

## Machine learning- Ex4, Osher Elhadad 318969748

### 1. (20 pts) Logistic regression

- (a) Prove that the logistic regression classifier is equivalent to a softmax over a linear multiclass classifier for two classes  $y = "a", y = "b"$ , where their separating hyperplanes obey  $w_a = -w_b$ .

#### תשובה-

נוכיח שקילות של *logistic regression classifier* ל-*softmax over linear multiclass classifier* ע"י הוכחת הטענה הבאה-

לכל קלט מתאים  $x$  עבור 2 classes "a" ו"b" אזי מתקיים כי קיים סט משקולות כלשהו

$w'$  שמייצג *logistic regression classifier* עבורו מתקיים  
 $\hat{y}_{\text{logistic regression classifier}} = "a"$  אמ"מ קיימים 2 סטים של משקולות כלשהם  $w_a, w_b$   
שמייצגים *softmax over linear multiclass classifier* כך ש-  
 $\hat{y}_{\text{softmax over linear multiclass classifier}} = "a"$  כאשר מתקיים  
their separating hyperplanes obey  $w_a = -w_b$

(מכיוון שאנו מבצעים *classification* עבור 2 classes לכן אם נוכיח טענה זו אזי זה גם יוכיח את הטענה השקולה עבור הערך "b" מכיוון שינם 2 ערכים- נתון *classifier binary*.)

#### כיוון ראשון:

יהי קלט  $x$  עבור 2 classes "a" ו"b" ועבור *logistic regression classifier* בהינתן תווית ליניארית  $y \in \{0,1\}$  (כאשר 1 מייצג את האות "a" ו0 את האות "b") וסט משקולות  $w'$  נניח כי  $\hat{y}_{\text{logistic regression classifier}} = "a"$  ונוכיח כי עבור 2 סטים של משקולות כלשהם  $w_a, w_b$  שמייצגים *softmax over linear multiclass classifier* מתקיים כי  $\hat{y}_{\text{softmax over linear multiclass classifier}} = "a"$  כאשר מתקיים *the classes separating hyperplanes obey*  $w_a = -w_b$   
נשים לב כי ב-*logistic regression classifier* בהינתן תווית ליניארית  $y \in \{0,1\}$  (כאשר 1 מייצג את האות "a" ו0 את האות "b") וסט משקולות  $w'$  אזי נקבל כי

$$p_{\text{logistic regression}}(y = 1|x) = \sigma(w'^T x) = \frac{1}{1 + e^{-w'^T x}} = \frac{e^{w'^T x}}{1 + e^{w'^T x}}$$

$$p_{\text{logistic regression}}(y = 0|x) = 1 - \sigma(w'^T x) = 1 - \frac{1}{1 + e^{-w'^T x}} = \frac{e^{-w'^T x}}{1 + e^{-w'^T x}}$$

ולכן למעשה ניתן להסתכל על הסתברות זו כפונקציה של  $w$ , נגדיר

$$f(w) = \frac{e^{w^T x}}{1 + e^{w^T x}}$$

ונקבל כי  $p_{\text{logistic regression}}(y = 0|x) = f(-w')$  וגם  $p_{\text{logistic regression}}(y = 1|x) = f(w')$  מכיוון שנתון  $\hat{y}_{\text{logistic regression classifier}} = "a"$  אזי לפי הגדרת *logistic regression classifier* מכיוון שה-*classifier* חזה כי התווית היא "a" לכן

למעשה ההסתברות לקבל אותו הייתה גדולה מההסתברות לקבל את "b", כלומר -  
 $p_{\text{logistic regression}}(y = 1|x) = f(w') > p_{\text{logistic regression}}(y = 0|x) = f(-w')$

$$\begin{aligned} f(w') > f(-w') &\rightarrow \frac{e^{w'^T x}}{1 + e^{w'^T x}} > \frac{e^{-w'^T x}}{1 + e^{-w'^T x}} \rightarrow e^{w'^T x}(1 + e^{-w'^T x}) \\ &> e^{-w'^T x}(1 + e^{w'^T x}) \\ &\rightarrow e^{w'^T x} - e^{-w'^T x} + e^{w'^T x}e^{-w'^T x} - e^{-w'^T x}e^{w'^T x} > 0 \\ &\rightarrow e^{w'^T x} - e^{-w'^T x} + e^0 - e^0 > 0 \rightarrow e^{w'^T x} > e^{-w'^T x} \\ &\rightarrow \frac{e^{w'^T x}}{e^{-w'^T x} + e^{w'^T x}} > \frac{e^{-w'^T x}}{e^{-w'^T x} + e^{w'^T x}} \end{aligned}$$

כאשר המעבר האחרון מוצדק מכיוון שכפלנו את 2 הצדדים בכולל חיובי  $\frac{1}{e^{-w'^T x} + e^{w'^T x}}$ .

