

# **Intro to Machine Learning: Assignment #4**

May 17, 2023

**Dor Bourshan**

# Theory Questions

## 1 SVM with multiple classes.

יהיו  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$  ויהיו  $y_1, \dots, y_n \in [K]$  labels מתוך קבוצה בת  $K$  labels.

נמצא מפריד לכל אחד מהמחלקות  $j \in [K]$ .

נגדיר את פונקציית loss הבאה (multiclass hinge – loss):

$$\ell(\mathbf{w}_1, \dots, \mathbf{w}_K, \mathbf{x}_i, y_i) = \max_{j \in [K]} (\langle \mathbf{w}_j, \mathbf{x}_i \rangle - \langle \mathbf{w}_{y_i}, \mathbf{x}_i \rangle + 1 \cdot \mathbf{1}(j \neq y_i))$$

נגדיר את בעיית ה multiclass SVM:

$$f(\mathbf{w}_1, \dots, \mathbf{w}_K) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}_1, \dots, \mathbf{w}_K, \mathbf{x}_i, y_i)$$

אחרי תהליך הלמידה אנחנו נסווג נקודות באופן הבא  $\arg \max_{j \in [K]} \langle \mathbf{w}_j, \mathbf{x}_i \rangle$ . נאמר כי אנחנו במקרה בו ה data מופרד. כלומר, קיימים  $\mathbf{w}_1^*, \dots, \mathbf{w}_K^*$  כך ש  $\langle \mathbf{w}_y^*, \mathbf{x}_i \rangle = \arg \max_y \langle \mathbf{w}_y^*, \mathbf{x}_i \rangle$  לכל  $i$ . נראה שכל  $\mathbf{w}_1, \dots, \mathbf{w}_K$  שממנמים את  $f(\mathbf{w}_1, \dots, \mathbf{w}_K)$  משיגים שגיאה אפס על הקלספיקציה.

אם כן, אנו נדרשים לבצע

$$\min_{\mathbf{w}_1, \dots, \mathbf{w}_K} f(\mathbf{w}_1, \dots, \mathbf{w}_K)$$

תחילה נפשט את הפונקציה להצגה אינפורמטיבית:

$$\begin{aligned} f(\mathbf{w}_1, \dots, \mathbf{w}_K) &= \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}_1, \dots, \mathbf{w}_K, \mathbf{x}_i, y_i) = \frac{1}{n} \sum_{i=1}^n \max_{j \in [K]} (\langle \mathbf{w}_j, \mathbf{x}_i \rangle - \langle \mathbf{w}_{y_i}, \mathbf{x}_i \rangle + 1 \cdot \mathbf{1}(j \neq y_i)) \\ &= \frac{1}{n} \sum_{i=1}^n \max \left\{ \max_{\substack{j \in [K] \\ j \neq y_i}} (\langle \mathbf{w}_j, \mathbf{x}_i \rangle - \langle \mathbf{w}_{y_i}, \mathbf{x}_i \rangle + 1), 0 \right\} \end{aligned}$$

כעת אנו יכולים לראות כי

$$f(\mathbf{w}_1, \dots, \mathbf{w}_K) \geq 0 \quad (1)$$

נראה כי

$$\ell(\mathbf{w}_1, \dots, \mathbf{w}_K, \mathbf{x}_i, y_i) \geq \ell_{1-0}(h(\mathbf{x}_i), y_i)$$

כלומר אנו נדרשים להראות כי:

$$\max_{j \in [K]} (\langle \mathbf{w}_j, \mathbf{x}_i \rangle - \langle \mathbf{w}_{y_i}, \mathbf{x}_i \rangle + 1 \cdot \mathbf{1}(j \neq y_i)) \geq 1 \cdot \mathbf{1}(\arg \max_{j \in [K]} \langle \mathbf{w}_j, \mathbf{x}_i \rangle \neq y_i)$$

נפריד למקרים:

• עבור  $\arg \max_{j \in [K]} \langle \mathbf{w}_j, \mathbf{x}_i \rangle = y_i$  אנו נקבל כי

$$\max_{j \in [K]} (\langle \mathbf{w}_j, \mathbf{x}_i \rangle - \langle \mathbf{w}_{y_i}, \mathbf{x}_i \rangle + 1 (j \neq y_i)) \stackrel{1}{\geq} 0$$

