

# 3. מינימיזציה של פונקציית האפס-

1 נסכה

1. (15 points) Step-size Perceptron. Consider the modification of Perceptron algorithm with the following update rule:

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

whenever  $\hat{y}_t \neq y_t$  ( $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t$  otherwise). Assume that data is separable with margin  $\gamma > 0$  and that  $\|\mathbf{x}_t\| = 1$  for all  $t$ . For simplicity assume that the algorithm makes  $M$  mistakes at the first  $M$  rounds, after which it has no mistakes. For  $\eta_t = \frac{1}{\sqrt{t}}$ , show that the number of mistakes step-size Perceptron makes is at most  $\frac{4}{\gamma^2} \log(\frac{1}{\gamma})$ . (Hint: use the fact that if  $x \leq a \log(x)$  then  $x \leq 2a \log(a)$ ). It's okay if you obtain a bound with slightly different constants, but the asymptotic dependence on  $\gamma$  should be tight.

$$\begin{aligned} \mathbf{w}_{t+1} \cdot \mathbf{w}^* &= (\mathbf{w}_t + \eta_t y_t \mathbf{x}_t) \cdot \mathbf{w}^* = (\mathbf{w}_t + \frac{1}{\sqrt{t}} y_t \mathbf{x}_t) \cdot \mathbf{w}^* \\ &= \mathbf{w}_t \cdot \mathbf{w}^* + \frac{1}{\sqrt{t}} y_t \mathbf{x}_t \cdot \mathbf{w}^* \geq \mathbf{w}_t \cdot \mathbf{w}^* + \frac{1}{\sqrt{t}} \cdot \gamma \\ &\quad \downarrow \\ &y_t \mathbf{x}_t \cdot \mathbf{w}^* \geq \gamma \end{aligned}$$

מכ. מוגאית נסכה:

$$\mathbf{w}_t \cdot \mathbf{w}^* \geq \sum_{t=1}^M \frac{\gamma}{\sqrt{t}} = \gamma \cdot \sum_{t=1}^M \frac{1}{\sqrt{t}} \geq \gamma \cdot \frac{M}{\sqrt{M}} = \gamma \cdot \sqrt{M}$$

:  $\|\mathbf{w}_t\|_2 \leq \gamma$  סע  
: גודל מדויק

$$\begin{aligned} \|\mathbf{w}_{t+1}\|_2^2 &= \|\mathbf{w}_t + \frac{1}{\sqrt{t}} y_t \mathbf{x}_t\|_2^2 = \|\mathbf{w}_t\|_2^2 + \underbrace{\frac{2}{\sqrt{t}} y_t \mathbf{x}_t \cdot \mathbf{w}_t + \frac{1}{t} \|\mathbf{x}_t\|^2}_{\text{טכ. שגיאה}} \\ &\leq \|\mathbf{w}_t\|_2^2 + \frac{1}{t} \|\mathbf{x}_t\|^2, \quad \|\mathbf{x}_t\|=1 \end{aligned}$$

$$\begin{aligned} \|\mathbf{w}_{t+1}\|_2^2 &\leq \sum_{t=1}^M \frac{1}{t} \quad \text{מכ. מוגאית נסכה:} \\ &\leq \sum_{t=1}^M \frac{1}{t} \quad \text{: סכום} \end{aligned}$$

$$\|w_{t+1}\|_2 \leq \sqrt{\sum_{t=1}^T \frac{1}{t}} \underset{\text{using } \sum t \geq T}{\approx} \sqrt{\log(\mu)}$$

$$w_t \cdot w^* \geq \gamma \sqrt{\mu}, \quad \|w_t\|_2 \leq \sqrt{\log(\mu)} \quad \|w^*\|_2 = 1$$

$$\gamma \sqrt{\mu} \leq w_t \cdot w^* \leq \|w_t\|_2 \|w^*\|_2 \leq \sqrt{\log(\mu)} \cdot 1$$

$$\gamma \sqrt{\mu} \leq \sqrt{\log(\mu)} \Rightarrow \sqrt{\mu} \leq \frac{1}{\gamma} \sqrt{\log(\mu)}$$

$$M \leq \frac{1}{\gamma^2} \log(\mu) \Rightarrow M \leq \frac{2}{\gamma^2} \cdot \log\left(\frac{1}{\gamma^2}\right)$$

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

$$\log_b(x^2) = 2\log_b(x)$$

correct

2 inde

2. (15 points) Convex functions.

- (a) Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  a convex function,  $A \in \mathbb{R}^{n \times n}$  and  $b \in \mathbb{R}^n$ . Show that,  $g(\mathbf{x}) = f(A\mathbf{x} + b)$  is convex.
- (b) Consider  $m$  convex functions  $f_1(\mathbf{x}), \dots, f_m(\mathbf{x})$ , where  $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ . Now define a new function  $g(\mathbf{x}) = \max_i f_i(\mathbf{x})$ . Prove that  $g(\mathbf{x})$  is a convex function. (Note that from (a) and (b) you can conclude that the hinge loss over linear classifiers is convex.)
- (c) Let  $\ell_{\log} : \mathbb{R} \rightarrow \mathbb{R}$  be the log loss, defined by

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

Show that  $\ell_{\log}$  is convex, and conclude that the function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  defined by  $f(\mathbf{w}) = \ell_{\log}(y\mathbf{w} \cdot \mathbf{x})$  is convex with respect to  $\mathbf{w}$ .

