

# **Intro to Machine Learning: Assignment #3**

April 30, 2023

**Dor Bourshan**

315780122

# Theory Questions

## 1 Step-size Perceptron.

נאמר שהאלגוריתם Perceptron מבצע את העדכון הבא בעת בשגיאה  $(\hat{y}_t \neq y_t)$ :

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta_t y_t \mathbf{x}_t$$

נניח שהדטא מופרד עם מרווח  $\gamma > 0$  וכי מתקיים  $\|\mathbf{x}_t\|_2 = 1$  לכל  $t$ <sup>1</sup>. בנוסף נניח כי האלגוריתם מבצע  $M$  שגיאות וכולן מתרחשות בתחילתו.

נראה שעבור  $\eta_t = \frac{1}{\sqrt{t}}$  מספר השגיאות של האלגוריתם חסום ע"י  $\mathcal{O}\left(\frac{4}{\gamma^2} \ln\left(\frac{1}{\gamma}\right)\right)$  כאשר אנו מניחים כי  $\gamma < 1$  אחרת אין לחסם משמעות.

מכך שהדטא מופרד, קיים מסווג מושלם שנשמנו  $\mathbf{w}^*$ , ונוכל להניח בלי הגבלת הכלליות כי  $\|\mathbf{w}_t\|_2 = \frac{1}{2}$ .

לכל  $1 \leq t \leq M$  האלגוריתם מבצע טעות, ולכן:

$$\langle \mathbf{w}_{t+1}, \mathbf{w}^* \rangle = \langle \mathbf{w}_t + \eta_t y_t \mathbf{x}_t, \mathbf{w}^* \rangle = \langle \mathbf{w}_t, \mathbf{w}^* \rangle + \langle \eta_t y_t \mathbf{x}_t, \mathbf{w}^* \rangle = \langle \mathbf{w}_t, \mathbf{w}^* \rangle + \eta_t \langle y_t \mathbf{x}_t, \mathbf{w}^* \rangle \geq \langle \mathbf{w}_t, \mathbf{w}^* \rangle + \frac{1}{\sqrt{t}} \gamma$$

כלומר עבור  $t > M$  מתקיים כי

$$\langle \mathbf{w}_{t+1}, \mathbf{w}^* \rangle \geq \gamma \cdot \sum_{t=1}^M \frac{1}{\sqrt{t}} \geq \gamma \sqrt{M} \quad (1)$$

בנוסף, מתקיים כי לכל  $1 \leq t \leq M$

$$\|\mathbf{w}_{t+1}\|_2^2 = \|\mathbf{w}_t + \eta_t y_t \mathbf{x}_t\|_2^2 = \|\mathbf{w}_t\|_2^2 + 2 \langle \mathbf{w}_t, \eta_t y_t \mathbf{x}_t \rangle + \|\eta_t y_t \mathbf{x}_t\|_2^2$$

נשים לב כי לכל  $1 \leq t \leq M$  האלגוריתם מבצע טעות ולכן  $\langle \mathbf{w}_t, y_t \mathbf{x}_t \rangle < 0$  ולכן  $2 \langle \mathbf{w}_t, \eta_t y_t \mathbf{x}_t \rangle < 0$  כלומר אנו מקבלים כי

$$\|\mathbf{w}_{t+1}\|_2^2 = \|\mathbf{w}_t\|_2^2 + 2 \langle \mathbf{w}_t, \eta_t y_t \mathbf{x}_t \rangle + \|\eta_t y_t \mathbf{x}_t\|_2^2 \leq \|\mathbf{w}_t\|_2^2 + \frac{1}{t} \|\mathbf{x}_t\|_2^2 \stackrel{1}{=} \|\mathbf{w}_t\|_2^2 + \frac{1}{t}$$

כלומר עבור  $t' > M$  מתקיים כי

$$\|\mathbf{w}_t\|_2^2 \leq \sum_{p=1}^M \frac{1}{p} = H_M \leq 1 + \ln(M) \quad (2)$$

אם כן עבור  $t > M$  מתקיים כי:

$$\gamma \sqrt{M} \leq \langle \mathbf{w}_t, \mathbf{w}^* \rangle \stackrel{1}{\leq} \|\mathbf{w}_t\|_2 \|\mathbf{w}^*\|_2 \stackrel{2}{\leq} \sqrt{\frac{1}{4} + \frac{1}{4} \ln(M)}$$

Cauchy-Schwarz

אם כן

$$\gamma\sqrt{M} \leq \sqrt{\frac{1}{4} + \frac{1}{4} \ln(M)}$$

• אם  $M = 0$  ברור כי מתקיים כי  $M \leq \frac{4}{\gamma^2} \ln\left(\frac{1}{\gamma}\right)$ .

• אם  $M = 1$  מתקיים כי  $\gamma \leq \frac{1}{4}$  וכן  $1 < 16 \ln(4) \leq \frac{4}{\gamma^2} \ln\left(\frac{1}{\gamma}\right)$ .

• עבור  $M \geq 2$  :

עבור  $x \geq 2$  מתקיים כי

$$1 + \ln(x) \leq 4 \ln(x)$$

כיוון ש

$$1 < \ln(8) \leq \ln(x^3) = 4 \ln(x) - \ln(x)$$

ולכן

$$\frac{1}{4} + \frac{1}{4} \ln(x) \leq \ln(x)$$

אם כן מתקיים כי

$$\gamma\sqrt{M} \leq \sqrt{\frac{1}{4} + \frac{1}{4} \ln(M)} \leq \sqrt{\ln(M)}$$

ולכן

$$\gamma^2 M \leq \ln(M)$$

ומהרמז מתקיים כי

$$M \leq \frac{2}{\gamma^2} \ln\left(\frac{1}{\gamma^2}\right) = \frac{4}{\gamma^2} \ln\left(\frac{1}{\gamma}\right)$$

לכן סה"כ מתקיים כי  $\frac{4}{\gamma^2} \ln\left(\frac{1}{\gamma}\right)$  חוסם את  $M$  כנדרש.

