|\} \stackrel{(***)}{\leq} \sum_i \max\{0, 0\} = \sum_i 0 = 0 \end{aligned}$$

כאשר השוויון המסומן ב-**(**)**, נובע מכך ש- $\varepsilon_i < 0$, ולכן $\varepsilon_i = -|\varepsilon_i|$. אי-השוויון המסומן ב-**(***)**, נובע מהטיעונים הבאים –

$$\varepsilon = \min_j |\varepsilon_j| \leq |\varepsilon_i| \quad \Rightarrow \quad \frac{1}{\varepsilon} \geq \frac{1}{|\varepsilon_i|} \quad \Rightarrow \quad -\frac{1}{\varepsilon} \leq -\frac{1}{|\varepsilon_i|} \quad \Rightarrow \quad -\frac{1}{\varepsilon} \cdot |\varepsilon_i| \leq -\frac{1}{|\varepsilon_i|} \cdot |\varepsilon_i| \quad \Rightarrow \quad 1 - \frac{1}{\varepsilon} \cdot |\varepsilon_i| \leq 1 - \underbrace{\frac{1}{|\varepsilon_i|} \cdot |\varepsilon_i|}_{=1} = 0$$

כלומר, מצאנו כי –

$$\sum_i \Delta_{zo}(f(x_i; w_1^{opt}, \dots, w_L^{opt}), y_i) \leq 0$$

ובנו כן, הוכחנו קודם לכן כי $0 \leq \Delta_{zo}(f(x_i; w_1^{opt}, \dots, w_L^{opt}), y_i)$, כלומר שמדובר בסכום של גורמים אי-שליליים החסום ע"י 0, ולכן בפרט ניתן להסיק את הרצוי –

$$\Delta_{zo}(f(x_i; w_1^{opt}, \dots, w_L^{opt}), y_i) = 0$$

■

לכל $i \in [L]$, כנדרש.

שאלה 4. Growth Function of Composition.

נתונות $F_1 \subseteq \mathcal{Y}_1^{\mathcal{X}}, F_2 \subseteq \mathcal{Y}_2^{\mathcal{Y}_1}$, ונגדיר $F = F_2 \circ F_1$, המוגדרת ע"י הרכבה של פונק' מ- F_1 ומ- F_2 , כך ש-
 $F = F_2 \circ F_1 = \{f_2 \circ f_1 \mid f_1 \in F_1, f_2 \in F_2\}$

נדרש להוכיח כי מתקיים – $\pi_F(m) \leq \pi_{F_1}(m) \cdot \pi_{F_2}(m)$

פיתרון. תהי $C = \{x_1, \dots, x_m\} \subseteq \mathcal{X}$ תת-קבוצה בגודל m . נגדיר בנוסף $f_1^C = \{f_1(x_1), \dots, f_1(x_m)\} = \{y_1, \dots, y_m \mid s.t. y_i = f_1(x_i)\} \subseteq \mathcal{Y}_1$
 תת-קבוצה של \mathcal{Y}_1 בגודל m , שהינה תלויה כאמור בבחירת C וכן תלויה ב- $f_1 \in F_1$.
 לפי ההגדרה, מתקיים –