: $\exists \lambda \in [0, 1] \quad x, y \in \mathbb{R}^n$  וו (a)

$$\begin{aligned} g(\lambda x + (1-\lambda)y) &= f(A(\lambda x + (1-\lambda)y) + b) \\ &= f(A\lambda x + A(1-\lambda)y + b) = f(\lambda Ax + (1-\lambda)Ay + b) \\ &= f(\lambda Ax + (1-\lambda)Ay + ((1-\lambda)+\lambda)b) \\ &= f(\lambda(Ax+b) + (1-\lambda)(Ay+b)) \end{aligned}$$

הנ"מ הינה פונקציית  $Ax+b, Ay+b \in \mathbb{R}^n$  ו "פונקציית

$$\leq \lambda f(Ax+b) + (1-\lambda)f(Ay+b) = \lambda g(x) + (1-\lambda)g(y)$$

: $\lambda \in [0, 1] \quad x, y \in \mathbb{R}^n$  הינה פונקציית

$$g(\lambda x + (1-\lambda)y) \leq \lambda g(x) + (1-\lambda)g(y)$$

הנ"מ הינה פונקציית

:  $f_i$  סדר גיאומטרי  $\forall \lambda \in [0, 1] \times y \in \mathbb{R}^n$  (b)

$$\forall i \in [n] \quad f_i(\lambda x + (1-\lambda)y) \leq \lambda f_i(x) + (1-\lambda) f_i(y)$$

: יסוד

$$g(\lambda x + (1-\lambda)y) = \max_i f_i(\lambda x + (1-\lambda)y)$$

$$\leq \max_i \{\lambda f_i(x) + (1-\lambda) f_i(y)\}$$

$$\leq \lambda \max_i \{f_i(x)\} + (1-\lambda) \max_i \{f_i(y)\}$$

$$\leq \lambda g(x) + (1-\lambda) g(y)$$

. הינה  $g$  פולינומיאלית

: פונקציית  $\ln$  מתקיימת (c)

$$(\ln)'(2) = \frac{1}{\ln(2) \cdot (1+e^{-2})} \circ -e^{-2} = -\frac{e^{-2}}{\ln(2) \cdot (1+e^{-2})}$$

$$= -\ln(2)^{-1} \cdot \frac{e^{-2}}{1+e^{-2}}$$

$$(\ln)''(2) = -\ln(2)^{-1} \cdot \left[ \frac{-e^{-2} (1+e^{-2})^{-2} - e^{-2} \cdot (-e^{-2})}{(1+e^{-2})^2} \right]$$

$$= -\ln(2)^{-1} \cdot \left[ \frac{-e^{-2} - e^{-2} + e^{-2}}{(1+e^{-2})^2} \right] = -\ln(2)^{-1} \cdot \frac{-e^{-2}}{(1+e^{-2})^2}$$

$$-\ln(2)^{-1} < 0, -e^{-2} < 0, (1+e^{-2})^2 > 0$$

פונקציית

הנגמיה כפולה פולינומיאלית

$$(\ln)''(2) > 0$$

(יעילו לנו)  $\ln$  פולינומיאלי

CND:

$$f(w) = \log(y \cdot w \cdot x) = \log((yx) \cdot w)$$

$x, w \in \mathbb{R}^n$

:  $\exists \lambda \in [0, 1] z_1, z_2 \in \mathbb{R}^n$

$$f(\lambda z_1 + (1-\lambda)z_2) = \log((yx) \cdot (\lambda z_1 + (1-\lambda)z_2))$$

$$= \log(yx \lambda z_1 + yx(1-\lambda)z_2)$$

$$= \log(\lambda(yxz_1) + (1-\lambda)(yxz_2))$$

לNCI מילוי הדרישה ל-NLI:

$$\leq \lambda \cdot \log(yxz_1) + (1-\lambda) \log(yxz_2)$$

$$= \lambda \cdot f(z_1) + (1-\lambda) \cdot f(z_2)$$

לפנינו קיומם כפונקציית f.

3 נס

3. (20 points) **Ranking.** In this question, we consider a new learning task in which the objective is to rank items. Assume items are elements of  $\mathcal{X} \subseteq \mathbb{R}^d$ , and you are given a training set of  $n$  lists of  $k$  items each, and for each list you receive a “label” vector corresponding to the correct ranking of its items. More formally, you receive a training set

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

such that for all  $1 \leq i \leq n$ ,  $\mathbf{y}^i \in \mathbb{R}^k$  assigns a value for each item in  $\bar{\mathbf{x}}^i = (\mathbf{x}_1^i, \mathbf{x}_2^i, \dots, \mathbf{x}_k^i)$ , interpreted as a ranking of the items. Your goal is to learn a ranking function  $h : \mathcal{X}^k \rightarrow \mathbb{R}^k$  which correctly ranks the lists of items from  $S$ . The *Kendall-Tau* loss between two rankings  $\mathbf{y}', \mathbf{y}$  is defined as follows:

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

Note that this function averages the total number of pairs of items which are in different order in  $\mathbf{y}'$  compared to  $\mathbf{y}$ . Assume you are trying to learn a linear ranking function, i.e. a function of the form