• אם  $\arg \max_{j \in [K]} \langle \mathbf{w}_j, \mathbf{x}_i \rangle \neq y_i$  מתקיים כי קיים  $j \neq y_i$  כך ש  $\arg \max_{j \in [K]} \langle \mathbf{w}_j, \mathbf{x}_i \rangle = j$  כלומר

$$\begin{aligned} \langle \mathbf{w}_j, \mathbf{x}_i \rangle &\geq \langle \mathbf{w}_{y_i}, \mathbf{x}_i \rangle \implies \langle \mathbf{w}_j, \mathbf{x}_i \rangle - \langle \mathbf{w}_{y_i}, \mathbf{x}_i \rangle \geq 0 \implies \langle \mathbf{w}_j, \mathbf{x}_i \rangle - \langle \mathbf{w}_{y_i}, \mathbf{x}_i \rangle + 1 \geq 1 \\ &\implies \max_{j \in [K]} (\langle \mathbf{w}_j, \mathbf{x}_i \rangle - \langle \mathbf{w}_{y_i}, \mathbf{x}_i \rangle + 1 (j \neq y_i)) \geq 1 \\ &\implies \max_{j \in [K]} (\langle \mathbf{w}_j, \mathbf{x}_i \rangle - \langle \mathbf{w}_{y_i}, \mathbf{x}_i \rangle + 1 (j \neq y_i)) \geq 1 \left( \arg \max_{j \in [K]} \langle \mathbf{w}_j, \mathbf{x}_i \rangle \neq y_i \right) \end{aligned}$$

כלומר הראנו כי

$$\ell(\mathbf{w}_1, \dots, \mathbf{w}_K, \mathbf{x}_i, y_i) \geq \ell_{1-0}(h(\mathbf{x}_i), y_i) \quad (2)$$

נראה כי קיימים  $\mathbf{w}_1^*, \dots, \mathbf{w}_K^*$  כך ש  $f(\mathbf{w}_1^*, \dots, \mathbf{w}_K^*) = 0$  ונקבל מ 1 כי ה"ל מנמנים את  $f$ .

נזכור כי ה data מופרד. כלומר, קיימים  $\mathbf{w}_1^*, \dots, \mathbf{w}_K^*$  כך ש  $y_i = \arg \max_y \langle \mathbf{w}_y^*, \mathbf{x}_i \rangle$  לכל  $i$ .

נסמן

$$\max_{1 \leq i \leq n} \max_{\substack{j \in [K] \\ j \neq y_i}} (\langle \mathbf{w}_j^*, \mathbf{x}_i \rangle - \langle \mathbf{w}_{y_i}^*, \mathbf{x}_i \rangle) = M$$

ונשים לב כי  $M$  מספר שלילי שכן ה data מופרד וכן  $y_i = \arg \max_y \langle \mathbf{w}_y^*, \mathbf{x}_i \rangle$  נסמן  $\mathbf{w}_i^{**} = -\frac{1}{M} \mathbf{w}_i^*$ . מתקיים כי:

$$\max_{\substack{j \in [K] \\ j \neq y_i}} (\langle \mathbf{w}_j^{**}, \mathbf{x}_i \rangle - \langle \mathbf{w}_{y_i}^{**}, \mathbf{x}_i \rangle) \leq -\frac{1}{M} M = -1$$

לפיכך

$$\max_{\substack{j \in [K] \\ j \neq y_i}} (\langle \mathbf{w}_j^{**}, \mathbf{x}_i \rangle - \langle \mathbf{w}_{y_i}^{**}, \mathbf{x}_i \rangle) + 1 \leq 0$$

ולכן

$$\max \left\{ \max_{\substack{j \in [K] \\ j \neq y_i}} (\langle \mathbf{w}_j^{**}, \mathbf{x}_i \rangle - \langle \mathbf{w}_{y_i}^{**}, \mathbf{x}_i \rangle + 1), 0 \right\} = 0$$

כלומר

$$f(\mathbf{w}_1^{**}, \dots, \mathbf{w}_K^{**}) = \frac{1}{n} \sum_{i=1}^n \max \left\{ \max_{\substack{j \in [K] \\ j \neq y_i}} (\langle \mathbf{w}_j^{**}, \mathbf{x}_i \rangle - \langle \mathbf{w}_{y_i}^{**}, \mathbf{x}_i \rangle + 1), 0 \right\} = \frac{1}{n} \sum_{i=1}^n 0 = 0$$