וקיבלנו כי לפי הגדרת *softmax over linear multiclass classifier* נקבל כי מכיוון שישנם 2 classes - "a" ו"b" (בהתאמה לפי הסדר 1 ו 0) בעצם מכיוון שהראנו כי מתקיימים-  
 $p_{\text{logistic regression}}(y = 1|x) = f(w')$  וגם  $p_{\text{logistic regression}}(y = 0|x) = f(-w')$  למעשה עבור  $w'$  נקבל את ההסתברות של *logistic regression classifier* שיחזה את הערך "a" (ששקול 1 כפי שהגדרנו) ועבור  $-w'$  נקבל את ההסתברות של *logistic regression classifier* שיחזה את הערך "b" (ששקול 0 כפי שהגדרנו) בדיוק כמו שמוגדר *softmax over linear multiclass classifier* ונגדיר אותו כך- נסמן כי  $w'_a = w_a$  ו  $w'_b = -w'$  (אכן מתקיים הנתון בשאלה-  $w_a = -w_b$ ) שהם וקטורי המשקולות שלאחר *softmax* על התוצאה של מכפלת כל אחד מהם בנפרד ( $w_b^T$  ו  $w_a^T$ ) בקלט  $x$  נקבל למעשה כי  

$$p_{\text{Soft max}}(y = 0|x) = \frac{e^{-w'^T x}}{e^{-w'^T x} + e^{w'^T x}} p_{\text{Soft max}}(y = 1|x) = \frac{e^{w'^T x}}{e^{-w'^T x} + e^{w'^T x}}$$

$$p_{\text{Soft max}}(y = 1|x) > p_{\text{Soft max}}(y = 0|x) \text{ לכן } \frac{e^{w'^T x}}{e^{-w'^T x} + e^{w'^T x}} > \frac{e^{-w'^T x}}{e^{-w'^T x} + e^{w'^T x}}$$
ש *softmax over linear multiclass classifier* יחזה את הערך "a" (ששקול 1), כלומר  $\hat{y}_{\text{softmax over linear multiclass classifier}} = "a"$  כאשר מתקיים *the classes separating hyperplanes obey*  $w_a = -w_b$ .  
**כיוון שני:**

יהי קלט  $x$  עבור 2 classes "a" ו"b" ועבור *softmax over linear multiclass classifier* בהינתן תווית ליניארית  $y \in \{0,1\}$  (כאשר 1 מייצג את האות "a" ו 0 את האות "b") ו 2 סטים של משקולות  $w_a, w_b$  כך ש  $w_a = -w_b$  נניח כי  $\hat{y}_{\text{softmax over linear multiclass classifier}} = "a"$  ונוכיח כי עבור סט של משקולות  $w'$  כלשהו שמייצג *logistic regression classifier* מתקיים כי  $\hat{y}_{\text{logistic regression classifier}} = "a"$ .  
יהי קלט  $x$  עבור 2 classes "a" ו"b" נניח כי  $\hat{y}_{\text{softmax over linear multiclass classifier}} = "a"$  כאשר מתקיים *the classes separating hyperplanes obey*  $w_a = -w_b$  ונוכיח כי  $\hat{y}_{\text{logistic regression classifier}} = "a"$

נשים לב כי ב *softmax over linear multiclass classifier* בהינתן תווית ליניארית  $y \in \{0,1\}$  (כאשר 1 מייצג את האות "a" ו 0 את האות "b") וסט משקולות  $w'$  אזי נקבל כי

$$p_{\text{Soft max}}(y = 1|x) = \frac{e^{w_a^T x}}{e^{w_a^T x} + e^{w_b^T x}} = \frac{e^{w_a^T x}}{e^{w_a^T x} + e^{-w_a^T x}}$$

$$p_{Soft\ max}(y = 0|x) = \frac{e^{w_b^T x}}{e^{w_a^T x} + e^{w_b^T x}} = \frac{e^{-w_a^T x}}{e^{w_a^T x} + e^{-w_a^T x}}$$

לכן לפי הגדרת ה *softmax over linear multiclass classifier* מכיון שה *classifier* חזה כי התווית היא "a" לכן למעשה ההסתברות לקבל אותו הייתה גדולה מההסתברות לקבל את "b", כלומר  $p_{Soft\ max}(y = 1|x) > p_{Soft\ max}(y = 0|x)$  ולכן

$$\begin{aligned} \frac{e^{w_a^T x}}{e^{w_a^T x} + e^{-w_a^T x}} &> \frac{e^{-w_a^T x}}{e^{w_a^T x} + e^{-w_a^T x}} \xrightarrow{e^{w_a^T x} + e^{-w_a^T x} \geq 0} e^{w_a^T x} > e^{-w_a^T x} \\ &\rightarrow e^{w_a^T x} - e^{-w_a^T x} + e^0 - e^0 > 0 \\ &\rightarrow e^{w_a^T x} - e^{-w_a^T x} + e^{w_a^T x} e^{-w_a^T x} - e^{-w_a^T x} e^{w_a^T x} > 0 \\ &\rightarrow e^{w_a^T x} (1 + e^{-w_a^T x}) > e^{-w_a^T x} (1 + e^{w_a^T x}) \rightarrow \frac{e^{w_a^T x}}{1 + e^{w_a^T x}} > \frac{e^{-w_a^T x}}{1 + e^{-w_a^T x}} \end{aligned}$$