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

for some  $\mathbf{w} \in \mathbb{R}^d$ , and your goal is to minimize the Kendall-Tau loss over  $S$ :  $\sum_{i=1}^n \Delta(h_{\mathbf{w}}(\bar{\mathbf{x}}^i), \mathbf{y}^i)$ . Since this function is hard to optimize, you instead optimize the surrogate “hinge” loss  $\sum_{i=1}^n \ell(h_{\mathbf{w}}(\bar{\mathbf{x}}^i), \mathbf{y}^i)$  where:

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

- (a) Prove that the hinge loss described above for the ranking objective is convex in  $\mathbf{w}$ .
- (b) Prove that the hinge loss upper-bounds the Kendall-Tau loss, i.e. that  $\Delta(h_{\mathbf{w}}(\bar{\mathbf{x}}), \mathbf{y}) \leq \ell(h_{\mathbf{w}}(\bar{\mathbf{x}}), \mathbf{y})$  for all  $\mathbf{w} \in \mathbb{R}^d$ ,  $\bar{\mathbf{x}} \in \mathcal{X}^k$ ,  $\mathbf{y} \in \mathbb{R}^k$ .
- (c) Prove that if the data is separable with a margin  $\gamma > 0$  (i.e. when there exists  $\mathbf{w}^* \in \mathbb{R}^d$  and  $\gamma > 0$  such that  $sgn(y_j^i - y_r^i) \mathbf{w}^* \cdot (\mathbf{x}_j^i - \mathbf{x}_r^i) \geq \gamma$  for all  $1 \leq i \leq n$  and all  $1 \leq j < r \leq k$ ), minimizing the hinge loss will result in a ranking function which minimizes the Kendall-Tau loss.

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

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

ההinge loss מוגדר כ  $\max\{0, 1 - sgn(y_j - y_r) \mathbf{w} \cdot (\mathbf{x}_j - \mathbf{x}_r)\}$

$$\max\{0, 1 - sgn(y_j - y_r) \mathbf{w} \cdot (\mathbf{x}_j - \mathbf{x}_r)\}$$

ההinge loss מוגדר כ  $\max\{0, 1 - sgn(y_j - y_r) \mathbf{w} \cdot (\mathbf{x}_j - \mathbf{x}_r)\}$

$$1 - \text{sgn}(y_j - y_r) w \cdot (x_j - x_r)$$

$$= (x_j - x_r) \cdot (-\text{sgn}(y_j - y_r)) w + 1$$

ז"ה מילוי המודולוס נקבע על ידי ק.ג.ה.

פונקציית ה- $\text{sgn}$  מוגדרת כ

פונקציית ניטול אנטוואט של הטענה ב"ט כפולה  
בנוסף ל $h_w(x), y$  הינה מוגדרת כ $h_w(x^i), y'$ .

לפניהם ניקח סכום על כל ה- $y'$  וונציאנו

(ב) כוכב נספחים:

$$\Delta(h_w(\bar{x}), y) = \frac{2}{K(K-1)} \sum_{j=1}^{K-1} \sum_{r=j+1}^{K-1} \mathbb{1} \{ \text{sgn}(h_w(\bar{x})_j - h_w(\bar{x})_r) \neq \text{sgn}(y_j - y_r) \}$$

$$l(h_w(\bar{x}, y)) = \frac{2}{K(K-1)} \sum_{j=1}^{K-1} \sum_{r=j+1}^{K-1} \max \{0, 1 - \text{sgn}(y_j - y_r) w \cdot (h_w(\bar{x})_j - h_w(\bar{x})_r) \}$$

כעת נראה גנרייט:

$$\text{מ长时间 } \text{sgn}(y_j - y_r) w \cdot (h_w(\bar{x})_j - h_w(\bar{x})_r) \geq 0 \text{ פlc(1)}$$

$$0 \leq \text{sgn}(y_j - y_r) = \text{sgn}(h_w(\bar{x})_j - h_w(\bar{x})_r)$$

.0 \leq l(h\_w(\bar{x}), y) \leq \Delta(h\_w(\bar{x}), y)

$$\text{מ长时间 } x = \text{sgn}(y_j - y_r) w \cdot (h_w(\bar{x})_j - h_w(\bar{x})_r) < 0 \text{ px(2)}$$

$$1 - \text{sgn}(y_j - y_r) \neq \text{sgn}(h_w(\bar{x})_j - h_w(\bar{x})_r)$$

:  
לפנינו  $l(h_w(\bar{x}), y) \leq \Delta(h_w(\bar{x}), y)$

$$\max \{0, 1 - x\} = \begin{cases} 0, & x \leq 0 \\ 1 + y, & x > 0 \end{cases} = 1 + y$$

פונקציית סטטיסטיקות.

נוכיח כי ניתן ללו נחיהה בסכום

נוכיח שסכום  $\frac{2}{K(K-1)} \sum_{j=1}^{K-1} \sum_{r=j+1}^{K-1} l(h_w(\bar{x}), y)$  הוא מינימום של פונקציית האנרגיה.