כלומר הראנו כי מנמום הפונקציה, משמעו כי קיימים  $\mathbf{w}_1, \dots, \mathbf{w}_K$  כך ש  $f(\mathbf{w}_1, \dots, \mathbf{w}_K) = 0$

אם כן יהיו  $\mathbf{w}_1, \dots, \mathbf{w}_K$  כך ש  $f(\mathbf{w}_1, \dots, \mathbf{w}_K) = 0$  מתקבל כי :

$$0 = f(\mathbf{w}_1, \dots, \mathbf{w}_K) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}_1, \dots, \mathbf{w}_K, \mathbf{x}_i, y_i) \stackrel{2}{\geq} \frac{1}{n} \sum_{i=1}^n \ell_{1-0}(h(\mathbf{x}_i), y_i) \geq 0$$

לפיכך

$$\frac{1}{n} \sum_{i=1}^n \ell_{1-0}(h(\mathbf{x}_i), y_i) = 0$$

כנדרש.

□

## 2 Expressivity of ReLU networks.

מגדירים את ReLU פונקציה אקטיבציה, באופן הבא :

$$h(x) = x^+ = \max\{0, x\}$$

נראה כי ניתן לממש את הפונקציה  $f(x_1, x_2) = \max\{x_1, x_2\}$  באמצעות one hidden layer המורכבת מפונקציית האקטיבציה ReLU. ניתן להניח כי אין פונקציית אקטיבציה לאחר השכבה האחרונה.

אם כן

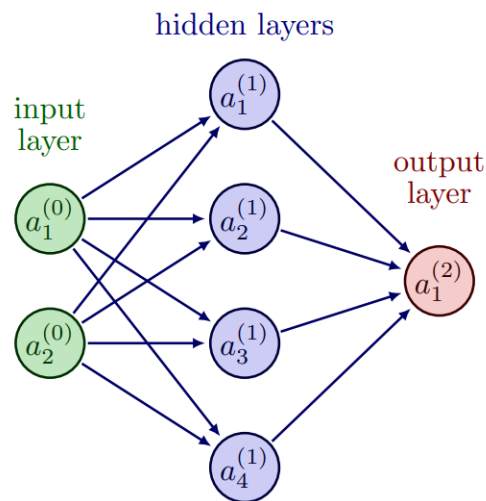


Figure 1: neural network

כאשר נגדיר :

$$\begin{aligned} a_1^{(1)} &= \frac{x_1 + x_2}{2} \\ a_2^{(1)} &= \frac{x_1 - x_2}{2} \\ a_3^{(1)} &= \frac{-x_1 + x_2}{2} \\ a_4^{(1)} &= \frac{-x_1 - x_2}{2} \end{aligned}$$

וכן

$$\begin{aligned} z_1^{(1)} &= \max\left\{\frac{x_1 + x_2}{2}, 0\right\} \\ z_2^{(1)} &= \max\left\{\frac{x_1 - x_2}{2}, 0\right\} \\ z_3^{(1)} &= \max\left\{\frac{-x_1 + x_2}{2}, 0\right\} \\ z_4^{(1)} &= \max\left\{\frac{-x_1 - x_2}{2}, 0\right\} \end{aligned}$$

וכן

$$a_1^{(2)} = \sum_{i=1}^3 z_i^{(1)} - z_4^{(1)}$$

נראה עבור חלוקה למקרים :

1. מתקיים כי  $x_1, x_2 \geq 0$ ולכן מתקיים כי  $z_1^{(1)} = \frac{x_1+x_2}{2}$  וכן  $z_4^{(1)} = 0$ (א) מתקיים כי  $x_1 \geq x_2$ מתקיים כי  $z_3^{(1)} = 0$  וכן  $z_2^{(1)} = \frac{x_1-x_2}{2}$  כמו כן  $z_2^{(1)} = \frac{x_1-x_2}{2} = \frac{|x_1-x_2|}{2}$ (ב) מתקיים כי  $x_1 < x_2$ מתקיים כי  $z_2^{(1)} = 0$  וכן  $z_3^{(1)} = \frac{-x_1+x_2}{2}$  כמו כן  $z_3^{(1)} = \frac{-x_1+x_2}{2} = \frac{|x_1-x_2|}{2}$ 