כעת נביט ב *logistic regression classifier* עם תווית ליניארית  $y \in \{0,1\}$  (כאשר 1 מייצג את האות "a" ו-0 את האות "b") וסט המשקולות  $w_a$  אזי נקבל כי

$$p_{logistic\ regression}(y = 1|x) = \sigma(w_a^T x) = \frac{1}{1 + e^{-w_a^T x}} = \frac{e^{w_a^T x}}{1 + e^{w_a^T x}}$$

$$p_{logistic\ regression}(y = 0|x) = 1 - \sigma(w_a^T x) = 1 - \frac{1}{1 + e^{-w_a^T x}} = \frac{e^{-w_a^T x}}{1 + e^{-w_a^T x}}$$

$$\text{ומכיון ש } \frac{e^{w_a^T x}}{1 + e^{w_a^T x}} > \frac{e^{-w_a^T x}}{1 + e^{-w_a^T x}} \text{ לכן}$$

$$p_{logistic\ regression}(y = 1|x) > p_{logistic\ regression}(y = 0|x)$$

ולכן ה *logistic regression classifier* יחזה את הערך "a" (ששקול ל1), כלומר  $\hat{y}_{logistic\ regression\ classifier} = "a"$  כדרוש.

- (b) Consider a softmax function for  $k$  classes that is scaled by a constant  $T \in \mathbb{R}$

$$f_i(\mathbf{x}) = \frac{\exp(\frac{1}{T} \mathbf{w}_i^T \mathbf{x})}{\sum_k \exp(\frac{1}{T} \mathbf{w}_k^T \mathbf{x})} \quad (1)$$

Show that when  $T \rightarrow 0$ , the softmax converges to the max function. What happens when  $T \rightarrow \infty$ ?

### תשובה-

נראה כי עבור  $T \in \mathbb{R}$  כך ש  $T \rightarrow 0$  אזי  $\text{softmax scaled by } T$  מתכנס ל  $\text{max function}$ , כלומר  $f_i(x)$  (כאשר  $1 \leq i \leq k$ ) מתכנס ל 1 עבור האינדקס  $i = j_1 \vee \dots \vee i = j_n$  כך ש  $1 \leq j_h \leq l$  לכל  $j_h \in \{j_1, \dots, j_n\}$  עבורו מתקיים כי  $w_{j_h}^T x > w_r^T x \forall 1 \leq r \leq k \wedge r \notin \{j_1, \dots, j_n\}$  כלומר לכל  $j_h \in \{j_1, \dots, j_n\}$  הערך  $w_{j_h}^T x$  מקבל ערך מקסימלי ( $w_{j_1}^T x = \dots = w_{j_n}^T x$ ), אחרת  $f_i(x)$  מתכנס ל 0 (עבור כל שאר האינדקסים  $\{1 \leq r \leq k: r \notin \{j_1, \dots, j_n\}\}$ ) ולכן למעשה נקבל שערכי  $\text{argmax}_i (w_i^T x)$  יקבלו ערך 1 ושאר האינדקסים עם ערך 0- כלומר פונקציית ה  $\text{max}$  שרצינו.

הפונקציה  $\text{softmax scaled by } T$  הינה מהצורה-

$$f_i(x) = \frac{e^{\frac{1}{T} w_i^T x}}{\sum_{k'} e^{\frac{1}{T} w_{k'}^T x}}$$

לכל  $i = j_h \in \{j_1, \dots, j_n\}$  נקבל:

$$\begin{aligned} \lim_{T \rightarrow 0} f_{j_h}(x) &= \lim_{T \rightarrow 0} \frac{e^{\frac{1}{T} w_{j_h}^T x}}{\sum_{k'} e^{\frac{1}{T} w_{k'}^T x}} = \lim_{T \rightarrow 0} \frac{e^{\frac{1}{T} w_{j_h}^T x}}{\sum_{k'} e^{\frac{1}{T} w_{k'}^T x}} \cdot \frac{e^{-\frac{1}{T} w_{j_h}^T x}}{e^{-\frac{1}{T} w_{j_h}^T x}} = \lim_{T \rightarrow 0} \frac{e^{\frac{1}{T} w_{j_h}^T x - \frac{1}{T} w_{j_h}^T x}}{\sum_{k'} e^{\frac{1}{T} w_{k'}^T x - \frac{1}{T} w_{j_h}^T x}} \\ &= \lim_{T \rightarrow 0} \frac{1}{e^{\frac{1}{T} w_{j_h}^T x - \frac{1}{T} w_{j_h}^T x} + \sum_{k' \neq j_h} e^{\frac{1}{T} w_{k'}^T x - \frac{1}{T} w_{j_h}^T x}} \\ &= \lim_{T \rightarrow 0} \frac{1}{1 + \sum_{k' \neq j} e^{\frac{1}{T} (w_{k'}^T x - w_{j_h}^T x)}} = 1 \end{aligned}$$