זהו:

$$\Delta(h_w(\bar{x}), y) \leq l(h_w(\bar{x}), y) \quad \forall \bar{x} \in \mathbb{R}^C, \bar{x} \in \mathcal{X}^K, y \in \mathbb{R}^K$$

done

לפניכם נתונים (C)

$\exists w \in R^k$ ,  $\gamma > 0$  and:

$$\operatorname{sgn}(y_j^i - y_r^i) w^* \cdot (x_j^i - x_r^i) \geq \gamma > 0 \quad (\text{ונתן } \gamma > 0)$$

$\forall 1 \leq i \leq n, 1 \leq j < r < k$

$\forall 1 \leq i \leq n, 1 \leq j < r < k$

$$\text{לפניכם } \bar{w} = \frac{1}{\gamma} w^*$$

$$\operatorname{sgn}(y_j^i - y_r^i) \bar{w} \cdot (x_j^i - x_r^i) \geq 1$$

ו今 הינה מושג ה-3 יריעתית כפולה:

$$\forall i: \max \{0, 1 - \operatorname{sgn}(y_j^i - y_r^i) \bar{w} \cdot (x_j^i - x_r^i)\} = 0$$

$$\operatorname{sgn}(y_j^i - y_r^i) \bar{w} \cdot (x_j^i - x_r^i) \geq 1$$

$$1 - \operatorname{sgn}(y_j^i - y_r^i) \bar{w} \cdot (x_j^i - x_r^i) \leq 0$$

ולכן

$$\max \{0, 1 - \operatorname{sgn}(y_j^i - y_r^i) \bar{w} \cdot (x_j^i - x_r^i)\} = 0$$

ה- hinge loss הינו תואם  $w_{\min}$  (zero loss). על כן יתאפשר (zero loss)  $w_{\min}$ .

: (P) נתונים  $\gamma$  ו-  $w_{\min}$

$$\text{zero hinge loss} \Rightarrow \operatorname{sgn}(y_j^i - y_r^i) w_{\min} \cdot (x_j^i - x_r^i) \geq 1$$

$$\Rightarrow \operatorname{sgn}(y_j^i - y_r^i) \cdot (x_j^i - x_r^i) \cdot w_{\min} \geq 1$$

כל נר או מינימום של  $(y_j^i - y_r^i)$

daten separation  $w_{\min}$  מוגדר על ידי  $(x_j^i - x_r^i) \cdot w_{\min}$

המינימום של המרחק בין נר לנתון

$$\Delta(h_w(\bar{x}), y) \leq \ell(h_w(\bar{x}), y) \text{ for all } w \in \mathbb{R}^d, \bar{x} \in \mathcal{X}^k, y \in \mathbb{R}^k.$$

kendall-tau loss of  $y'$  ו-  $y$  מוגדר כפונקציית מילוי

hinge loss = רציך:

$$\sum_{i=1}^n \ell(h_w(\bar{x}^i), y^i) \text{ where:}$$

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

kendall-tau loss =

$$\sum_{i=1}^n \Delta(h_w(\bar{x}^i), y^i). \text{ where:}$$

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

4. (15 points) Gradient Descent on Smooth Functions.

We say that a continuously differentiable function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $\beta$ -smooth if for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$

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

In words,  $\beta$ -smoothness of a function  $f$  means that at every point  $\mathbf{x}$ ,  $f$  is upper bounded by a quadratic function which coincides with  $f$  at  $\mathbf{x}$ .

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a  $\beta$ -smooth and non-negative function (i.e.,  $f(\mathbf{x}) \geq 0$  for all  $\mathbf{x} \in \mathbb{R}^n$ ). Consider the (non-stochastic) gradient descent algorithm applied on  $f$  with constant step size  $\eta > 0$ :

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

Assume that gradient descent is initialized at some point  $\mathbf{x}_0$ . Show that if  $\eta < \frac{2}{\beta}$  then

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

(Hint: Use the smoothness definition with points  $\mathbf{x}_{t+1}$  and  $\mathbf{x}_t$  to show that  $\sum_{t=0}^{\infty} \|\nabla f(\mathbf{x}_t)\|^2 < \infty$  and recall that for a sequence  $a_n \geq 0$ ,  $\sum_{n=1}^{\infty} a_n < \infty$  implies  $\lim_{n \rightarrow \infty} a_n = 0$ . Note that  $f$  is not assumed to be convex!)

הה הינה מושג-ב גודל וריאנט של פונקציית  $f$  ב-

ב:

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

$$\therefore \mathbf{x}_{t+1} = \mathbf{x}_t - \eta \nabla f(\mathbf{x}_t) \text{ ס. 3)$$

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

$$\text{נאר}$$

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

$$= -\eta \nabla f(\mathbf{x}_t)^T (\nabla f(\mathbf{x}_t)) + \eta^2 \cdot \frac{\beta}{2} \|\nabla f(\mathbf{x}_t)\|^2$$

$$= -\eta \|\nabla f(\mathbf{x}_t)\|^2 + \eta^2 \cdot \frac{\beta}{2} \|\nabla f(\mathbf{x}_t)\|^2$$

רניל:

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

$$f(x_t) - f(x_{t+1}) \geq \|\nabla f(x_t)\|^2 \cdot n \left( 1 - \frac{n\beta}{2} \right)$$

$$\sum_{t=0}^{\infty} \|\nabla f(x_t)\|^2 \cdot n \left( 1 - \frac{n\beta}{2} \right) = n \left( 1 - \frac{n\beta}{2} \right) \cdot \sum_{t=0}^{\infty} \nabla f(x_t)$$

$$\leq \sum_{t=0}^{\infty} f(x_t) - f(x_{t+1}) = f(x_0) - f(x_{T+1}) \quad T \in \mathbb{N}$$

