



מבוא למערכות לומדות (236756)

סמסטר אביב תשפ"ד – 3 בספטמבר 2024

מרצה: ד"ר ניר רוזנפלד

## מבחן מסכם מועד א'

### הנחיות הבחינה:

- **משך הבחינה:** שלוש שעות.
- **חומר עזר:** המבחן בחומר סגור (ללא ספרים, מחברות, דפי נוסחאות).
- מחשבון: מותר.
- כלי כתיבה: עט בלבד.
- יש לכתוב את התשובות **על גבי שאלון זה**.
- מותר לענות בעברית או באנגלית.
- הוכחות והפרכות צריכות להיות פורמליות.
- קריאות:
  - סימונים לא ברורים בשאלות רב-ברירה ו/או תשובות מילוליות בכתב יד לא קריא יובילו לפסילת התשובה.
  - לא יתקבלו ערעורים בנושא.
- במבחן 19 עמודים ממוספרים סה"כ, כולל עמוד שער זה שמספרו 1 ושלושה עמודי טיוטה בסוף הגיליון.
- נא לכתוב רק את המבוקש ולצרף הסברים קצרים עפ"י ההנחיות.
- **בתום המבחן יש להגיש את שאלון זה בלבד.**
- **לזכאים להערכה חלופית מתאפשרת בחירה בין שאלות 3 ו-4.**
- **זכרו: Less is more.** אל תכתבו פרטים מיותרים.

**בהצלחה!**

## שאלה 1: SVM, Imbalanced sampling [35 נק']

בבעיה זאת נעסוק בבעיית סיווג בינארי  $\mathcal{Y} = \{-1, 1\}$ , עם מרחב דוגמאות  $\mathcal{X} = \mathbb{R}^d$ .

א. [15 נק'] בסעיף זה  $d = 2$ . דגמו באקראי  $S \subset \mathcal{X}$ .

מימין נתונים 4 איורים המראים decision boundaries שונים. באזורים אפורים המודל מנבא +1.

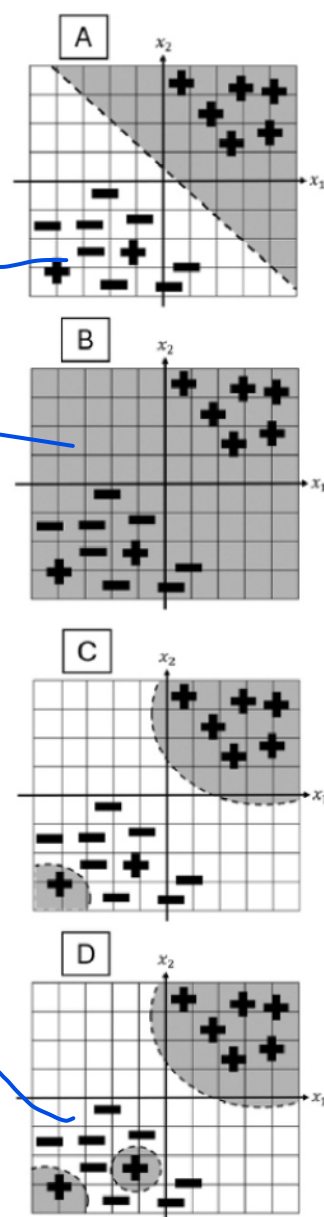
משמאל נתונות 4 גרסאות שונות של בעיית SVM. בנוסף נתונה פונקציה למיפוי פיצורים פולינומיאלית מממד גבוה  $\phi$ .

1) ע"י מתיחת קו, לכל בעיית SVM התאימו איור אחד שמייצג את ה-decision boundary של המודל שיתקבל מבעיה זו. אם בעיית ה-SVM לא יכולה להגיע לפתרון, סמנו עליה [X].

תזכורת: כלל ההחלטה עבור מפריד לינארי ללא מיפוי  $\phi$   $h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$ .

עם מיפוי  $\phi$   $h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \phi(\mathbf{x}) + b)$ .

I) $\underset{\mathbf{w}, b}{\text{argmin}} \ \mathbf{w}\ _2^2 + b$
II) $\underset{\mathbf{w}, b}{\text{argmin}} \ \mathbf{w}\ _2^2$ such that: $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \forall i$
III) $\underset{\mathbf{w}, b}{\text{argmin}} \ \mathbf{w}\ _2^2$ such that: $y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1, \forall i$
IV) $\underset{\mathbf{w}, b}{\text{argmin}} \ \mathbf{w}\ _2^2 + \sum_i \max\{0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b)\}$



(2) נמקו בקצרה את בחירתכם.