$$\begin{aligned} |F_C| &= |\{[f(x_1), \dots, f(x_m)]: f \in F, s.t. C = \{x_1, \dots, x_m\}\}| = \\ &= |\{[f_2(f_1(x_1)), \dots, f_2(f_1(x_m))]: f_1 \in F_1, f_2 \in F_2, s.t. C = \{x_1, \dots, x_m\}\}| = \\ &= |\{[f_2(y_1), \dots, f_2(y_m)]: f_1 \in F_1, f_2 \in F_2, s.t. C = \{x_1, \dots, x_m\} \text{ and } y_i = f_1(x_i)\}| = \\ &\stackrel{\substack{= \\ \text{def} \\ \text{of} \\ F_{1C}}}{=} \left| \bigcup_{[y_1, \dots, y_m] \in F_{1C}} \{[f_2(y_1), \dots, f_2(y_m)]: f_2 \in F_2\} \right| \leq \\ &\stackrel{\substack{\leq \\ \text{Union} \\ \text{Bound}}}{=} \sum_{[y_1, \dots, y_m] \in F_{1C}} |\{[f_2(y_1), \dots, f_2(y_m)]: f_2 \in F_2\}| = \\ &= \sum_{[y_1, \dots, y_m] \in F_{1C}} |F_{2_{f_1^C}}| \leq \sum_{[y_1, \dots, y_m] \in F_{1C}} \max_{\substack{C' \subseteq \mathcal{Y}_1 \\ s.t. |C'|=m}} |F_{2_{C'}}| = \sum_{[y_1, \dots, y_m] \in F_{1C}} \pi_{F_2}(m) = \\ &= |F_{1C}| \cdot \pi_{F_2}(m) \leq \max_{\substack{C' \subseteq \mathcal{X} \\ s.t. |C'|=m}} |F_{1_{C'}}| \cdot \pi_{F_2}(m) = \pi_{F_1}(m) \cdot \pi_{F_2}(m) \end{aligned}$$

נבחין כי הקבוצה $C = \{x_1, \dots, x_m\} \subseteq \mathcal{X}$ הינה קבוצה כללית מגודל m , ולכן לכל קבוצה $C \subseteq \mathcal{X}$ מגודל m , מתקיים – $|F_C| \leq \pi_{F_1}(m) \cdot \pi_{F_2}(m)$
 ולפיכך ניתן להסיק כי הא"ש מתקיים גם עבור תת-קבוצה C מגודל m שעבורה $|F_C|$ מקסימלי, כלומר –

$$\max_{\substack{C \subseteq \mathcal{X} \\ s.t. |C|=m}} |F_C| \leq \pi_{F_1}(m) \cdot \pi_{F_2}(m)$$

$\underbrace{\hspace{10em}}_{=\pi_F(m)}$

■

ואכן קיבלנו את הנדרש – $\pi_F(m) \leq \pi_{F_1}(m) \cdot \pi_{F_2}(m)$ כרצוי.

שאלה 5. Gradient Descent on Smooth Function.

נתונה $l: \mathbb{R}^n \rightarrow \mathbb{R}$ פונק' שהיא β -חלקה (β -Smooth) ואי-שלילית. נסתכל על אלגוריתם GD עבור הפונק' l עם גודל צעד קבוע $\eta > 0$ –

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \nabla l(\mathbf{x}_t)$$

עם נק' התחלה של \mathbf{x}_0 . נוכיח שאם $\eta < \frac{2}{\beta}$, אזי מתקיים –

$$\lim_{t \rightarrow \infty} \|\nabla l(\mathbf{x}_t)\| = 0$$

פיתרון. נתון כי l פונק' שהיא β -חלקה, ולכן לכל $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ מתקיים –

$$l(\mathbf{y}) \leq l(\mathbf{x}) + \nabla l(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|^2$$

ובפרט עבור הנקודות $\mathbf{x}_t, \mathbf{x}_{t+1}$ יתקיים – $\mathbf{x}_{t+1} - \mathbf{x}_t = -\eta \nabla l(\mathbf{x}_t)$ וע"י שימוש בכך ש- l פונק' שהיא β -חלקה, נקבל –

$$\begin{aligned} l(\mathbf{x}_{t+1}) &\leq l(\mathbf{x}_t) + \nabla l(\mathbf{x}_t)^T (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{\beta}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 = l(\mathbf{x}_t) - \eta \underbrace{\nabla l(\mathbf{x}_t)^T \nabla l(\mathbf{x}_t)}_{=\|\nabla l(\mathbf{x}_t)\|^2} + \eta^2 \frac{\beta}{2} \|\nabla l(\mathbf{x}_t)\|^2 = \\ &= l(\mathbf{x}_t) - \eta \|\nabla l(\mathbf{x}_t)\|^2 + \eta^2 \frac{\beta}{2} \|\nabla l(\mathbf{x}_t)\|^2 = l(\mathbf{x}_t) + (\eta^2 \frac{\beta}{2} - \eta) \|\nabla l(\mathbf{x}_t)\|^2 \end{aligned}$$