וילג' פול'

$$n \left( 1 - \frac{n\beta}{2} \right) \geq 0 \text{ ו } 1 - \frac{n\beta}{2} > 0 \text{ מכאן } 0 \leq n < \frac{2}{\beta}$$

: סכין

$$\sum_{t=0}^{\infty} n \left( 1 - \frac{n\beta}{2} \right) \|\nabla f(x_t)\|^2 \leq f(x_0) - f(x_{T+1})$$

$$\Rightarrow \sum_{t=0}^{\infty} \|\nabla f(x_t)\|^2 \leq \frac{f(x_0) - f(x_{T+1})}{n \left( 1 - \frac{n\beta}{2} \right)}$$

ככל שהפונקציה יגדל

$$\leq \frac{f(x_0)}{n \left( 1 - \frac{n\beta}{2} \right)} = C \underset{n \rightarrow \infty}{\sim} < \infty$$

וכאן חישוב.

כואז  $\lim_{t \rightarrow \infty} \|\nabla f(x_t)\|^2 \geq 0$  (ווארה כראיה) יפל ג'ס. הראז

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

(ו.ג)

סolutions

1 exercise

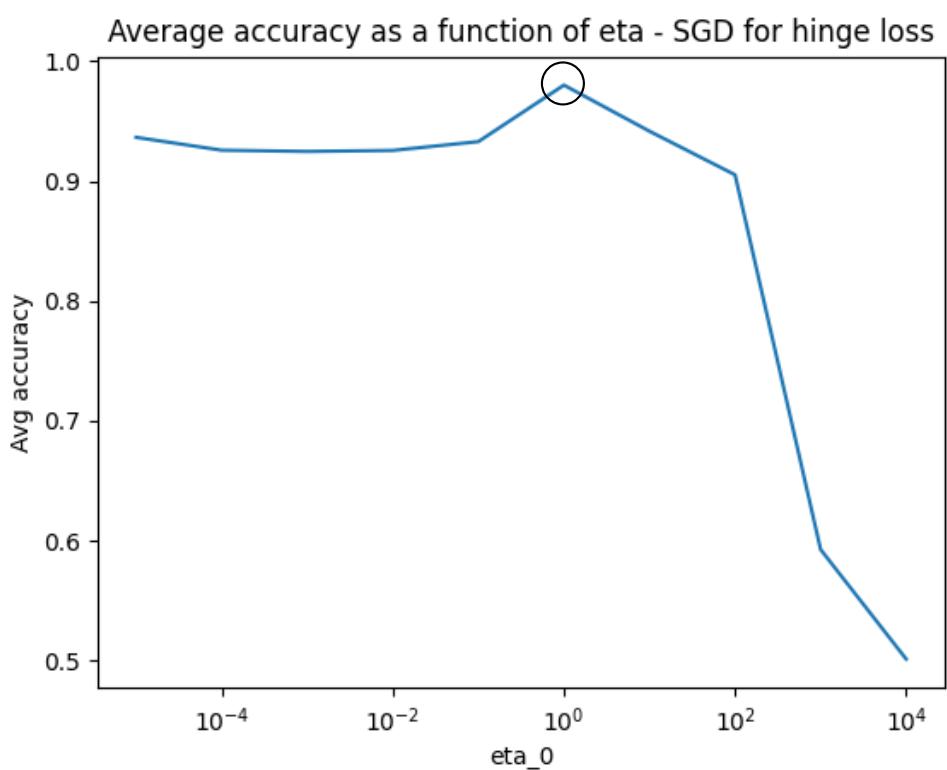
1. (20 points) SGD for Hinge loss. We will continue working with the MNIST data set. The file template (`skeleton_sgd.py`), contains the code to load the training, validation and test sets for the digits 0 and 8 from the MNIST data. In this exercise we will optimize the Hinge loss with  $L_2$ -regularization ( $\ell(\mathbf{w}, \mathbf{x}, y) = C(\max\{0, 1 - y\langle \mathbf{w}, \mathbf{x} \rangle\}) + 0.5\|\mathbf{w}\|^2$ ), using the stochastic gradient descent implementation discussed in class. Namely, we initialize  $\mathbf{w}_1 = 0$ , and at each iteration  $t = 1, \dots$  we sample  $i$  uniformly; and if  $y_i \mathbf{w}_t \cdot \mathbf{x}_i < 1$ , we update:

$$\mathbf{w}_{t+1} = (1 - \eta_t)\mathbf{w}_t + \eta_t C y_i \mathbf{x}_i$$

and  $\mathbf{w}_{t+1} = (1 - \eta_t)\mathbf{w}_t$  otherwise, where  $\eta_t = \eta_0/t$ , and  $\eta_0$  is a constant. Implement an SGD function that accepts the samples and their labels,  $C$ ,  $\eta_0$  and  $T$ , and runs  $T$  gradient updates as specified above. In the questions that follow, make sure your graphs are meaningful. Consider using `set_xlim` or `set_ylim` to concentrate only on a relevant range of values.

