

1. (15 points) Step-size Perceptron.

בהנחה: $\|w^*\| = 1$ ו- w^* הוא מרכז עיגול S היחיד.

בהנחה: $\gamma > 0$ הוא קבוע, $w^* \in S$, $w_t \in S$, $\eta_t \in [0, 1]$.

$$y_t \cdot w^* \cdot x_t = \frac{|w^* \cdot x_t|}{\|w^*\|} = \text{dist}(x_t, \text{hyperplane}(w^*)) \geq \gamma$$

$$\|w^*\|$$

לפיכך γ הוא המרחק בין w^* ל- S .

$$w_{t+1} \cdot w^* = (w_t + \eta_t y_t x_t) \cdot w^* = w_t \cdot w^* + \eta_t y_t x_t \cdot w^* \geq$$

$$w_t \cdot w^* + \eta_t \gamma = w_t \cdot w^* + \frac{1}{\sqrt{t}} \gamma$$

אם, לפיכך w_{t+1} הוא מרכז עיגול S , נקבל:

$$w_{t+1} \cdot w^* \geq w_m \cdot w^* + \frac{1}{\sqrt{m}} \cdot \gamma \geq \dots \geq \gamma \sum_{i=1}^m \frac{1}{\sqrt{i}} \geq \gamma \int_1^{m+1} \frac{1}{\sqrt{x}} dx$$

$$= \gamma (2\sqrt{m+1} - 2) \geq \gamma \sqrt{em}$$

$$\lim_{m \rightarrow \infty} \frac{\sqrt{em}}{2\sqrt{m+1} - 2} = \frac{\sqrt{e}}{2} < 1$$

... - עדיין m מספיק קטן.

$$m \rightarrow \infty \quad 2\sqrt{m+1} - 2$$

בהנחה: לפיכך w_{t+1} הוא מרכז עיגול S , נקבל:

$$\|w_{t+1}\|^2 \leq \|w_t\|^2 + \frac{1}{t} \leq \dots \leq \sum_{i=1}^m \frac{1}{i} \leq 1 + \int_1^m \frac{1}{x} dx = 1 + \ln(m) = \ln(em)$$

$$\gamma \sqrt{em} \leq w_{t+1} \cdot w^* \leq \|w_{t+1}\| \cdot \|w^*\| = \|w_{t+1}\| \leq \sqrt{\ln(em)}$$

$$0 < \gamma \sqrt{em} \leq \sqrt{\ln(em)} \rightarrow \gamma^2 em \leq \ln(em) \rightarrow em \leq \frac{1}{\gamma^2} \ln(em)$$

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

... - עדיין m מספיק קטן.

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.

פתרון:

(3 נקודות) \mathbb{R}^n הוא קטע קטע.

יהי $x_1, x_2 \in \mathbb{R}^n$, ויהי $\lambda \in [0, 1]$. אז נקיים:

$$g(\lambda x_1 + (1-\lambda)x_2) = f(A(\lambda x_1 + (1-\lambda)x_2) + b) \leq \lambda f(Ax_1 + b) + (1-\lambda)f(Ax_2 + b) = \lambda g(x_1) + (1-\lambda)g(x_2)$$

נראה כי

- (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.)

פתרון: יהי $x_1, x_2 \in \mathbb{R}^d$, ויהי $\lambda \in [0, 1]$. אז נקיים:

$$g(\lambda x_1 + (1-\lambda)x_2) = \max_i f_i(\lambda x_1 + (1-\lambda)x_2) \leq \lambda \max_i f_i(x_1) + (1-\lambda) \max_i f_i(x_2)$$

$$= \lambda g(x_1) + (1-\lambda)g(x_2) \quad \square$$

(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 ℓ_{\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} .

פתרון:

נגלה כי ℓ_{\log} היא פונקציה קמורה.

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

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

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

לכן ℓ_{\log} היא פונקציה קמורה.

נראה כי f היא פונקציה קמורה. נניח $f : \mathbb{R}^d \rightarrow \mathbb{R}$ מתקיים $f(x) \geq 0$.

אנחנו יודעים כי f היא פונקציה קמורה. נניח $x \in \mathbb{R}^d$ ו- $y \in \mathbb{R}$.