ומכאן נקבל –

$$(1) \quad -(\eta^2 \frac{\beta}{2} - \eta) \|\nabla l(\mathbf{x}_t)\|^2 \leq l(\mathbf{x}_t) - l(\mathbf{x}_{t+1})$$

ונבחין כי

$$\eta^2 \frac{\beta}{2} - \eta \underset{\eta < \frac{2}{\beta}}{\leq} \eta \cdot \underbrace{\frac{2}{\beta} \cdot \frac{\beta}{2}}_{=1} - \eta = \eta - \eta = 0 \quad \Rightarrow \quad -(\eta^2 \frac{\beta}{2} - \eta) < 0$$

ולכן ע"י חלוקה של אי-השוויון (1) בגורם $-(\eta^2 \frac{\beta}{2} - \eta)$, נקבל כי מתקיים –

$$\|\nabla l(\mathbf{x}_t)\|^2 \leq \underbrace{(-(\eta^2 \frac{\beta}{2} - \eta))}_{const.} \cdot (l(\mathbf{x}_t) - l(\mathbf{x}_{t+1}))$$

ואי-שוויון זה מתקיים כאמור לכל $t \geq 0$ טבעי. נוכיח כי $\sum_{t=0}^{\infty} \|\nabla l(\mathbf{x}_t)\|^2 < \infty$ כלומר מתכנסת, ולכן מכאן לפי ההתכנסות הטור הנ"ל, נסיק כי האיבר הכללי $\|\nabla l(\mathbf{x}_t)\|^2 \xrightarrow{t \rightarrow \infty} 0$, ולפיכך בהכרח מתקיים $\|\nabla l(\mathbf{x}_t)\| \xrightarrow{t \rightarrow \infty} 0$, כפי שנדרש. על מנת להוכיח התכנסות של הטור, נוכיח התכנסות של סדרת הסכומים החלקיים. מתקיים –

$$\sum_{t=0}^N \|\nabla l(\mathbf{x}_t)\|^2 \leq \underbrace{\left(-(\eta^2 \frac{\beta}{2} - \eta)\right)}_{const.} \cdot \sum_{t=0}^N l(\mathbf{x}_t) - l(\mathbf{x}_{t+1}) = \underbrace{\left(-(\eta^2 \frac{\beta}{2} - \eta)\right)}_{const.} \cdot (l(\mathbf{x}_0) - l(\mathbf{x}_{N+1})) \underset{l(\cdot) \geq 0}{\leq} \underbrace{\left(-(\eta^2 \frac{\beta}{2} - \eta)\right) \cdot l(\mathbf{x}_0)}_{const. \forall N \in \mathbb{N}}$$

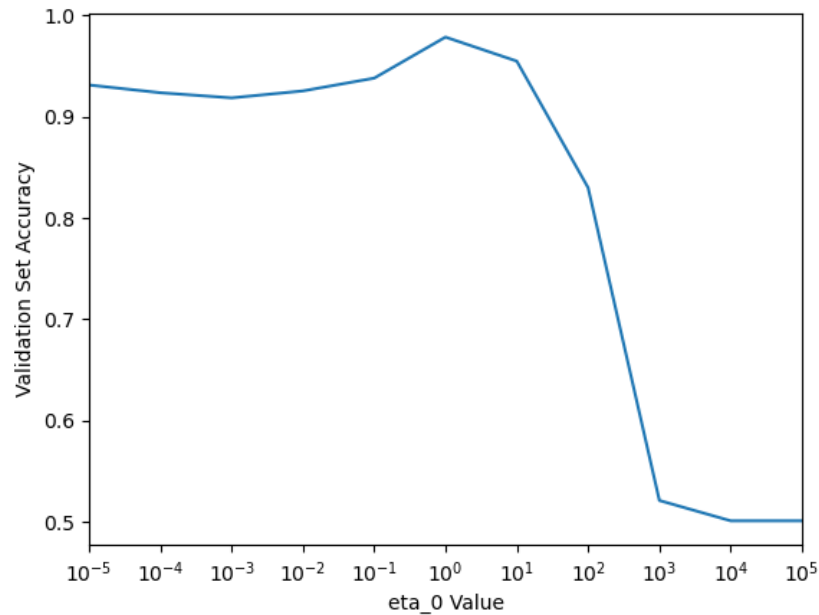
כלומר, מצאנו כי סדרת הסכומים החלקיים חסומה ע"י קבוע. נבחין בנוסף כי סדרת הסכומים החלקיים הינה סדרה מונוטונית עולה, שכן $\|\cdot\| \geq 0$, ולפיכך ניתן להסיק כי סדרת הסכומים מתכנסת שכן היא מונו' עולה וחסומה מלעיל.
לכן, נסיק כי –