המעבר האחרון מוצדק לפי המעברים הבאים-

$$T \rightarrow 0 \Rightarrow \frac{1}{T} \rightarrow \infty$$

$$\forall k' \notin \{j_1, \dots, j_n\}: w_{j_h}^T x > w_{k'}^T x \Rightarrow w_{k'}^T x - w_{j_h}^T x < 0$$

ונקבל-

$$\frac{1}{T} (w_{k'}^T x - w_{j_h}^T x) \rightarrow -\infty \Rightarrow e^{\frac{1}{T} (w_{k'}^T x - w_{j_h}^T x)} \rightarrow 0 \Rightarrow \frac{1}{1 + \sum_{k' \neq j} e^{\frac{1}{T} (w_{k'}^T x - w_{j_h}^T x)}} \rightarrow 1$$

לכל  $i \notin \{j_1, \dots, j_n\}$  נקבל:

$$\begin{aligned}\lim_{T \rightarrow 0} f_i(x) &= \lim_{T \rightarrow 0} \frac{e^{\frac{1}{T} w_i^T x}}{\sum_{k'} e^{\frac{1}{T} w_{k'}^T x}} = \lim_{T \rightarrow 0} \frac{e^{\frac{1}{T} w_i^T x}}{\sum_{k'} e^{\frac{1}{T} w_{k'}^T x}} \cdot \frac{e^{-\frac{1}{T} w_i^T x}}{e^{-\frac{1}{T} w_i^T x}} = \lim_{T \rightarrow 0} \frac{e^{\frac{1}{T} w_i^T x - \frac{1}{T} w_i^T x}}{\sum_{k'} e^{\frac{1}{T} w_{k'}^T x - \frac{1}{T} w_i^T x}} \\ &= \lim_{T \rightarrow 0} \frac{1}{\sum_{k_1} e^{\frac{1}{T}(w_{k_1}^T x - w_i^T x)} + \sum_{k_2} e^{\frac{1}{T}(w_{k_2}^T x - w_i^T x)} + \sum_{k_3} e^{\frac{1}{T}(w_{k_3}^T x - w_i^T x)}}\end{aligned}$$

ונגדיר כל  $k_1 \in \{1, \dots, k\}$  כך שמתקיים  $w_{k_1}^T x - w_i^T x = 0$  בנוסף כל  $k_2 \in \{1, \dots, k\}$  שמתקיים  $w_{k_2}^T x - w_i^T x < 0$  וכל  $k_3 \in \{1, \dots, k\}$  כך שמתקיים  $w_{k_3}^T x - w_i^T x > 0$ . ונקבל-

$$\begin{aligned}& \lim_{T \rightarrow 0} \frac{1}{\sum_{k_1} e^{\frac{1}{T}(w_{k_1}^T x - w_i^T x)} + \sum_{k_2} e^{\frac{1}{T}(w_{k_2}^T x - w_i^T x)} + \sum_{k_3} e^{\frac{1}{T}(w_{k_3}^T x - w_i^T x)}} \\ &= \lim_{T \rightarrow 0} \frac{1}{\sum_{k_1} e^{\frac{1}{T} \cdot 0} + \sum_{k_2} e^{\frac{1}{T}(w_{k_2}^T x - w_i^T x)} + \sum_{k_3} e^{\frac{1}{T}(w_{k_3}^T x - w_i^T x)}} \left( = \frac{1}{1+0+\infty} \right) = 0\end{aligned}$$

המעבר האחרון מוצדק לפי המעברים הבאים-

$$T \rightarrow 0 \Rightarrow \frac{1}{T} \rightarrow \infty$$

$$\forall k_2 \in \{1, \dots, k\}: w_{k_2}^T x - w_i^T x < 0$$

ונקבל-

$$\frac{1}{T} (w_{k_2}^T x - w_i^T x) \rightarrow -\infty \Rightarrow e^{\frac{1}{T}(w_{k_2}^T x - w_i^T x)} \rightarrow 0$$

$$\forall k_3 \in \{1, \dots, k\}: w_{k_3}^T x - w_i^T x > 0$$

ונקבל-

$$\frac{1}{T} (w_{k_3}^T x - w_i^T x) \rightarrow \infty \Rightarrow e^{\frac{1}{T}(w_{k_3}^T x - w_i^T x)} \rightarrow \infty$$