□

## 2 Convex functions.

### 2.1 section (a)

תהי  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  פונקציה קמורה,  $A \in \mathbb{R}^{n \times n}$  ו  $\mathbf{b} \in \mathbb{R}^n$ . נראה כי  $g(\mathbf{x}) = f(A\mathbf{x} + \mathbf{b})$  היא פונקציה קמורה.

ניזכר כי מטריצה  $A$  מגדירה העתקה ליניארית, כלומר  $A$  היא פונקציה  $T_A: \mathbb{R}^n \rightarrow \mathbb{R}^n$ . אם כן נסמן את הקבוצה הבאה:

$$C = \{A\mathbf{x} + \mathbf{b} \mid \mathbf{x} \in \mathbb{R}^n\} \subseteq \mathbb{R}^n$$

$C$  קבוצה קמורה. יהי  $\lambda \in (0, 1)$ , ויהיו  $\mathbf{c}_1, \mathbf{c}_2 \in C$ , קיימים  $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$  כך ש  $\mathbf{c}_i = A\mathbf{x}_i + \mathbf{b}$  עבור  $i \in \{1, 2\}$ .

$$\lambda \mathbf{c}_1 + (1 - \lambda) \mathbf{c}_2 = \lambda (A\mathbf{x}_1 + \mathbf{b}) + (1 - \lambda) (A\mathbf{x}_2 + \mathbf{b}) = A\lambda \mathbf{x}_1 + A(1 - \lambda) \mathbf{x}_2 + \mathbf{b} = A(\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2) + \mathbf{b}$$

ברור כי  $\mathbf{y} = \lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2 \in \mathbb{R}^n$  ולכן  $A\mathbf{y} + \mathbf{b} \in C$  כנדרש.

אם כן מתקיים כי  $C$  קבוצה קמורה ולכל  $\mathbf{c}_1, \mathbf{c}_2 \in C$  ולכל  $\lambda \in (0, 1)$ , קיימים  $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$  כך ש  $\mathbf{c}_i = A\mathbf{x}_i + \mathbf{b}$  עבור  $i \in \{1, 2\}$  וכן  $A(\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2) + \mathbf{b} \in C$  אם כן

$$\begin{aligned} g(\lambda \mathbf{c}_1 + (1 - \lambda) \mathbf{c}_2) &= f(A(\lambda \mathbf{c}_1 + (1 - \lambda) \mathbf{c}_2) + \mathbf{b}) = f(\lambda (A\mathbf{c}_1 + \mathbf{b}) + (1 - \lambda) (A\mathbf{c}_2 + \mathbf{b})) \\ &\leq \lambda f(A\mathbf{c}_1 + \mathbf{b}) + (1 - \lambda) f(A\mathbf{c}_2 + \mathbf{b}) \end{aligned}$$

כלומר

$$g(\lambda \mathbf{c}_1 + (1 - \lambda) \mathbf{c}_2) \leq \lambda f(A\mathbf{c}_1 + \mathbf{b}) + (1 - \lambda) f(A\mathbf{c}_2 + \mathbf{b}) = \lambda g(\mathbf{c}_1) + (1 - \lambda) g(\mathbf{c}_2)$$

כנדרש מפונקציה קמורה.

□

## 2.2 section (b)

תהינה  $m$  פונקציות קמורות  $f_1(\mathbf{x}), \dots, f_m(\mathbf{x})$  כך ש  $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ . נגדיר פונקציה חדשה  $g : \mathbb{R}^d \rightarrow \mathbb{R}$

$$g(\mathbf{x}) = \max_i f_i(\mathbf{x})$$

נראה כי  $g$  פונקציה קמורה.

ראשית  $\mathbb{R}^d$  קבוצה קמורה (מהיותה שדה).

יהיו  $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d$  ויהי  $\lambda \in (0, 1)$ . מתקיים כי

$$\begin{aligned} g(\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2) &= \max_i f_i(\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2) \leq \max_i (\lambda f_i(\mathbf{x}_1) + (1 - \lambda) f_i(\mathbf{x}_2)) \\ &\leq \lambda \max_i f_i(\mathbf{x}_1) + (1 - \lambda) \max_i f_i(\mathbf{x}_2) \end{aligned}$$

כלומר

$$g(\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2) \leq \lambda \max_i f_i(\mathbf{x}_1) + (1 - \lambda) \max_i f_i(\mathbf{x}_2) = \lambda g(\mathbf{x}_1) + (1 - \lambda) g(\mathbf{x}_2)$$

כנדרש מפונקציה קמורה.

□

## 2.3 section (c)

תהי  $\ell_{\log} : \mathbb{R} \rightarrow \mathbb{R}$  פונקציית  $\log$  loss המוגדרת ע"י

$$\ell_{\log}(z) = \log_2(1 + e^{-z})$$

נוכיח תחילה כי  $\ell_{\log}$  היא פונקציה קמורה.