$$\sum_{t=0}^{\infty} \|\nabla l(\mathbf{x}_t)\|^2 < \infty$$

ולכן האיבר הכללי שואף ל-אפס, כלומר $\|\nabla l(\mathbf{x}_t)\|^2 \xrightarrow{t \rightarrow \infty} 0$ ולפיכך נסיק כי $\|\nabla l(\mathbf{x}_t)\| \xrightarrow{t \rightarrow \infty} 0$, כנדרש. ■

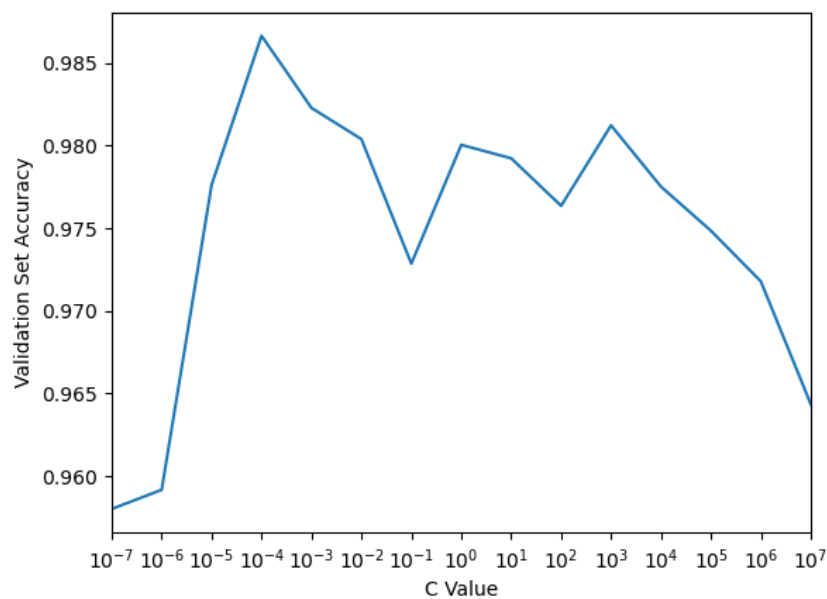
שאלה 1. SGD for Hinge Loss

סעיף א. נגדיר $T = 1000, C = 1$, ונבצע SGD על גבי $\eta_0 \in \{10^{-5}, 10^{-4}, \dots, 10^4, 10^5\}$, ונבצע cross-validation, כדי לאמוד את טיב ההיפוטזה שנלמדה, כדי לקבל את הגרף הבא של דיוק הממוצע על גבי סט הוואלידציה (accuracy on validation set), כפונק' של η_0 . התוצאות –



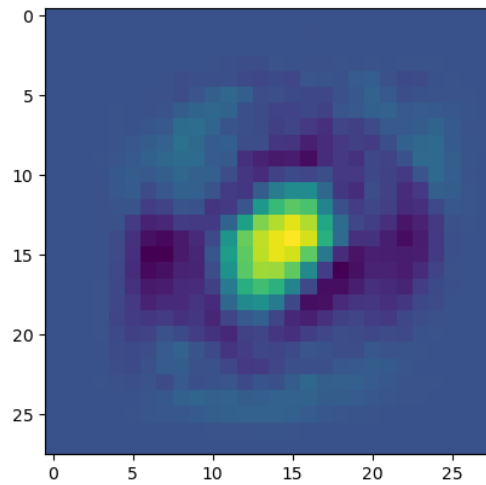
נבחין כי η_0 הטוב ביותר התקבל עבור $\eta_0 = 10^0 = 1$, עפ"י הגרף שהתקבל.