נראה כי f היא פונקציה קמורה. נניח $x \in \mathbb{R}^d$ ו- $y \in \mathbb{R}$.

נניח $w \in \mathbb{R}^d$. מתקיים $w \cdot x = w \cdot (y \cdot x) = A w$ כאשר $A = (y \cdot x)^T$.

לכן $f(w) = \ell_{\log}(A w) = f(w)$ ו- f היא פונקציה קמורה.

3. (20 points) Ranking

- (a) Prove that the hinge loss described above for the ranking objective is convex in w .

(b) הכוונה היתה hinge loss

$$l(h_w(\bar{x}), y) = \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{i=j+1}^k \max\{0, 1 - \text{sgn}(y_j - y_i) w \cdot (x_j - x_i)\}$$

הפונקציה l קמורה ביחס ל- w .

בהינתן $w_1, w_2 \in \mathbb{R}^d$, $d \in \mathbb{N}$, $\alpha \in [0, 1]$, נגד:

$$l(h_{\alpha w_1 + (1-\alpha)w_2}(\bar{x}), y) = \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{i=j+1}^k \max\{0, 1 - \text{sgn}(y_j - y_i) (\alpha w_1 + (1-\alpha)w_2) \cdot (x_j - x_i)\}$$

$$= \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{i=j+1}^k \max\{0, \alpha(1 - \text{sgn}(y_j - y_i) w_1 \cdot (x_j - x_i)) + (1-\alpha)(1 - \text{sgn}(y_j - y_i) w_2 \cdot (x_j - x_i))\}$$

$$\leq \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{i=j+1}^k \max\{0, \alpha(1 - \text{sgn}(y_j - y_i) w_1 \cdot (x_j - x_i)) + (1-\alpha)(1 - \text{sgn}(y_j - y_i) w_2 \cdot (x_j - x_i))\}$$

$$= \alpha l(h_{w_1}(\bar{x}), y) + (1-\alpha) l(h_{w_2}(\bar{x}), y)$$

כלומר, הפונקציה l היא קמורה ביחס ל- w .

כלומר, l קמורה ביחס ל- w .

(b) Prove that the hinge loss upper-bounds the Kendall-Tau loss, i.e. that $\Delta(h_w(\bar{x}), y) \leq \ell(h_w(\bar{x}), y)$ for all $w \in \mathbb{R}^d$, $\bar{x} \in \mathcal{X}^k$, $y \in \mathbb{R}^k$.

הוכחה: Kendall-Tau

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

הי. $w \in \mathbb{R}^d$, $\bar{x} \in \mathcal{X}^k$, $y \in \mathbb{R}^k$.

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

$$= \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{r=j+1}^k \mathbb{1} \{ \text{sgn}(w \cdot x_j - w \cdot x_r) \neq \text{sgn}(y_j - y_r) \}$$

$$= \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{r=j+1}^k \mathbb{1} \{ \text{sgn}(w(x_j - x_r)) \neq \text{sgn}(y_j - y_r) \}$$

$$\leq \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{r=j+1}^k \mathbb{1} \{ \text{sgn}(w(x_j - x_r)) \text{sgn}(y_j - y_r) \in \{0, -1\} \}$$

$$\leq \frac{2}{k(k-1)} \sum_{j=1}^{k-1} \sum_{r=j+1}^k \max \{ 0, 1 - \text{sgn}(y_j - y_r) w \cdot (x_j - x_r) \} = \ell(h_w(\bar{x}), y)$$

• הנחנו גם במקרה זה את הנגזרת של ϕ - 0.

• אם $\gamma = -1$ $\text{sgn}(y_j - y_r) \text{sgn}(x_j - x_r) = 1$

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

$$\leq 1 \leq (x_j - x_r) \cdot \text{sgn}(y_j - y_r) \leq 0, \text{ אם } \gamma = 1.$$

במקרה השני - לא הכחנו את הסיוע γ : לקחת ϕ כפונקציה.

□

4. (15 points) **Gradient Descent on Smooth Functions.** We say that a continuously differentiable function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is β -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, β -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 β -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!)

בהנחה:

$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n \quad f$ הוא β -Smooth.