תחילה  $\mathbb{R}$  היא קבוצה קמורה (מהיותה שדה). נסמן לשם הנוחות

$$g(z) = 1 + e^{-z}$$

$$f(x) = \log_2(x)$$

אם כן

$$\ell_{\log} = h = f \circ g$$

ניזכר במשפט שהוכח בקורס חדו"א

**משפט:** אם  $f$  גזירה פעמיים בקטע, והנגזרת השנייה אי שלילית בכל הקטע, אז הפונקציה קמורה בקטע.

נגזור את  $h$  פעמיים ונראה כי הנגזרת אי שלילית ונסיק כי פונקציית  $\log$  loss היא פונקציה קמורה.

אם כן לפי כלל השרשרת:

$$h' = f'(g(x)) g'(x) = -\frac{e^{-x}}{(1+e^{-x}) \ln(2)}$$

ושוב

$$h'' = \frac{1}{\ln(2)} \left( \frac{e^{-x}(1+e^{-x}) - e^{-2x}}{(1+e^{-x})^2} \right) = \frac{1}{\ln(2)} \left( \frac{e^{-x}((1+e^{-x}) - e^{-x})}{(1+e^{-x})^2} \right) = \frac{1}{\ln(2)} \left( \frac{e^{-x}}{(1+e^{-x})^2} \right) > 0$$

לפיכך  $\ell_{log}$  קמורה כנדרש.

□

כעת נסיק כי  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  המוגדרת ע"י

$$f(\mathbf{w}) = \ell_{log}(y\mathbf{w} \cdot \mathbf{x})$$

היא פונקציה קמורה ביחס ל  $\mathbf{w}$ .

ראשית  $\mathbb{R}^d$  קבוצה קמורה (מהיותה שדה).

יהי  $\lambda \in (0, 1)$  ויהי  $\mathbf{w}_1, \mathbf{w}_2 \in \mathbb{R}^d$

$$\begin{aligned} f(\lambda \mathbf{w}_1 + (1-\lambda) \mathbf{w}_2) &= \ell_{log}(\langle \lambda \mathbf{w}_1 + (1-\lambda) \mathbf{w}_2, y\mathbf{x} \rangle) = \ell_{log}(\langle \lambda \mathbf{w}_1 + (1-\lambda) \mathbf{w}_2, y\mathbf{x} \rangle) \\ &= \ell_{log}(\lambda y \langle \mathbf{w}_1, \mathbf{x} \rangle + (1-\lambda) y \langle \mathbf{w}_2, \mathbf{x} \rangle) \end{aligned}$$

מכך שהראנו כי  $\ell_{log}$  פונקציה קמורה וכן  $\langle \mathbf{w}_i, \mathbf{x} \rangle \in \mathbb{R}$

$$\begin{aligned} f(\lambda \mathbf{w}_1 + (1-\lambda) \mathbf{w}_2) &= \ell_{log}(\lambda y \langle \mathbf{w}_1, \mathbf{x} \rangle + (1-\lambda) y \langle \mathbf{w}_2, \mathbf{x} \rangle) \leq \lambda \ell_{log}(y \langle \mathbf{w}_1, \mathbf{x} \rangle) + (1-\lambda) \ell_{log}(y \langle \mathbf{w}_2, \mathbf{x} \rangle) \\ &= \lambda f(\mathbf{w}_1) + (1-\lambda) f(\mathbf{w}_2) \end{aligned}$$

כנדרש מפונקציה קמורה.

□

### 3 Ranking.

נניח כי  $\mathcal{X} \subseteq \mathbb{R}^d$  וניתנת קבוצת אימון של  $n$  רשימות של  $k$  איברים, לכל רשימה ניתן וקטור label. במפורש

$$S = \left\{ \left( (\mathbf{x}_1^i, \mathbf{x}_2^i, \dots, \mathbf{x}_k^i), \mathbf{y}^i \right) \right\}_{i=1}^n$$

כך שלכל  $\mathbf{y}^i \in \mathbb{R}^k$ ,  $1 \leq i \leq n$ , מושם ערך לכל ערך ב  $(\bar{\mathbf{x}}^i) = (\mathbf{x}_1^i, \dots, \mathbf{x}_k^i)$ .

המטרה היא ללמוד פונקציית ranking  $h : \mathcal{X}^k \rightarrow \mathbb{R}^k$  שמבצעת rank נכון לרשימת האיברים שמגיעים מ  $S$ .

מגדירים את  $Kendall - Tau$  loss בין שני וקטורי rank  $\mathbf{y}', \mathbf{y} : \text{rank}$  ע"י :

$$\Delta(\mathbf{y}', \mathbf{y}) = \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{r=j+1}^k \mathbf{1} \{ \text{sgn}(y'_j - y'_r) \neq \text{sgn}(y_j - y_r) \}$$

נניח שאנו מעוניינים ללמוד פונקציית rank לינארית, כלומר פונקציה מהצורה

$$h_{\mathbf{w}}((\mathbf{x}_1, \dots, \mathbf{x}_k)) = (\langle \mathbf{w}, \mathbf{x}_1 \rangle, \dots, \langle \mathbf{w}, \mathbf{x}_k \rangle)$$

עבור איזושהו  $\mathbf{w} \in \mathbb{R}^d$ , והמטרה היא למנמם את Kendall-Tau loss על  $S$

$$\sum_{i=1}^n \Delta(h_{\mathbf{w}}(\bar{\mathbf{x}}^i), \mathbf{y}^i)$$

כיוון שפונקציה זו קשה לאופטימיזציה, אנו נשתמש בפונקציית hinge loss

$$\sum_{i=1}^n \ell(h_{\mathbf{w}}(\bar{\mathbf{x}}), \mathbf{y})$$

כאשר

$$\ell(h_{\mathbf{w}}(\bar{\mathbf{x}}), \mathbf{y}) = \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{r=j+1}^k \max\{0, 1 - \text{sgn}(y_j - y_r) \langle \mathbf{w}, \mathbf{x}_j - \mathbf{x}_r \rangle\}$$

#### 3.1 section (a)

תחילה נוכיח כי פונקציית hinge loss לעיל היא פונקציה קמורה ביחס ל  $\mathbf{w}$ .