ניתן לראות אם כן כי

$$a_1^{(2)} = \sum_{i=1}^3 z_i^{(1)} - z_4^{(1)} = \frac{x_1+x_2}{2} + \frac{|x_1-x_2|}{2} = \max\{x_1, x_2\} = f(x_1, x_2)$$

2. מתקיים כי  $x_1, -x_2 \geq 0$ ולכן מתקיים כי  $z_2^{(1)} = \frac{x_1-x_2}{2} = \frac{|x_1-x_2|}{2}$  וכן  $z_3^{(1)} = 0$ (א) מתקיים כי  $x_1 \geq -x_2$ מתקיים כי  $z_1^{(1)} = \frac{x_1+x_2}{2}$  וכן  $z_4^{(1)} = 0$ (ב) מתקיים כי  $x_1 < -x_2$ מתקיים כי  $z_1^{(1)} = 0$  וכן  $z_4^{(1)} = \frac{-x_1-x_2}{2}$  ולכן  $z_4^{(1)} = \frac{x_1+x_2}{2}$ 

ניתן לראות אם כן כי

$$a_1^{(2)} = \sum_{i=1}^3 z_i^{(1)} - z_4^{(1)} = \frac{x_1+x_2}{2} + \frac{|x_1-x_2|}{2} = \max\{x_1, x_2\} = f(x_1, x_2)$$

3. מתקיים כי  $-x_1, -x_2 \geq 0$ ולכן מתקיים כי  $z_4^{(1)} = \frac{-x_1-x_2}{2} = \frac{x_1+x_2}{2}$  וכן  $z_1^{(1)} = 0$ (א) מתקיים כי  $-x_1 \geq -x_2$ מתקיים כי  $z_3^{(1)} = \frac{-x_1+x_2}{2} = \frac{|x_1-x_2|}{2}$  וכן  $z_2^{(1)} = 0$ (ב) מתקיים כי  $-x_1 < -x_2$ מתקיים כי  $z_3^{(1)} = 0$  וכן  $z_2^{(1)} = \frac{x_1-x_2}{2} = \frac{|x_1-x_2|}{2}$ 

ניתן לראות אם כן כי

$$a_1^{(2)} = \sum_{i=1}^3 z_i^{(1)} - z_4^{(1)} = \frac{x_1+x_2}{2} + \frac{|x_1-x_2|}{2} = \max\{x_1, x_2\} = f(x_1, x_2)$$

4. מתקיים כי  $-x_1, x_2 \geq 0$

ולכן מתקיים כי  $z_2^{(1)} = 0$  וכן  $z_3^{(1)} = \frac{-x_1+x_2}{2} = \frac{|x_1-x_2|}{2}$

(א) מתקיים כי  $-x_1 \geq x_2$

מתקיים כי  $z_4^{(1)} = \frac{-x_1-x_2}{2}$  וכן  $z_1^{(1)} = 0$  וכן  $-z_4^{(1)} = \frac{x_1+x_2}{2}$

(ב) מתקיים כי  $-x_1 < x_2$

מתקיים כי  $z_4^{(1)} = 0$  וכן  $z_1^{(1)} = \frac{x_1+x_2}{2}$

ניתן לראות אם כן כי

$$a_1^{(2)} = \sum_{i=1}^3 z_i^{(1)} - z_4^{(1)} = \frac{x_1+x_2}{2} + \frac{|x_1-x_2|}{2} = \max\{x_1, x_2\} = f(x_1, x_2)$$

סה"כ עברנו על כל המקרים וקיבלנו כי

$$a_1^{(2)} = \max\{x_1, x_2\} = f(x_1, x_2)$$

□

### 3 Soft SVM with $\ell^2$ penalty.

נניח שקיימת הבעיה הבאה :

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}, \xi \in \mathbb{R}^n} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + \frac{C}{2} \sum_{i=1}^n \xi_i^2 \\ \text{s.t.} \quad & y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i \quad \forall i = 1, \dots, n \end{aligned}$$

#### 3.1 section (a)

נראה תחילה כי אם נוסף את האילוץ  $\xi_i \geq 0$  לא נשנה את הפיתרון המינימלי.