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

$$\mathbf{x}_t - \mathbf{x}_{t+\eta} = \eta \nabla f(\mathbf{x}_t) \quad \text{נבחר } \mathbf{y} = \mathbf{x}_{t+\eta}, \mathbf{x} = \mathbf{x}_t \text{ ב-} f$$

אולי ננסה $t \geq 0$

$$0 \leq \eta \|\nabla f(\mathbf{x}_t)\|^2 = \nabla f(\mathbf{x}_t) \cdot (\eta \nabla f(\mathbf{x}_t)) = \nabla f(\mathbf{x}_t) \cdot (\mathbf{x}_t - \mathbf{x}_{t+\eta})$$

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

$$\sum_{t=0}^T \eta \|\nabla f(\mathbf{x}_t)\|^2 \leq \frac{1}{\eta} \sum_{t=0}^T f(\mathbf{x}_t) - f(\mathbf{x}_{t+\eta}) + \frac{\beta}{2} \sum_{t=0}^T \|\mathbf{x}_t - \mathbf{x}_{t+\eta}\|^2 \quad \text{סכום}$$

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

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

$$f(\mathbf{x}_{T+\eta}) \geq 0 \quad \text{כי } f \geq 0, \quad 0 < 1 - \frac{\beta\eta}{2} \quad \text{כל } \eta < \frac{2}{\beta} \quad \text{אז } \leftarrow$$

$$\sum_{t=0}^T \eta \|\nabla f(\mathbf{x}_t)\|^2 \leq \frac{1}{1 - \frac{\beta\eta}{2}} \cdot \frac{1}{\eta} (f(\mathbf{x}_0) - f(\mathbf{x}_{T+\eta})) \leq \frac{1}{1 - \frac{\beta\eta}{2}} \cdot \frac{1}{\eta} f(\mathbf{x}_0) \quad \leftarrow$$

$$\sum_{t=0}^{\infty} \|\nabla f(x_t)\|^2 \leq \frac{1}{1-\frac{\mu\eta}{2}} \cdot \frac{1}{\eta} f(x_0) < \infty : f(x_0) < \infty \quad T \rightarrow \infty \quad \text{אזכור צורך}$$

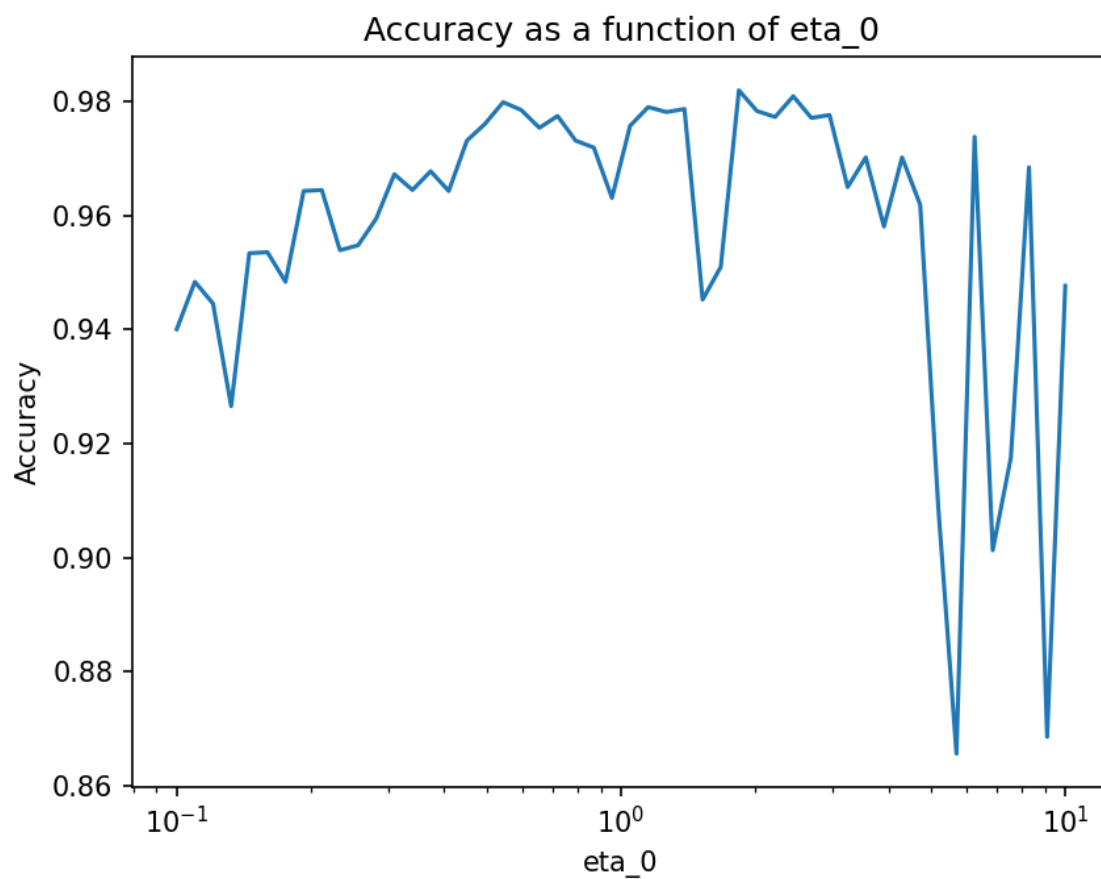
$$\lim_{t \rightarrow \infty} \|\nabla f(x_t)\| = 0 \quad \text{כי} \quad \lim_{t \rightarrow \infty} \|\nabla f(x_t)\|^2 = 0 \quad \Leftarrow$$