אנחנו נראה כי  $1 - \text{sgn}(y_j - y_r) \langle \mathbf{w}, \mathbf{x}_j - \mathbf{x}_r \rangle$  היא פונקציה קמורה ביחס ל  $\mathbf{w}$  וע"י שימוש ב 2.2 נסיק כי

$\max\{0, 1 - \text{sgn}(y_j - y_r) \langle \mathbf{w}, \mathbf{x}_j - \mathbf{x}_r \rangle\}$  פונקציה קמורה (שכן פונקציית האפס היא פונקציה קמורה). כדי להשלים את התמונה אנו נצטרך להוכיח כי סכום של פונקציות קמורות היא פונקציה קמורה.

אם כן ראשית נוכיח כי סכום של פונקציות קמורות היא פונקציה קמורה באינדוקציה.

• **בסיס :** נניח כי  $f, g$  הן פונקציות קמורות על הקבוצה הקמורה  $C$  ונסמן  $h = f + g$ . יהי  $\lambda \in (0, 1)$  ויהיו  $\mathbf{c}_1, \mathbf{c}_2 \in C$

$$\begin{aligned} h(\lambda \mathbf{c}_1 + (1 - \lambda) \mathbf{c}_2) &= f(\lambda \mathbf{c}_1 + (1 - \lambda) \mathbf{c}_2) + g(\lambda \mathbf{c}_1 + (1 - \lambda) \mathbf{c}_2) \leq \\ &\lambda(f(\mathbf{c}_1) + g(\mathbf{c}_1)) + (1 - \lambda)(f(\mathbf{c}_2) + g(\mathbf{c}_2)) = \lambda h(\mathbf{c}_1) + (1 - \lambda) h(\mathbf{c}_2) \end{aligned}$$

כנדרש

• **צעד :** נניח כי הנ"ל נכון עבור  $n$  ונראה עבור  $n + 1$ . יהיו  $f_1, \dots, f_{n+1}$  פונקציות קמורות על הקבוצה הקמורה  $C$  : נגדיר  $h = \sum_{i=1}^{n+1} f_i$ . מהנחת האינדוקציה  $\sum_{i=1}^n f_i$  היא פונקציה קמורה, ובדומה להוכחה מקרה הבסיס  $\sum_{i=1}^n f_i + f_{n+1}$  כלומר  $h = \sum_{i=1}^{n+1} f_i$  קמורה כנדרש.

כעת נראה כי  $f(\mathbf{w}) = 1 - \text{sgn}(y_j - y_r) \langle \mathbf{w}, \mathbf{x}_j - \mathbf{x}_r \rangle$  היא פונקציה קמורה ביחס ל  $\mathbf{w}$ . נסמנה

ראשית  $\mathbf{w} \in \mathbb{R}^d$  ו  $\mathbb{R}^d$  קבוצה קמורה (מהיותה שדה). יהי  $\lambda \in (0, 1)$  ויהיו  $\mathbf{w}_1, \mathbf{w}_2 \in \mathbb{R}^d$

$$\begin{aligned} f(\lambda \mathbf{w}_1 + (1 - \lambda) \mathbf{w}_2) &= 1 - \text{sgn}(y_j - y_r) \langle \lambda \mathbf{w}_1 + (1 - \lambda) \mathbf{w}_2, \mathbf{x}_j - \mathbf{x}_r \rangle \\ &= 1 - \text{sgn}(y_j - y_r) (\langle \lambda \mathbf{w}_1, \mathbf{x}_j - \mathbf{x}_r \rangle + \langle (1 - \lambda) \mathbf{w}_2, \mathbf{x}_j - \mathbf{x}_r \rangle) \\ &= 1 - \lambda \cdot \text{sgn}(y_j - y_r) \langle \mathbf{w}_1, \mathbf{x}_j - \mathbf{x}_r \rangle + (1 - \lambda) \text{sgn}(y_j - y_r) \langle \mathbf{w}_2, \mathbf{x}_j - \mathbf{x}_r \rangle \\ &= \lambda + 1 - \lambda - \lambda \cdot \text{sgn}(y_j - y_r) \langle \mathbf{w}_1, \mathbf{x}_j - \mathbf{x}_r \rangle - (1 - \lambda) \text{sgn}(y_j - y_r) \langle \mathbf{w}_2, \mathbf{x}_j - \mathbf{x}_r \rangle \\ &= \lambda (1 - \text{sgn}(y_j - y_r) \langle \mathbf{w}_1, \mathbf{x}_j - \mathbf{x}_r \rangle) + (1 - \lambda) (1 - \text{sgn}(y_j - y_r) \langle \mathbf{w}_2, \mathbf{x}_j - \mathbf{x}_r \rangle) \\ &= \lambda f(\mathbf{w}_1) + (1 - \lambda) f(\mathbf{w}_2) \end{aligned}$$

כנדרש מפונקציה קמורה.