אם כן תהיי  $I \subseteq [n]$  קבוצת אינדקסים כך שלכל  $i \in I$  מתקיים

$$\xi_i < 0$$

יהיי  $\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}, \xi \in \mathbb{R}^n$  פתרון אופטימלי. מתקיים אם כן כי

$$\begin{aligned} y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) &\geq 1 - \xi_i \geq 1 + \xi_i \quad \forall i \in I \\ y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) &\geq 1 - \xi_i \quad \forall i \in [n] \setminus I \end{aligned}$$

נשים לב כי בפונקציה המנומנמת, אין חשיבות לסימן של  $\xi_i$  ולכן עבור  $i \in I$  אם נחליף את  $\xi_i$  ב  $-\xi_i$  נקבל שהאילוץ מתקיימים וכן ערך הפונקציה המנומנמת זהה.

כלומר נוכל לעבור לבעיה הבאה :

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}, \xi \in \mathbb{R}^n} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + \frac{C}{2} \sum_{i=1}^n \xi_i^2 \\ \text{s.t.} \quad & y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i \quad \forall i = 1, \dots, n \\ & \xi_i \geq 0 \quad \forall i = 1, \dots, n \end{aligned}$$

□

#### 3.2 section (b)

נרשום את פונקציית הלגרנג'יאן של הבעיה :

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}, \xi \in \mathbb{R}^n} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + \frac{C}{2} \sum_{i=1}^n \xi_i^2 \\ \text{s.t.} \quad & 0 \geq 1 - \xi_i - y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \quad \forall i = 1, \dots, n \\ & 0 \geq -\xi_i \quad \forall i = 1, \dots, n \end{aligned}$$



אם כן

$$\mathcal{L}(\mathbf{w}, b, \xi, \alpha, \beta) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \frac{C}{2} \sum_{i=1}^n \xi_i^2 + \sum_{i=1}^n \alpha_i (1 - \xi_i - y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b)) + \sum_{i=1}^n \beta_i (-\xi_i)$$

כאשר אין בבעיה אילוצי שוויון.

□

### 3.3 section (c)

נמנמם את הפונקציה ביחס ל  $\mathbf{w}, b, \xi$ . זוהי תהיה הפונקציה הדואלית.

כלומר

$$g(\alpha, \beta) = \min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}, \xi \in \mathbb{R}^n} \mathcal{L}(\mathbf{w}, b, \xi, \alpha, \beta)$$

היא הפונקציה הדואלית.

ראשית נסדר את הלגרנגיאן.

$$\mathcal{L}(\mathbf{w}, b, \xi, \alpha, \beta) = \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_{i=1}^n \alpha_i y_i \langle \mathbf{w}, \mathbf{x}_i \rangle + \frac{C}{2} \sum_{i=1}^n \left( \xi_i^2 - \frac{2}{C} \xi_i (\alpha_i + \beta_i) \right) + \sum_{i=1}^n \alpha_i y_i b + \sum_{i=1}^n \alpha_i$$

נשים לב כי אם

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

אין לבעיה מינימום. כיוון שלכל ערך קטן כרצוננו ניתן לבחור  $b$ .

כעת ניתן לראות כי זו פונקציה של סכומים של פונקציות קמורות, ולכן קמורה (מתרגיל 3). והפונקציה דיפרנציאבילית כסכום של פונקציות דיפרנציאביליות.

אם קיים מינימום הוא נמצא בנקודה בה ההגרדיאנט מתאפס.

$$\nabla_{\mathbf{w}} \mathcal{L} = \mathbf{w} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i = 0 \implies \mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$$

$$\frac{\partial \mathcal{L}}{\partial b} = \sum_{i=1}^n \alpha_i y_i = 0$$

$$\frac{\partial \mathcal{L}}{\partial \xi_i} = C \cdot \xi_i - (\alpha_i + \beta_i) = 0 \implies \xi_i = \frac{\alpha_i + \beta_i}{C}$$