- (a) (5 points) Train the classifier on the training set. Use cross-validation on the validation set to find the best  $\eta_0$ , assuming  $T = 1000$  and  $C = 1$ . For each possible  $\eta_0$  (for example, you can search on the log scale  $\eta_0 = 10^{-5}, 10^{-4}, \dots, 10^4, 10^5$  and increase resolution if needed), assess the performance of  $\eta_0$  by averaging the accuracy on the validation set across 10 runs. Plot the average accuracy on the validation set, as a function of  $\eta_0$ .

רשות (Q)

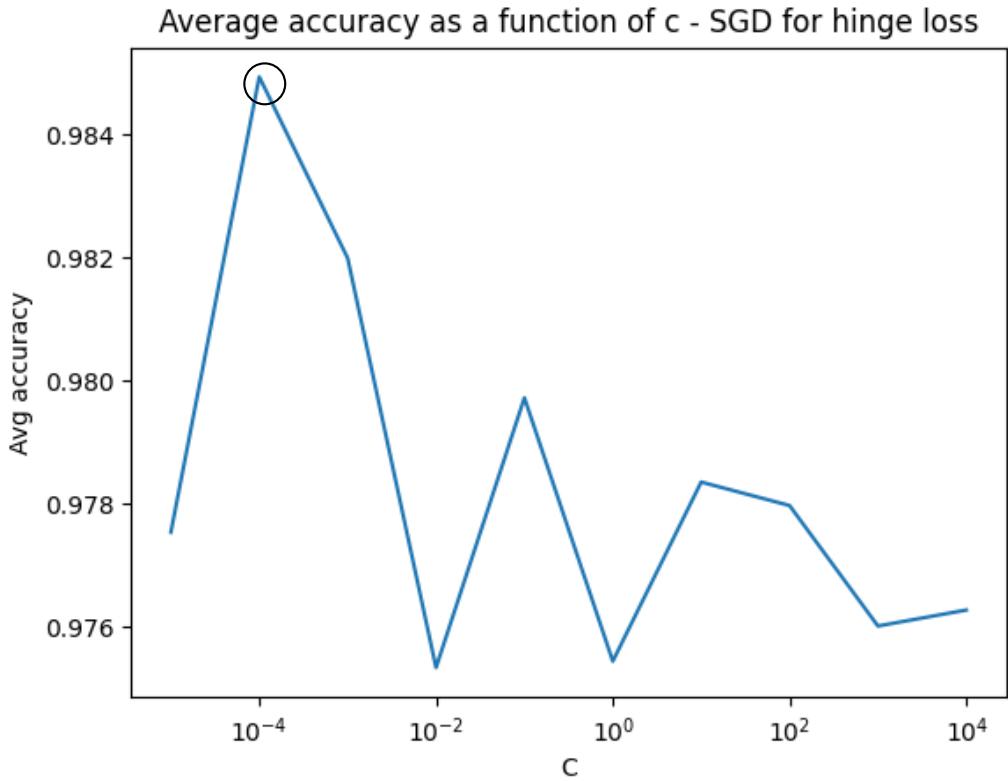


The best eta is: 1

השאלה - מתקן הדרישה לארון פולינומי (degree)  $n=1$  ותבונת מינימום נגדי (minimum)

- (b) (5 points) Now, cross-validate on the validation set to find the best  $C$  given the best  $\eta_0$  you found above. For each possible  $C$  (again, you can search on the log scale as in section (a)), average the accuracy on the validation set across 10 runs. Plot the average accuracy on the validation set, as a function of  $C$ .

:DR) (b)



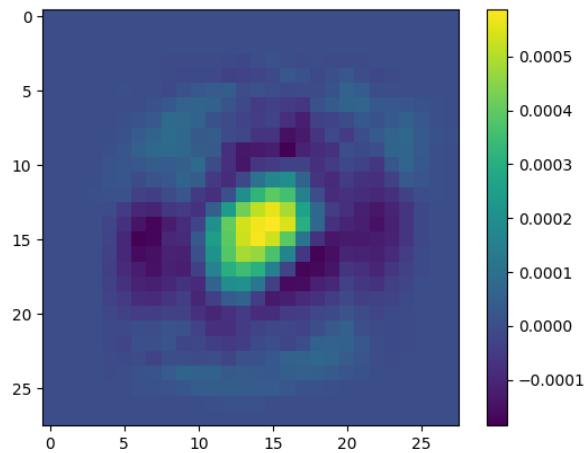
לירן מילס שיטות פרמיינט

The best  $C$  is: 0.0001

מונטגנו

- (c) (5 points) Using the best  $C$ ,  $\eta_0$  you found, train the classifier, but for  $T = 20000$ . Show the resulting  $\mathbf{w}$  as an image, e.g. using the following `matplotlib.pyplot` function: `imshow(reshape(image, (28, 28)), interpolation='nearest')`. Give an intuitive interpretation of the image you obtain.

(c) מבחן



אנו מודגש כי מבחן פותח על מנת לבדוק הדריך לא בדיקות  
סבירות פלטטן כבוד 8 אמינו הדריך הוכח והם מוגדרים  
סבירות פלטטן כבוד 0.

- (d) (5 points) What is the accuracy of the best classifier on the test set?