**מסקנה :**  $\ell(h_{\mathbf{w}}(\bar{\mathbf{x}}), \mathbf{y})$  פונקציה קמורה כנדרש.

□

## 3.2 section (b)

נוכיח כי פונקציית hinge loss לעיל חוסמת מלמעלה את Kendall-Tau loss

יהי  $\mathbf{w} \in \mathbb{R}^d$  וכן  $\bar{\mathbf{x}} \in \mathcal{X}^k$  ו  $\mathbf{y} \in \mathbb{R}^k$

נתבונן ב

$$\begin{aligned} \max \{0, 1 - \text{sgn}(y_j - y_r) \langle \mathbf{w}, \mathbf{x}_j - \mathbf{x}_r \rangle\} &= \max \{0, 1 - \text{sgn}(y_j - y_r) (y'_j - y'_r)\} \\ &= \begin{cases} t(y) & \text{sgn}(y_j - y_r) (y'_j - y'_r) < 1 \\ 0 & \text{otherwise} \end{cases} \\ &= \begin{cases} t(y) & \text{sgn}(y_j - y_r) \text{sgn}(y'_j - y'_r) |y'_j - y'_r| < 1 \\ 0 & \text{otherwise} \end{cases} \\ &= \begin{cases} 1 + p(y) & \text{sgn}(y_j - y_r) \neq \text{sgn}(y'_j - y'_r) \\ r(y) & \text{sgn}(y_j - y_r) = \text{sgn}(y'_j - y'_r) \wedge |y'_j - y'_r| < 1 \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$



כאשר  $p(y), r(y)$  פונקציות אי שליליות, ולכן מתקיים כי

$$\left( \begin{cases} 1+p(y) & \text{sgn}(y_j - y_r) \neq \text{sgn}(y'_j - y'_r) \\ r(y) & \text{sgn}(y_j - y_r) = \text{sgn}(y'_j - y'_r) \wedge |y'_j - y'_r| < 1 \\ 0 & \text{otherwise} \end{cases} \right) \leq \left( \begin{cases} 1 & \text{sgn}(y'_j - y'_r) \neq \text{sgn}(y_j - y_r) \\ 0 & \text{otherwise} \end{cases} \right)$$

אבל נשים לב כי

$$\mathbf{1} \{ \text{sgn}(y'_j - y'_r) \neq \text{sgn}(y_j - y_r) \} = \begin{cases} 1 & \text{sgn}(y'_j - y'_r) \neq \text{sgn}(y_j - y_r) \\ 0 & \text{otherwise} \end{cases}$$

כלומר מתקיים כי

$$\mathbf{1} \{ \text{sgn}(y'_j - y'_r) \neq \text{sgn}(y_j - y_r) \} \leq \max \{0, 1 - \text{sgn}(y_j - y_r) \langle \mathbf{w}, \mathbf{x}_j - \mathbf{x}_r \rangle \}$$

ולכן מתקיים :

$$\Delta(h_{\mathbf{w}}(\bar{\mathbf{x}}^i), \mathbf{y}^i) \leq \ell(h_{\mathbf{w}}(\bar{\mathbf{x}}), \mathbf{y})$$

כנדרש

□

### 3.3 section (c)

נניח כי הדטא מופרד במרווח  $\gamma > 0$  (כלומר קיים  $\mathbf{w}^* \in \mathbb{R}^d$  ו  $\gamma > 0$  כך ש  $\text{sgn}(y_j - y_r) \langle \mathbf{w}^*, \mathbf{x}_j - \mathbf{x}_r \rangle \geq \gamma$  לכל  $1 \leq j < r \leq k$ ).

נראה כי אם נמנמם את hinge loss יוביל להיפותזה שמנממת את Kendall-Tau loss

מהנתון מתקיים כי קיים  $\mathbf{w}^* \in \mathbb{R}^d$  ו  $\gamma > 0$  ולכן קיים  $\mathbf{w}^{**} = \frac{1}{\gamma} \mathbf{w}^*$  כלומר מתקיים כי  $\text{sgn}(y_j - y_r) \langle \mathbf{w}^{**}, \mathbf{x}_j - \mathbf{x}_r \rangle \geq 1 > 0$

כלומר מתקיים כי  $1 - \text{sgn}(y_j - y_r) \langle \mathbf{w}^{**}, \mathbf{x}_j - \mathbf{x}_r \rangle \leq 0$

ולכן

$$\max \{0, 1 - \text{sgn}(y_j - y_r) \langle \mathbf{w}^{**}, \mathbf{x}_j - \mathbf{x}_r \rangle \} = 0$$

כלומר למנמם את hinge loss על  $S$  פירושו

$$\sum_{i=1}^n \ell(h_{\mathbf{w}^{**}}(\bar{\mathbf{x}}), \mathbf{y}) = 0$$

מ 3.2.2 אנו מקבלים כי

$$0 \leq \sum_{i=1}^n \Delta(h_{\mathbf{w}^{**}}(\bar{\mathbf{x}}^i), \mathbf{y}^i) \leq \sum_{i=1}^n \ell(h_{\mathbf{w}^{**}}(\bar{\mathbf{x}}), \mathbf{y}) = 0$$

כלומר

$$\sum_{i=1}^n \Delta(h_{\mathbf{w}^{**}}(\bar{\mathbf{x}}^i), \mathbf{y}^i) = 0$$

עבור  $\mathbf{w}^{**}$  שנמצא ע"י המנמום של  $\sum_{i=1}^n \ell(h_{\mathbf{w}^{**}}(\bar{\mathbf{x}}), \mathbf{y})$ .

כלומר ע"י מנמום  $\sum_{i=1}^n \ell(h_{\mathbf{w}^{**}}(\bar{\mathbf{x}}), \mathbf{y})$  הצלחנו למנמום (עבור אותה פונקציה) את  $\sum_{i=1}^n \Delta(h_{\mathbf{w}^{**}}(\bar{\mathbf{x}}^i), \mathbf{y}^i)$

כנדרש.