כלומר  $f_i(x)$  (כאשר  $1 \leq i \leq k$ ) מתכנס ל-1 עבור האינדקס  $i = j_1 \vee \dots \vee i = j_n$  כך ש  $1 \leq j_h \leq k$  לכל  $j_h \in \{j_1, \dots, j_n\}$  עבורו מתקיים כי  $w_{j_h}^T x > w_r^T x \forall 1 \leq r \leq k \wedge r \notin \{j_1, \dots, j_n\}$  כלומר לכל  $j_h \in \{j_1, \dots, j_n\}$  הערך  $w_{j_h}^T x$  מקבל ערך מקסימלי ( $w_{j_1}^T x = \dots = w_{j_n}^T x$ ), אחרת  $f_i(x)$  מתכנס ל-0 (עבור כל שאר האינדקסים  $\{1 \leq r \leq k: r \notin \{j_1, \dots, j_n\}\}$  כדרוש).

כמו כן עבור  $T \rightarrow \infty$  נקבל:

$$\lim_{T \rightarrow 0} f_i(x) = \lim_{T \rightarrow 0} \frac{e^{\frac{1}{T} w_i^T x}}{\sum_k e^{\frac{1}{T} w_k^T x}} = \frac{e^0}{\sum_k e^0} = \frac{1}{k}$$

ונקבל במקרה זה כי  $f_i(x)$  מתכנסת ל- $\frac{1}{k}$  לכל  $i$  (כלומר למספר קבוע  $k$ - מספר  $classes$ ). כלומר זו תהיה פונקציה קבועה לכל  $i$  ו  $x$  שנקבל- מאבדים פה את כל ההבדלים בין כל  $w_i$  שאנו חיפשו מלכתחילה ב- $softmax$  כדי לקבל הסתברויות או מקסימום כמו שהצגנו פה למעלה.

- (c) Write the gradient update rule for a logistic regression model, when the usual loss of the negative log likelihood is now regularized with the square of the  $L_2$  norm over the weight vector  $\frac{1}{2}||w||^2$ .

### תשובה-

נראה את חוק העדכון של מודל ה *logistic regression* כאשר פונקציית ה *loss* הינה *negative log likelihood* בתוספת רגולריזציה של *square of L2 norm over weight vector*  $\frac{1}{2}||w||^2$

כלומר נקבל כי פונקציית ה *loss* היא  $Err_{new}(w) = -\log(\text{likelihood}) + \frac{1}{2}||w||^2$

$$Err_{new}(w) = -\log(\text{likelihood}) + \frac{1}{2}||w||^2$$

נזכיר מההרצאה כי-

$$\begin{aligned} \text{likelihood} &= \prod_{y_i=1} p(y=1|x_i) \prod_{y_i=0} p(y=0|x_i) \quad \forall z \in \mathbb{R} | z^0 = 1 \\ &= \prod_{i \in \{1, \dots, n\}} p(y=1|x_i)^{y_i} \prod_{i \in \{1, \dots, n\}} p(y=0|x_i)^{1-y_i} \\ p(y=1|x_i) &= \sigma(w^T x_i) \wedge p(y=0|x_i) = 1 - \sigma(w^T x_i) \\ &= \\ &= \prod_{i \in \{1, \dots, n\}} \sigma(w^T x_i)^{y_i} \prod_{i \in \{1, \dots, n\}} (1 - \sigma(w^T x_i))^{1-y_i} \end{aligned}$$

ולכן נקבל-

$$\begin{aligned} -\log(\text{likelihood}) &= \\ -\sum_{i \in \{1, \dots, n\}} y_i \log(\sigma(w^T x_i)) - \sum_{i \in \{1, \dots, n\}} (1 - y_i) \log(1 - \sigma(w^T x_i)) \end{aligned}$$

כעת כדי לחשב את הגרדיאנט נשתמש בכללים אלו-

$$\frac{d\sigma(z)}{dz} = \sigma(z)(1 - \sigma(z))$$

ולפי כלל השרשרת-

$$\frac{d \log(\text{likelihood})}{dw} = \frac{d \log(\sigma(w^T x_i))}{d\sigma(w^T x_i)} \frac{d\sigma(w^T x_i)}{dw^T x_i} \frac{dw^T x_i}{dw}$$

בנוסף כפי שראינו בהרצאה

$$\frac{d \frac{1}{2}||w||^2}{dw} = w$$

ולכן כעת נמצא את הגרדיאנט של  $Err_{new}$  נסמן את  $\sigma(w^T x_i)$  בקיצור  $\sigma_1$ -

$$\begin{aligned}
\nabla Err_{new}(w) &= w - \sum_{i \in \{1, \dots, n\}} y_i \cdot \frac{1}{\sigma_1} \cdot \sigma_1 \cdot (1 - \sigma_1) \cdot x_i + (1 - y_i) \cdot \frac{-1}{1 - \sigma_1} \cdot \sigma_1 \cdot (1 - \sigma_1) \cdot x_i \\
&= w - \sum_{i \in \{1, \dots, n\}} y_i \cdot (1 - \sigma_1) \cdot x_i - (1 - y_i) \cdot \sigma_1 \cdot x_i \\
&= w - \sum_{i \in \{1, \dots, n\}} x_i (y_i \cdot (1 - \sigma_1) - (1 - y_i) \cdot \sigma_1) \\
&= w - \sum_{i \in \{1, \dots, n\}} x_i (y_i - y_i \cdot \sigma_1 - \sigma_1 + y_i \cdot \sigma_1) = w - \sum_{i \in \{1, \dots, n\}} x_i (y_i - \sigma(w^T x_i))
\end{aligned}$$

$$\nabla Err_{new}(w) = \sum_{i \in \{1, \dots, n\}} \frac{1}{n} w - x_i (y_i - \hat{y}_i)$$