סעיף ב. בסעיף זה נשתמש ב- η_0 הטוב ביותר שהתקבל בסעיף א', $\eta_0 = 1$. נגדיר גם $T = 1000$, ונבצע SGD על גבי $C \in \{10^{-7}, 10^{-6}, \dots, 10^6, 10^7\}$, ונבצע cross-validation, כדי לאמוד את טיב ההיפוטזה שנלמדה, כדי לקבל את הגרף הבא של דיוק הממוצע על גבי סט הוואלידציה (accuracy on validation set), כפונק' של C . התוצאות –



נבחין כי C הטוב ביותר התקבל עבור $C = 10^{-4}$, עפ"י הגרף שהתקבל.

סעיף ג. בסעיף זה נשתמש ב- η_0 הטוב ביותר שהתקבל בסעיף א', $\eta_0 = 1$, וב- C הטוב ביותר שהתקבל בסעיף ב', $C = 10^{-4}$.
 נאמן את המסווג עבור $T = 20000$, עם הפרמטרים שמצאנו בסעיפים הקודמים. נציג את w , וקטור המשקלים, המתקבל כתמונה –



ניתן לראות כי ישנו ריכוז במרכז התמונה, כלומר הפיקסלים במרכז התמונה הם הפיקסלים המרכזיים יותר באבחנה בין המספר 0 למספר 8.

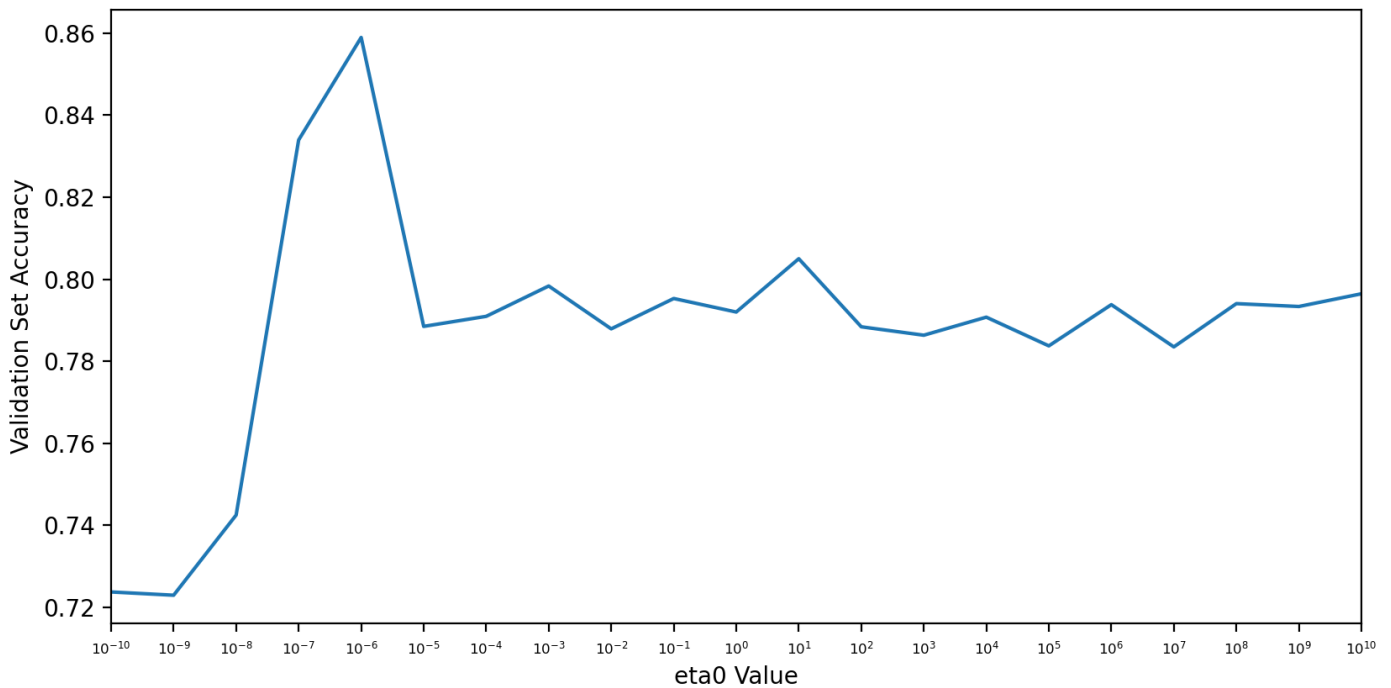
סעיף ד. נחשב את הדיוק של המסווג האופטימלי על גבי ה-*test set*.