□

## 4 Gradient Sedcent on Smooth Functions.

**הגדרה :** נאמר שפונקציה  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  המקיימת  $f \in C^1$  היא  $\beta$ -smooth אם לכל  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  מתקיים

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

תהיי  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  פונקציה  $\beta$ -smooth ואי שלילית.

נניח כי gradient descent algorithm מבוצע על  $f$  עם  $\eta > 0$ , כלומר

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

ונניח כי אנו מתחילים בנקודה  $\mathbf{x}_0$  כלשהי.

נראה כי אם מתקיים ש  $\eta < \frac{2}{\beta}$  אז מתקיים ש

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

אם כן, מהיות  $f$  פונקציה  $\beta$ -smooth ואי שלילית מתקיים עבור  $\mathbf{x}_t, \mathbf{x}_{t+1}$  כי :

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

נזכור כי  $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \eta \nabla f(\mathbf{x}_t)$  ולכן

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

כלומר

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

כלומר

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

נזכור את ההנחה כי  $0 < \eta < \frac{2}{\beta}$  ולכן מתקיים כי

$$\eta^2 \frac{\beta}{2} - \eta < 0$$

או לחלופין

$$^1 \quad \eta - \eta^2 \frac{\beta}{2} > 0$$

ולכן

$$\left(\frac{\beta}{2}\eta^2 - \eta\right)^{-1} f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) \geq \|\nabla f(\mathbf{x}_t)\|^2 \geq 0$$

הנ"ל נכון לכל  $t$ , ננצל זאת על מנת לקבל סכום טלסקופי

$$+ \begin{cases} \left(\frac{\beta}{2}\eta^2 - \eta\right)^{-1} f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) \geq \|\nabla f(\mathbf{x}_t)\|^2 \\ \left(\frac{\beta}{2}\eta^2 - \eta\right)^{-1} f(\mathbf{x}_t) - f(\mathbf{x}_{t-1}) \geq \|\nabla f(\mathbf{x}_{t-1})\|^2 \\ \vdots \\ \left(\frac{\beta}{2}\eta^2 - \eta\right)^{-1} f(\mathbf{x}_1) - f(\mathbf{x}_0) \geq \|\nabla f(\mathbf{x}_0)\|^2 \end{cases}$$

ונקבל

$$\left(\eta - \frac{\beta}{2}\eta^2\right)^{-1} (f(\mathbf{x}_0) - f(\mathbf{x}_{t+1})) \geq \sum_{t=0}^T \|\nabla f(\mathbf{x}_t)\|^2$$

כאשר נשים לב כי מהיותה של  $f$  אי שלילית וכן  $\left(\eta - \frac{\beta}{2}\eta^2\right)^{-1} > 0$  מתקיים כי

$$\left(\eta - \frac{\beta}{2}\eta^2\right)^{-1} f(\mathbf{x}_0) \geq \sum_{t=0}^T \|\nabla f(\mathbf{x}_t)\|^2$$

כלומר לכל  $T$ 

$$positive - const \geq \sum_{t=0}^T \|\nabla f(\mathbf{x}_t)\|^2$$

נסמן

$$S_T = \sum_{t=0}^T \|\nabla f(\mathbf{x}_t)\|^2$$

מכך ש  $\|\nabla f(\mathbf{x}_t)\|^2 \geq 0$  זו סדרה מונוטונית עולה, והראנו כי חסומה, ולכן מתכנסת. כלומר

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

ממשפט בקורס חדו"א אנו יודעים שכיוון ש  $\|\nabla f(\mathbf{x}_t)\|^2 \geq 0$  וכן הטור מתכנס, מתקיים כי

$$\lim_{t \rightarrow \infty} \|\nabla f(\mathbf{x}_t)\|^2 = 0$$

נשים לב כי

$$\|\nabla f(\mathbf{x}_t)\| = \sqrt{\|\nabla f(\mathbf{x}_t)\|^2} \geq 0$$

ולכן ממשפט בקורס חדוו"א :

**משפט :** אם  $f$  רציפה ומתקיים כי  $\lim_{n \rightarrow \infty} a_n = L$  אזי  $\lim_{n \rightarrow \infty} f(a_n) = f(L)$

ולכן, כיוון ש  $f(x) = \sqrt{x}$  פונקציה רציפה מתקיים כי

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

כנדרש.

□

# Programming Assignment

## 5 SGD for Hinge loss.

### 5.1 (a)

כדי לבחור את  $\eta_0$  הטוב ביותר, ניתן לאלגוריתם ערכים שהחלו ב

$$(10^{-5}, 10^{-4}, \dots, 10^4, 10^5)$$

$\eta_0$  הטוב ביותר שהתקבל הוא  $\eta_0 = 1$  עם  $\text{accuracy} = 0.979$

נתמקד באזור בשבו התקבלה התוצאה הטובה ביותר, והאלגוריתם הורץ שנית (ברזולוציית קפיצות קטנה יותר)

$$(-9999, -9899, \dots, 9801, 9901)$$

$\eta_0$  הטוב ביותר שהתקבל הוא  $\eta_0 = 1$  עם  $\text{accuracy} = 0.977$

נתמקד באזור בשבו התקבלה התוצאה הטובה ביותר, והאלגוריתם הורץ שלישית (ברזולוציית קפיצות קטנה יותר)