כלומר הגרדיאנט עבור דגימה  $i$  הוא  $\frac{1}{n} w - x_i (y_i - \hat{y}_i)$

כעת נגדיר את כלל העדכון של  $w$  (בדומה להרצאה לפי הדגימה  $i$ )-

$$w^{t+1} = w^t - \eta \left( \frac{1}{n} w^t - x_i (y_i - \hat{y}_i) \right) = w^t - \frac{\eta}{n} w^t + \eta x_i (y_i - \hat{y}_i)$$

וכלל העדכון של  $w$  לפי  $n$  דגימות יהיה-

$$w^{t+1} = w^t - \eta \left( w^t - \sum_{i \in \{1, \dots, n\}} x_i (y_i - \hat{y}_i) \right) = w^t - \eta w^t + \eta \sum_{i \in \{1, \dots, n\}} x_i (y_i - \hat{y}_i)$$

## 2. (10 pts) Support Vector Machines

At class we discussed the SVM formulation

$$\min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i \quad (2)$$

where  $C$  is a trade-off hyper parameter that weighs the cost of misclassification (through the slack variable). You are given a binary classification problem for which the cost of misclassifying positive samples is different than the cost of misclassifying negative examples. Formulate the problem as an SVM with two types of slack variables. How many variables are now optimized over in your optimization problem, and what are their dimensions?

### תשובה-

נמדל את בעיית האופטימיזציה ש-SVM פותר כאשר אנחנו מתבוננים בשני סוגי slack variables ( $\xi_i$ ). הראשון הינו עבור דגימות חיוביות שמסווגות לא נכון (כשליליות) או שמפספסות את הmargin (כפי שהוצג בהרצאה). והשני הינו עבור דגימות שליליות שמסווגות לא נכון (כחיוביות) או שמפספסות את הmargin (כפי שהוצג בהרצאה).

$$\min_{\mathbf{w} \in \mathbb{R}^d, \xi_i} \frac{1}{2} \|\mathbf{w}\|^2 + C \cdot \sum_{i=1}^n \xi_i$$

כאשר  $\forall i: \xi_i \geq 0$  וגם  $w^T x_i y_i \geq 1 - \xi_i$ .

למעשה הווקטורים  $w$  שיבחרו יקיימו שהמרחק לsupport vectors הוא 1 ולכן למעשה החלק השמאלי של הרגולציה בעצם דואג לכך שהנורמה של הווקטור  $w$  תהיה קטנה, עם ערכים קטנים  $w$  ושעדיין המרחק מהsupport vectors הוא 1. כמו כן המשתנה  $\xi_i$  עבור כל דגימה, והוגדר בהרצאה שמתקיים  $\xi_i = \max(0, 1 - w^T x_i y_i)$  כלומר אם הדגימה אכן סווגה נכונה וגם נמצאת מעל הטווח שהגדרנו (מעבר לsupport vector שבצד של הסיווג הנכון לדגימה זו) אזי  $\xi_i$  יהיה למעשה 0- כלומר לא יהיה עונש, אחרת יקבע ע"י המרחק מהsupport vector (שבצד של הסיווג הנכון לדגימה זו) כעונש לנקודה זו שלא נמצאת בטווח שרצינו. כלומר נקבל hinge loss בחלק הימני של הביטוי. כעת נשים לב כי  $C$  הוא למעשה היפר פרמטר שכאשר הוא גדול יותר הוא נותן יותר חשיבות לעונש שהגדרנו לנקודות שלא בטווח הנכון (מעבר לsupport vector בכיוון הסיווג הנכון של הדגימה) ופחות לנורמה של  $w$ , כלומר יאפשר  $w$  יותר גדולים. לעומת זאת כאשר  $C$  קטן יותר נקבל כי אנחנו נותנים בעיקר חשיבות לנורמה של  $w$  ופחות לטעויות של  $\xi_i$ . כעת נשים לב כי אנו משנים את הבעיה במעט- העונש על הטעות יהיה שונה לדגימות שליליות וחיוביות, כלומר נכפיל בערך  $C$  שונה לדגימות החיוביות והשליליות. נחלק את הדגימות כך-

לכל דגימה  $x_i^+$  חיובית נסמן את הslack variable המתאים עם  $\xi_i^+$  ואת ההיפר פרמטר  $C$  עם  $C^+$ , ולכל דגימה  $x_i^-$  שלילית נסמן את ה-slack variable המתאים עם  $\xi_i^-$  ואת ההיפר פרמטר  $C$  עם  $C^-$  (כך ש-  $C^+ \neq C^-$ ).