נציב ונקבל

$$\begin{aligned}
\mathcal{L}(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta}) &= \frac{1}{2} \|\mathbf{w}\|_2^2 - \underbrace{\mathbf{w}^\top \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i}_{\mathbf{w}} + \sum_{i=1}^n \alpha_i - \frac{1}{2C} \sum_{i=1}^n (\alpha_i + \beta_i)^2 + b \sum_{i=1}^n \alpha_i y_i \\
&= -\frac{1}{2} \|\mathbf{w}\|_2^2 - \frac{1}{2C} \sum_{i=1}^n (\alpha_i + \beta_i)^2 + \sum_{i=1}^n \alpha_i \\
&= \boxed{-\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i \mathbf{x}_j - \frac{1}{2C} \sum_{i=1}^n (\alpha_i + \beta_i)^2 + \sum_{i=1}^n \alpha_i}
\end{aligned}$$

### 3.4 section (d)

הפונקציה הדואלית:

$$g(\boldsymbol{\alpha}, \boldsymbol{\beta}) = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i \mathbf{x}_j - \frac{1}{2C} \sum_{i=1}^n (\alpha_i + \beta_i)^2 + \sum_{i=1}^n \alpha_i$$

הבעיה הדואלית:

$$\begin{aligned}
&\max g(\boldsymbol{\alpha}, \boldsymbol{\beta}) \\
&\sum_{i=1}^n \alpha_i y_i = 0 \\
&\alpha_i, \beta_i \geq 0 \quad \forall i = 1, \dots, n
\end{aligned}$$

□

## 4 Gradient of cross-entropy loss over softmax.

יהי  $\mathbf{y} \in \{0, 1\}^d$  וקטור בעל כניסה אחת של 1 והשאר 0. נגדיר את הפונקציה  $\ell_{\mathbf{y}} : \mathbb{R}^d \rightarrow \mathbb{R}$  הנקראת *cross-entropy loss* ע"י :

$$\ell_{\mathbf{y}}(\mathbf{w}) = -\mathbf{y} \ln(\text{softmax}(\mathbf{w}))$$

כאשר

$$\text{softmax}(\mathbf{w}) = \left( \frac{e^{w_1}}{\sum_{j=1}^d e^{w_j}}, \dots, \frac{e^{w_d}}{\sum_{j=1}^d e^{w_j}} \right)$$

וכן הלוגריתם מודר על וקטור כאיבר איבר.

נחשב את הדרדיאנט של  $\ell_{\mathbf{y}}(\mathbf{w})$  :

$$\begin{aligned} \ell_{\mathbf{y}}(\mathbf{w}) &= -\mathbf{y} \ln(\text{softmax}(\mathbf{w})) = -\sum_{i=1}^n y_i \ln \left( \frac{e^{w_i}}{\sum_{j=1}^d e^{w_j}} \right) = \sum_{i=1}^n \left( -y_i \ln(e^{w_i}) + y_i \ln \left( \sum_{j=1}^d e^{w_j} \right) \right) = \\ &= \sum_{i=1}^n -y_i w_i + \sum_{i=1}^n y_i \ln \left( \sum_{j=1}^d e^{w_j} \right) \end{aligned}$$

נגזור

$$\frac{\partial \ell}{\partial w_q} = -y_q + \sum_{i=1}^n y_i \frac{e^{w_q}}{\sum_{j=1}^d e^{w_j}}$$

נזכור כי יש רק כניסה אחת ל  $\mathbf{y}$  ששווה ל 1 ושאר הכניסות אפסים.

$$\frac{\partial \ell}{\partial w_q} = -y_q + \frac{e^{w_q}}{\sum_{j=1}^d e^{w_j}}$$

ולכן

$$\nabla \ell_{\mathbf{y}}(\mathbf{w}) = \text{softmax}(\mathbf{w}) - \mathbf{y}$$

□

# Programming Assignment

## Neural Networks

(a)

Code assignment

(b)