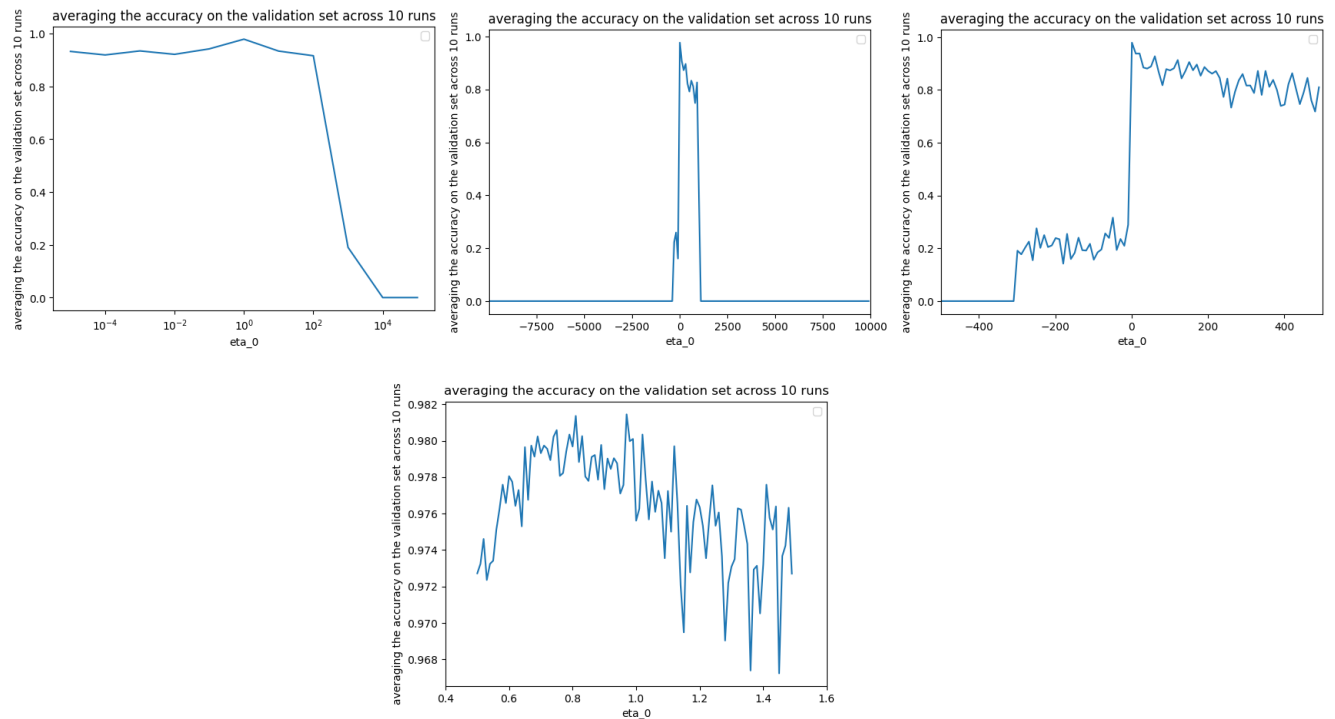
$$(-499, -489, \dots, 481, 491)$$

$\eta_0$  הטוב ביותר שהתקבל הוא  $\eta_0 = 1$  עם  $\text{accuracy} = 0.979$

לאחר שהובן כח הערכים מתכנסים ל  $\eta_0 = 1$  הורץ בפעם הרביעית על

$$(0.5, 0.51, \dots, 1.49, 1.5)$$

$\eta_0$  הטוב ביותר שהתקבל הוא  $\eta_0 = 1$  עם  $\text{accuracy} = 0.981$

Figure 1: accuracy as function of  $\eta_0$ 

## 5.2 (b)

כדי לבחור את  $C$  הטוב ביותר, ניתן לאלגוריתם ערכים שהחלו ב

$$(10^{-5}, 10^{-4}, \dots, 10^4, 10^5)$$

$C$  הטוב ביותר שהתקבל הוא  $C = 10^{-4}$  עם  $\text{accuracy} = 0.986$

נתמקד באזור בשבו התקבלה התוצאה הטובה ביותר, והאלגוריתם הורץ שנית (ברזולוציית קפיצות קטנה יותר)

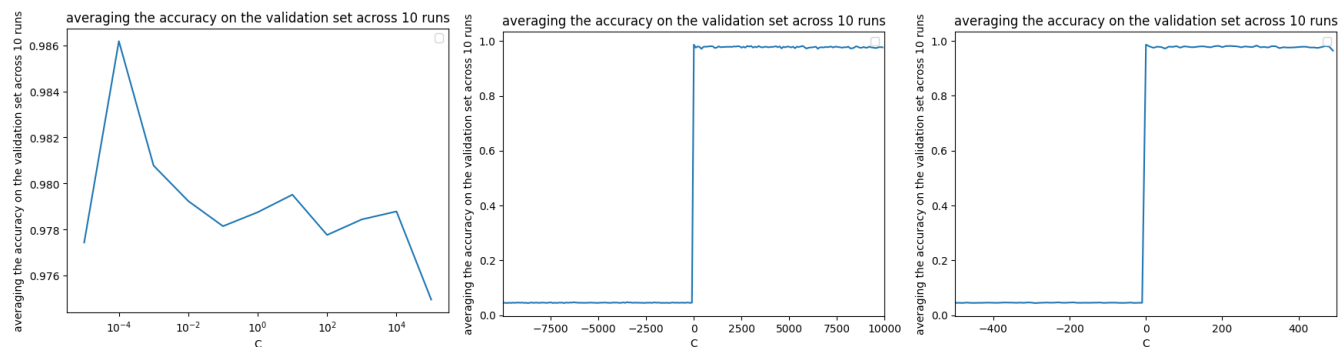
$$(-9999.9999, -9899.9999, \dots, 9800.0001, 9900.0001)$$