The accuracy of the best classifier is: 99.2835209825998 %

(d)

2. (15 points) SGD for log-loss. In this exercise we will optimize the log loss defined as follows:

$$\ell_{\log}(\mathbf{w}, \mathbf{x}, y) = \log(1 + e^{-y\mathbf{w}\cdot\mathbf{x}})$$

(in the lecture you defined the loss with  $\log_2(\cdot)$ , but for optimization purposes the logarithm base doesn't matter). Derive the gradient update for this case, and implement the appropriate SGD function.

- In your computations, it is recommended to use `scipy.special.softmax` to avoid numerical issues which arise from exponentiating very large numbers.

- (a) (5 points) Train the classifier on the training set. Use cross-validation on the validation set to find the best  $\eta_0$ , assuming  $T = 1000$ . For each possible  $\eta_0$  (for example, you can search on the log scale  $\eta_0 = 10^{-5}, 10^{-4}, \dots, 10^4, 10^5$  and increase resolution if needed), assess the performance of  $\eta_0$  by averaging the accuracy on the validation set across 10 runs. Plot the average accuracy on the validation set, as a function of  $\eta_0$ .

: גראף סלולרי (q)

$$\frac{\partial \ell_{\log}}{\partial w_i} = \frac{\partial (\log(1 + e^{-y_i w \cdot x}))}{\partial w_j}$$

$$= \frac{1}{1 + e^{-y_i w \cdot x}} \cdot e^{-y_i w \cdot x} \cdot (-y_i x_i)$$

כינן:

$$1 - \frac{1}{1 + e^{-y_i w \cdot x}} = 1 - \text{expit}(y_i w \cdot x) = \frac{1 + e^{-y_i w \cdot x}}{1 + e^{-y_i w \cdot x}} - 1$$

$$= \frac{1}{1 + e^{-y_i w \cdot x}} \cdot e^{-y_i w \cdot x}$$

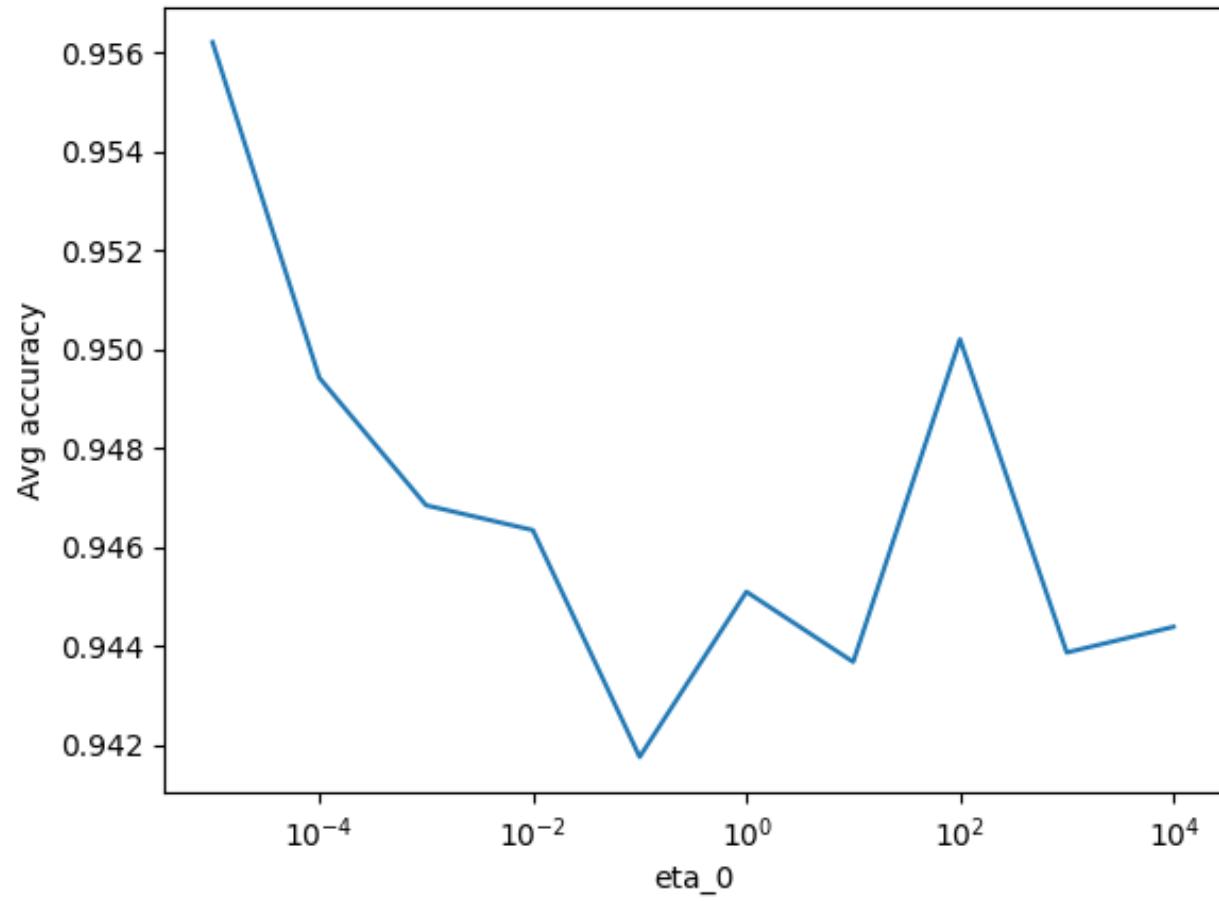
: סטוד, expit מושגתו אוסף של נספחים פלט

$$\nabla \ell_{\log} = \left( \frac{e^{-y_i w \cdot x}}{1 + e^{-y_i w \cdot x}} \cdot (-y_i x_i), \dots, \frac{e^{-y_t w \cdot x}}{1 + e^{-y_t w \cdot x}} \cdot (-y_t x_i) \right)$$