□ (נגרע).

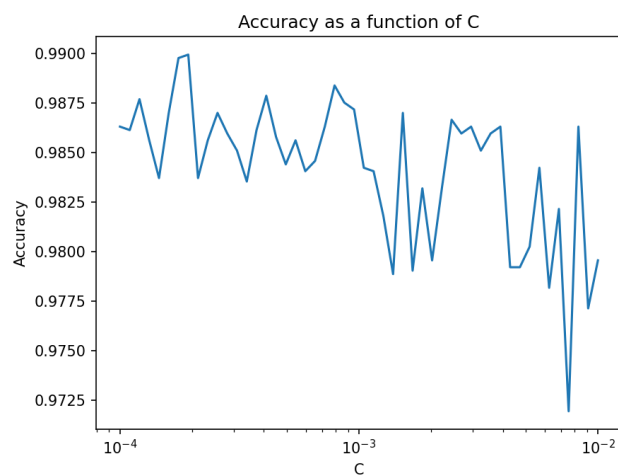
חלק תכנותי:

שאלה 1:

סעיף א:

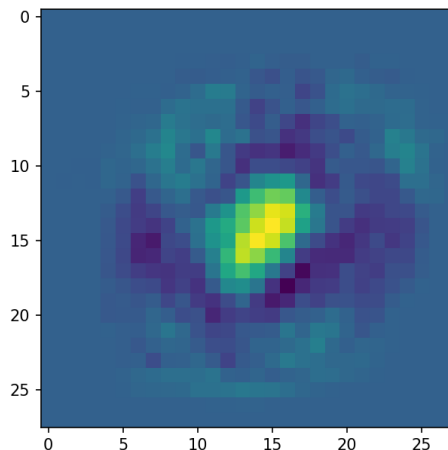


סעיף ב:



סעיף ג:

הסבר אינטואיטיבי – נשייך את הצבעים הבהירים לספרה 8, ואת הצבעים הכהים לספרה 0. הסימון הצהוב מבטא **דגש במרכז** (כלומר את התוכן שמופיע במרכז הספרה 8 לעומת הספרה 0). מנגד הסימון הכהה מבטא **דגש בצדדים** בצורה מעגלית, כלומר את התוכן של המעטפת של 0 לעומת המעטפת של הספרה 8.