The accuracy of the best classifier on the test set is 0.9923234390992836

הדיוק המתקבל הינו 0.9923234390992836.

שאלה 2. SGD for Multi Class Cross Entropy.

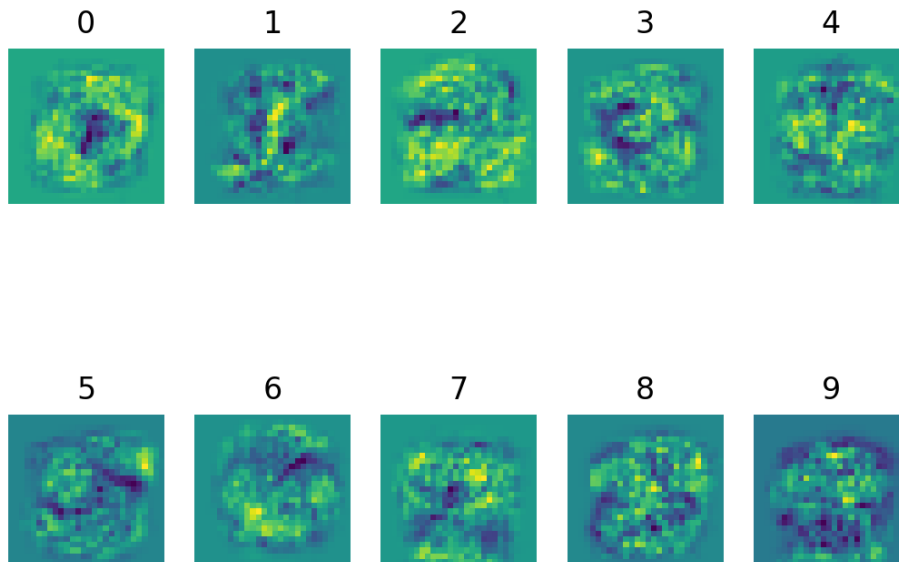
סעיף א. נגדיר $T = 1000$, ובבצע SGD על גבי $\eta_0 \in \{10^{-10}, 10^{-9}, \dots, 10^9, 10^{10}\}$, ובבצע $cross-validation$, כדי לאמוד את טיב ההיפותזה שנלמדה, כדי לקבל את הגרף הבא של דיוק הממוצע על גבי סט הוואלידציה ($accuracy\ on\ validation\ set$), כפונק' של η_0 . התוצאות –



נבחין כי η_0 הטוב ביותר התקבל עבור $\eta_0 = 10^{-6}$, עפ"י הגרף שהתקבל.

סעיף ב. בסעיף זה נשתמש ב- η_0 הטוב ביותר שהתקבל בסעיף א', $\eta_0 = 10^{-6}$. נאמן את המסווג עם $\eta_0 = 10^{-6}$ עבור $T = 20000$.

נציג את וקטורי המשקלים w_0, w_1, \dots, w_9 שהתקבלו כתמונות –



סעיף ג. נחשב את הדיוק של המסווג האופטימלי על גבי ה- $test\ set$.

The accuracy of the best classifier on the test set is 0.869

הדיוק המתקבל הינו 0.869 .