$$w_{t+1} = w_t - \eta_t \cdot \nabla \ell / \log$$

:መሆኑን ተግባር

Average accuracy as a function of eta - SGD for log loss

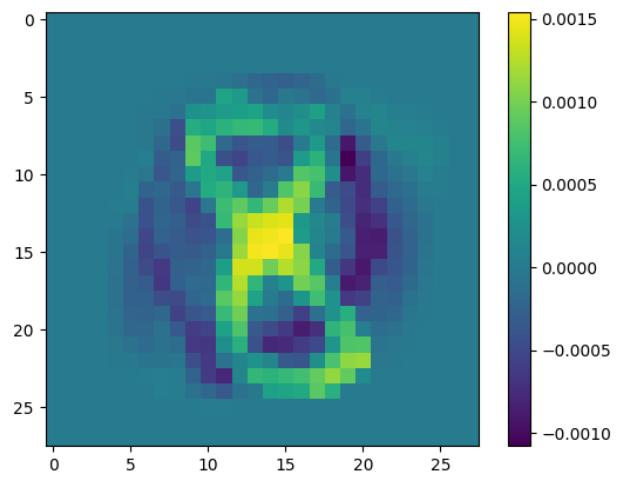


The best eta is: 1e-05

:መሆኑን የሚከተሉት አይነት

- (b) (5 points) Using the best  $\eta_0$  you found, train the classifier, but for  $T = 20000$ . Show the resulting  $\mathbf{w}$  as an image. What is the accuracy of the best classifier on the test set?

: מגדיר (ב)



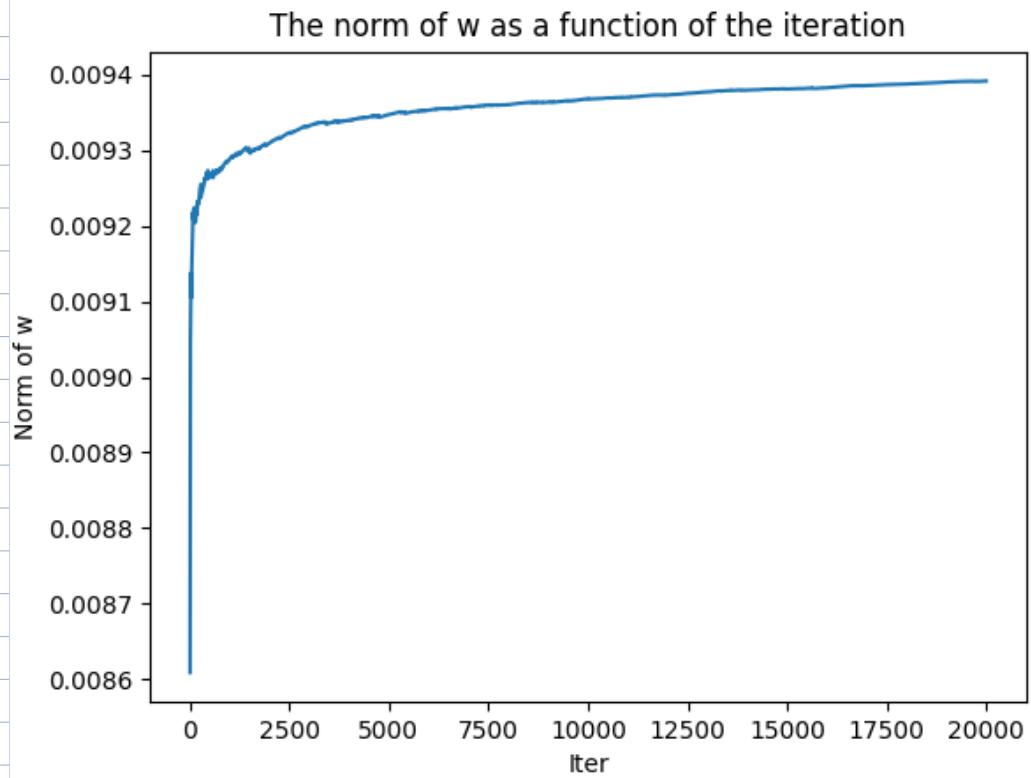
למגה ג נוֹטְרָלְנְדִי אֶת כְּבוֹד

: גַּיְשׁוּ כְּלָבִים

The accuracy of the best classifier is: 97.54350051177073 %

- (c) (5 points) Train the classifier for  $T = 20000$  iterations, and plot the norm of  $\mathbf{w}$  as a function of the iteration. How does the norm change as SGD progresses? Explain the phenomenon you observe.

(C)



מבחן בפועל נרמז, כי הגדלת תרומה כפולה נהייה  
יותר יעילה מאשר מילוי תרומה אחת. מילוי תרומה  
אחד פעמיים יתבצע בקצב מהיר יותר מאשר מילוי  
תרומה אחת פעמיים.