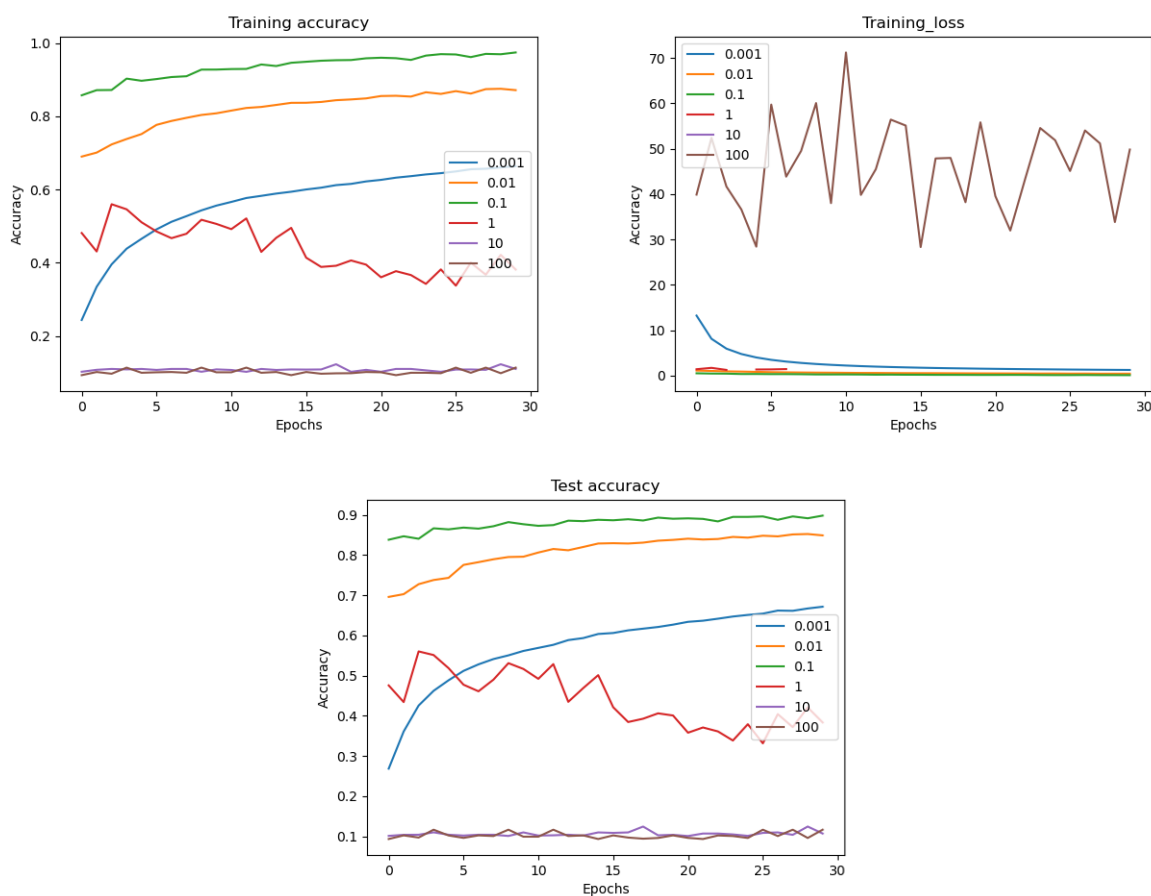


Figure 2: training accuracy, training loss and test accuracy

נוכח שה learning rate הוא גודל הצעד שאנו הולכים בכיוון הנגדי לגרדיאנט בנקודה.

אנו יכולים ראשית שהגרפים עבור Training accuracy and Test accuracy הם די זהים.

כאשר ה learning rate גדול מדי, אנחנו עתידים לא להגיע לנקודת מינימום לוקאלית (כלומר לרוב לא נגיע ממנה, ואם כן נצא ממנה במהירות). נצפה לראות את ה loss משתנה תכופות שכן אנו "שרירותית" משתנים את ערכי הרשת.

כאשר ה learning rate קטן מדי, אנחנו נוקטים בצעדים קטנים, נצפה אכן להגיע למינימום לוקאלי, אך ביותר איטרציות (וכתוצאה מכך בזמן). נצפה לראות את ה loss קטן אך לא בצורה דרסטית.

(c)

Test Accuracy : 0.9397
------------------------

## Training a deep Convolutional Neural Network.

(a)

	Loss	Accuracy on the validation	Accuracy on test	Learning rate
SGD	1.1383	48.7%	49.91%	0.005
Adam	0.6086	61.89%	62.49%	0.00005

(b)

	Loss	Accuracy on the validation	Accuracy on test	Learning rate
Adam	2.2486	13.44%	13.63%	0.005
Adam	0.6086	61.89%	62.49%	0.00005

(c)

	Loss	Accuracy on the validation	Accuracy on test	Learning rate	BatchNorm+Dropout
Adam	2.2486	60.62%	61.86%	0.00005	×
Adam	0.6086	61.89%	62.49%	0.00005	✓