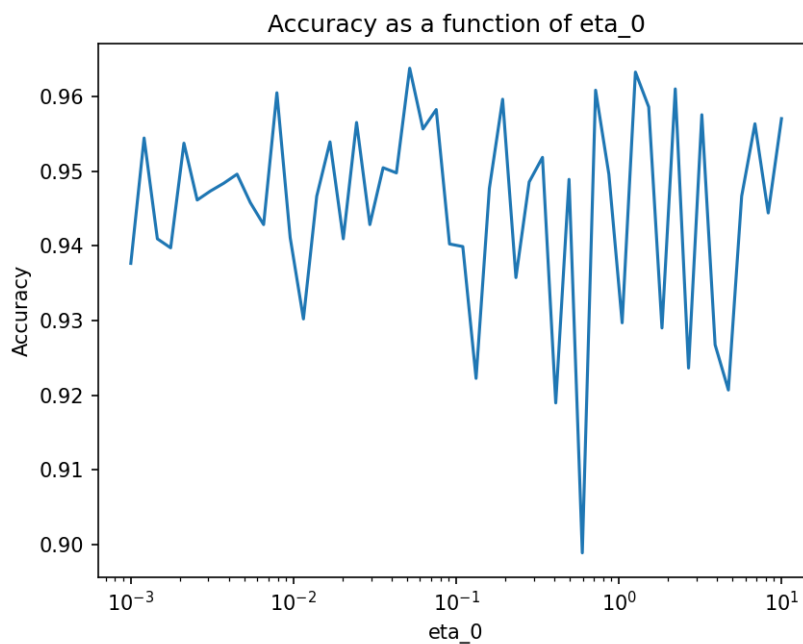


סעיף ד:

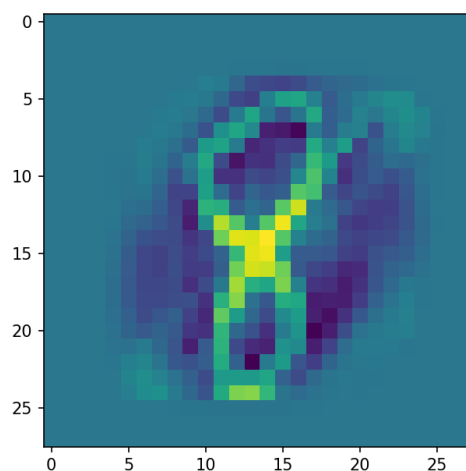
The accuracy of the best classifier on the test set is: 0.9923234390992836

שאלה 2:

סעיף א':

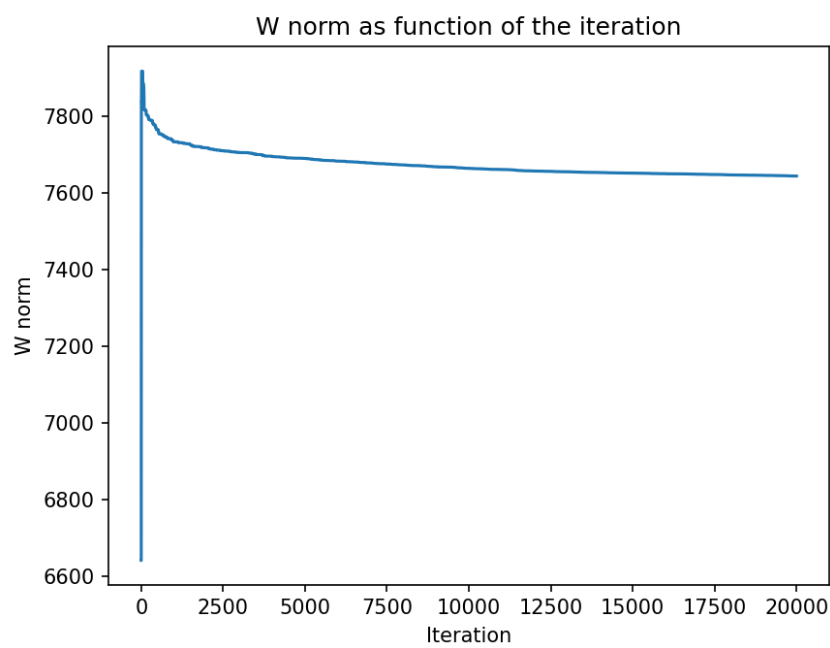


סעיף ב':



The accuracy of the best classifier on the test set is: 0.981064483111566

סעיף ג':



ככל ש-SGD מתקדם בתהליך – הנורמה קטנה (אמנם נשארת די יציבה). ההסבר לכך הוא שהאלגוריתם מגיע במהירות גבוהה יחסית למסווג טוב, ולאחר הגעה זאת ממשיך ומבצע רק שינויים קטנים (טיובים).