$C$  הטוב ביותר שהתקבל הוא  $C = 0.001$  עם  $\text{accuracy} = 0.987$

נתמקד באזור בשבו התקבלה התוצאה הטובה ביותר, והאלגוריתם הורץ שלישית (ברזולוציית קפיצות קטנה יותר)

$$(-499.9999, -489, \dots, 480.0001, 490.0001)$$

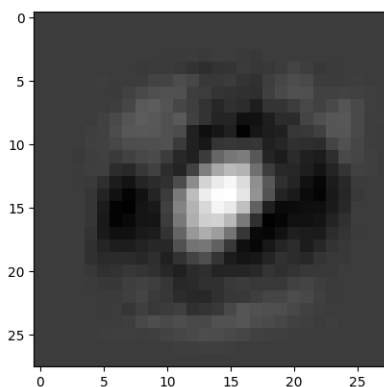
$C$  הטוב ביותר שהתקבל הוא  $C = 1$  עם  $\text{accuracy} = 0.986$

Figure 2: accuracy as function of  $C$ 

### 5.3 (c)

ניתן לראות כי הפיקסלים הלבנים באמצע יכולים להעיד על אמצע הספרה 8 והפיקסלים בצבעי לבן - אפור בצדדים מעידים על הספרה 8 גם כן.

הפיקסלים השחורים מעידים על הספרה 0.

Figure 3:  $w$  as an image

### 5.4 (d)

על סט הבדיקה התקבל :

$$accuracy = 0.992835$$



## 6 SGD for log-loss.

### 6.1 (a)

כדי לבחור את  $\eta_0$  הטוב ביותר, ניתן לאלגוריתם ערכים שהחלו ב

$$(10^{-5}, 10^{-4}, \dots, 10^4, 10^5)$$

$\eta_0$  הטוב ביותר שהתקבל הוא  $\eta_0 = 10^{-5}$  עם  $\text{accuracy} = 0.955$

נתמקד באזור בשבו התקבלה התוצאה הטובה ביותר, והאלגוריתם הורץ שנית (ברזולוציית קפיצות קטנה יותר)

$$(-9999, -9899, \dots, 9801, 9901)$$

$\eta_0$  הטוב ביותר שהתקבל הוא  $\eta_0 = 10^{-5}$  עם  $\text{accuracy} = 0.959$

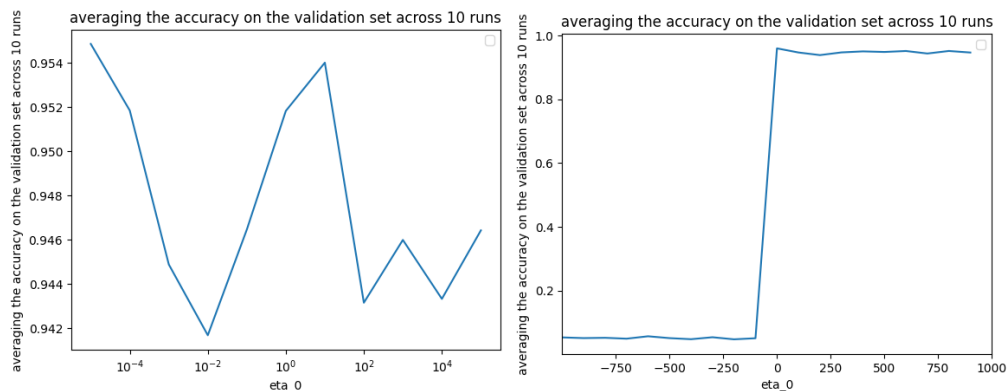


Figure 4: accuracy as function of  $\eta_0$

### 6.2 (b)

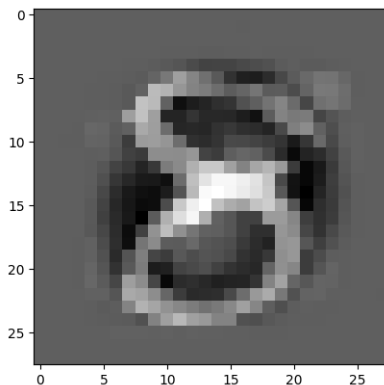


Figure 5:  $\mathbf{w}$  as an image

על סט הבדיקה התקבל :

$$accuracy = 0.973$$

### 6.3 (c)

ניתן הסבר :

- ההגעה לערך המינימלי של פונקציית log-loss ע"י אלגוריתם SGD מתבצע מהר, ולכן  $w$  "נלכד" בערך המינימלי ולכן  $\|w_t\|_2$  אינו משתנה לאחר האיטרציות הראשונות יחסית ל 20,000.
- בכל שלב  $\eta$  קטן, ולכן השינוי ב  $\|w_t\|_2$ , שתלוי בו, יקטן ככל ש  $t$  גדל.

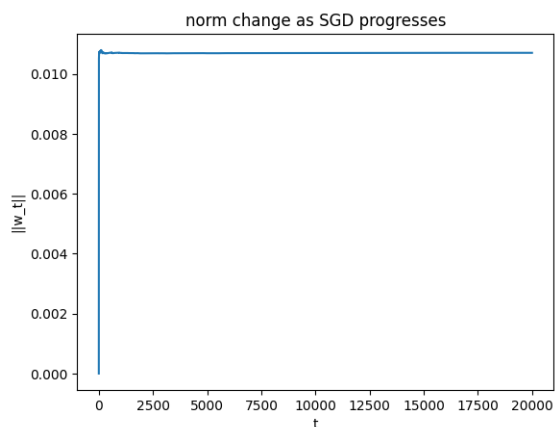


Figure 6:  $\|w_t\|_2$  as function of  $t$