כעת נגדיר את הבעיה המתאימה לתשובה כך-

$$\min_{\mathbf{w} \in \mathbb{R}^d, \xi_i^+, \xi_i^-} \frac{1}{2} \|\mathbf{w}\|^2 + C^+ \cdot \sum_{x_i^-} \xi_i^+ + C^- \cdot \sum_{x_i^+} \xi_i^-$$

כאשר  $\forall i: \xi_i^- \geq 0$ , בנוסף  $\forall i: \xi_i^+ \geq 0$  וגם  $w^T x_i y_i \geq 1 - \xi_i^+$



בבעיה המקורית בהרצאה נאמר כי עקב הוספת  $\xi$  הפרמטרים שאנו לומדים הפכו להיות  $w, \xi$  אך מכיוון ש  $\xi_i = \max(0, 1 - w^T x_i y_i)$  לכן למעשה אנו לומדים את  $w$  ויש לנו אילוצים על  $\xi$  (ויש לנו  $n$  משתנים כאלו- כלומר מימד  $n$ ) ומתוך זה גם נקבל את  $\xi_i$  (שהוא תלוי ב  $w$ ) שאנו מחפשים (את  $C$  אנו לא לומדים- היפר פרמטר). זאת בדיוק כמו בבעיה החדשה שעליה אנו מסתכלים שגם בה אנו מאפסמים את  $w$  שמימדו הוא  $n$  ויש לנו אילוצים על  $\xi$  (ויש לנו  $n$  משתנים כאלו- כלומר מימד  $n$ ) ומתוך כך גם נקבל את  $\xi_i$  (שהוא תלוי ב  $w$ ) שאנו מחפשים (גם פה אנו לא לומדים את ההיפר פרמטרים  $C^-, C^+$ ).

### 3. (10 pts) Polynomial kernels

Given a pair of samples in input space  $\mathbf{x}, \mathbf{z} \in R^2$ , show that the polynomial kernel of degree 3,  $K(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x}^T \mathbf{z})^3$  can be expressed as a dot product in feature space  $\phi(\mathbf{x})^T \phi(\mathbf{z})$  with the right choice of a feature map  $\phi : R^2 \rightarrow R^d$  ( $d > 2$ ). Write the explicit form of  $\phi$ .

How many operations are needed to compute  $K(\mathbf{x}, \mathbf{z})$ , and how many are needed for computing  $\phi(\mathbf{x})\phi(\mathbf{z})$ ? How does the number of operation changes for polynomial kernel of higher degrees?

**תשובה-**

$$z = \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}, x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \text{ - נסמן}$$

למעשה נמצא את הפונקציה  $\phi$  ע"י פתיחה של הביטוי  $(1 + x^T z)^3$  עד הסוף ולאחר מכן בעצם כל מחובר שיהיה בביטוי הסופי יהיה בעצם מכפלה של 2 ביטויים בוקטורים  $\phi(x), \phi(z)$  באותה השורה (כאשר אלו וקטורי עמודה) וכך למעשה נמצא כל ביטוי בכל שורה בוקטור שהפונקציה  $\phi$  מחזירה.

$$1 + 3x_1z_1 + 3x_2z_2 + 3x_1^2z_1^2 + 6x_1z_1x_2z_2 + 3x_2^2z_2^2 + x_1^3z_1^3 + 3x_1^2z_1^2x_2z_2 + 3x_1z_1x_2^2z_2^2 + x_2^3z_2^3 =$$

$$1 + 3(x_1z_1 + x_2z_2) + 3(x_1z_1 + x_2z_2)^2 + (x_1z_1 + x_2z_2)^3 =$$

$$1 + 3x^T z + 3(x^T z)^2 + (x^T z)^3 =$$

$$(1 + x^T z)^3$$

$$\phi(x)^T \cdot \phi(z) =$$

$$\begin{pmatrix} x_2^3 & \sqrt{3} \cdot x_1x_2^2 & \sqrt{3} \cdot x_1^2x_2 & x_1^3 & \sqrt{3} \cdot x_2^2 & \sqrt{6} \cdot x_1x_2 & \sqrt{3} \cdot x_1^2 & \sqrt{3} \cdot x_2 & \sqrt{3} \cdot x_1 & 1 \end{pmatrix} \cdot \begin{pmatrix} z_2^3 \\ \sqrt{3} \cdot z_1z_2^2 \\ \sqrt{3} \cdot z_1^2z_2 \\ z_1^3 \\ \sqrt{3} \cdot z_2^2 \\ \sqrt{6} \cdot z_1z_2 \\ \sqrt{3} \cdot z_1^2 \\ \sqrt{3} \cdot z_2 \\ \sqrt{3} \cdot z_1 \\ 1 \end{pmatrix} =$$

$$= 1 + 3x_1z_1 + 3x_2z_2 + 3x_1^2z_1^2 + 6x_1z_1x_2z_2 + 3x_2^2z_2^2 + x_1^3z_1^3 + 3x_1^2z_1^2x_2z_2 + 3x_1z_1x_2^2z_2^2 + x_2^3z_2^3$$

$$\phi(v) = \phi \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} = \begin{pmatrix} v_2^3 \\ \sqrt{3} \cdot v_1 v_2^2 \\ \sqrt{3} \cdot v_1^2 v_2 \\ v_1^3 \\ \sqrt{3} \cdot v_2^2 \\ \sqrt{6} \cdot v_1 v_2 \\ \sqrt{3} \cdot v_1^2 \\ \sqrt{3} \cdot v_2 \\ \sqrt{3} \cdot v_1 \\ 1 \end{pmatrix} \quad \text{לכן, נגדיר את הפונקציה } \phi: \mathbb{R}^2 \rightarrow \mathbb{R}^{10} \text{ לכל } v, \text{ מתקיים}$$

כעת נגדיר פעולה אחת ככפל סקלרי או חיבור סקלרי, ולכן חזקת  $i$  צורכת  $i$  פעולות.

נחשב כמה פעולות חישוב אנו צריכים לבצע על מנת לחשב את  $\phi(x)^T \cdot \phi(z)$  ישירות-

תחילה נשים לב שעל מנת להגיע מכפל הווקטורים לביטוי הבא-

$$1 + 3x_1z_1 + 3x_2z_2 + 3x_1^2z_1^2 + 6x_1z_1x_2z_2 + 3x_2^2z_2^2 + x_1^3z_1^3 + 3x_1^2z_1^2x_2z_2 + 3x_1z_1x_2^2z_2^2 + x_2^3z_2^3$$

למעשה עלינו לבצע הכפלה של כל המספרים  $1, \sqrt{3}, \sqrt{6}$  בין האיברים המתאימים בכפל הווקטורי כלומר יש לנו 8 הכפלות סקלריות (8 פעולות כפי שהגדרנו) כאלו שיביאו לנו את האיברים הבאים בתוך הביטוי הגדול שלנו (שביניהם יהיה חיבור בתוספת האיברים שלא היה להם מספרים ממשיים שכפלנו בהם)-

$$1, 3x_1z_1, 3x_2z_2, 3x_1^2z_1^2, 6x_1z_1x_2z_2, 3x_2^2z_2^2, 3x_1^2z_1^2x_2z_2, 3x_1z_1x_2^2z_2^2$$

$3x_1z_1$  – צורך 2 פעולות כפל,  $3x_2z_2$  – צורך 2 פעולות כפל,  $3x_1^2z_1^2$  – צורך 4 פעולות כפל,  $6x_1z_1x_2z_2$  – צורך 5 פעולות כפל,  $3x_2^2z_2^2$  – צורך 4 פעולות כפל,  $3x_1^2z_1^2x_2z_2$  – צורך 6 פעולות כפל,  $3x_1z_1x_2^2z_2^2$  – צורך 6 פעולות כפל,  $x_1^3z_1^3$  – צורך 5 פעולות כפל,  $x_2^3z_2^3$  – צורך 5 פעולות כפל  
סה"כ הגענו ל-10 מחוברים, כלומר 9 פעולות חיבור. נסכום הכל:

$$2 + 2 + 4 + 4 + 4 + 5 + 6 + 6 + 5 + 9 + 8 = 55$$

נחשב כמה פעולות חישוב אנו צריכים לבצע על מנת לחשב את הקרנל-

חישוב  $x^T z$  יעלה לנו 2 פעולות כפל ופעולת חיבור לכפל הווקטורי בין 2 ווקטורים ממימד 2. כמו כן הוספת 1 (עוד פעולה אחת), והעלאה בשלישית מצריכה 2 פעולות כפל. סה"כ ביצענו 6 פעולות.

נחשב כמה פעולות חישוב אנו צריכים לבצע אם  $d$  היה גדול יותר-

אם המימד  $d$  היה גדול יותר, הקרנל היה מהצורה  $(1 + x^T z)^d$ , ולכן בהנחה ש- $x, z \in \mathbb{R}^2$ , בהתאמה היינו נדרשים לבצע  $d + 4$  פעולות. כלומר כמות הפעולות שהקרנל מצריך היא לינארית למימד של  $\phi$ , בניגוד לחישוב ישיר שעולה לנו הרבה יותר פעולות כפי שראינו עבור  $d = 3$  (מכיוון שנקבל יותר מחוברים בסוף ההעלאה בחזקה כפי שעשינו בתחילת התרגיל).