- I) 0 vector
- II) no graph
- III) hard svm for high degree
- IV) soft SVM

(3) מבין ארבעת בעיות ה-SVM, עבור מי הביטוי  $\mathbb{E}_{S,x} \left[ \left( h_S(x) - \bar{h}(x) \right)^2 \right]$  הוא הגדול ביותר? נמקו בקצרה.  
 הערה: התוחלת כאן היא ביחס למשתנים המקריים  $S, x$ .  
 תזכורת: הפונקציה  $\bar{h}(x)$  היא הממוצע של מחלקת ההיפותוזות המוגדרת ע"י בעיית ה-SVM.

III) highest complexity -> highest Variance

ב. [20 נק'] בסעיפים הבאים ננתח תרחיש של imbalanced sampling, תרחיש שבו התפלגות התיוגים בקבוצת מדגם  $S$  היא לא מאוזנת.

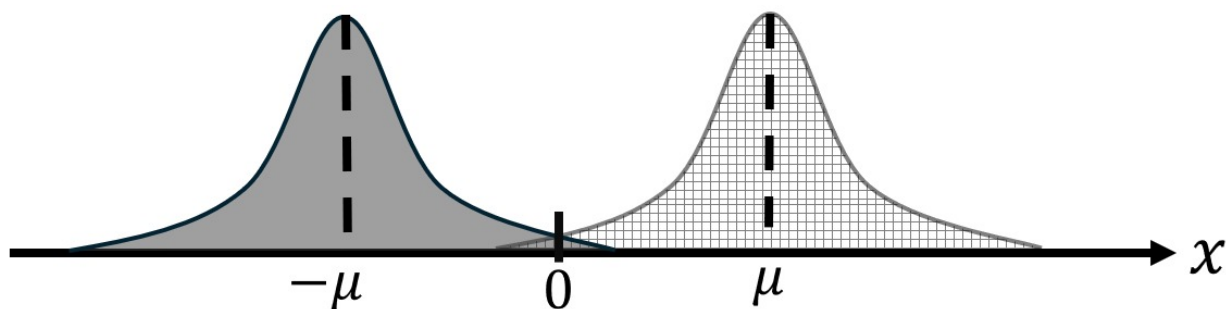
ננתח את המקרה שבו  $d = 1$ . ידוע כי הפיצ'ר  $x$  מתפלג כמו משתנה אקראי גאוס, כאשר תוחלת הגאוסיאן תלויה בערך של  $y$ . השונות זהה בשני המקרים. יהי  $\mu > 0$ , אז:

$$\mathbb{P}(x|y = 1) = N(\mu, \sigma^2) \quad \circ$$

$$\mathbb{P}(x|y = -1) = N(-\mu, \sigma^2) \quad \circ$$

$$\mathbb{P}(y = 1) = \mathbb{P}(y = -1) = 1/2 \quad \circ$$

(1) להלן שרטוט של התפלגות המשותפת (עם צביעה שונה להתפלגויות המותנות השונות):



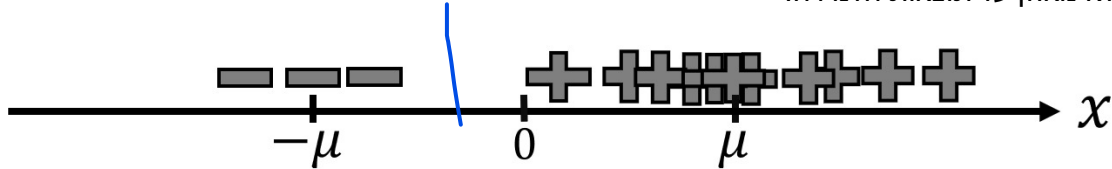
נניח כי יש לנו גישה להתפלגות עצמה, כלומר אנחנו יודעים את  $\mu, \sigma^2$ . מהו המפריד מסוג threshold (לינארי) אשר ישיג את שגיאת ההכללה המינימאלית על פילוג זה.

כתבו את כלל ההחלטה שלו (כלומר כלל הפרדיקציה  $\hat{y} = h(x)$ ) באופן מפורש, ונמקו בקצרה.

$$h(x) = \text{sign}(x)$$

נניח כעת כי יש לנו גישה רק לדגימה מההתפלגות (כלומר אנחנו לא יודעים את  $(\mu, \sigma^2)$ ). יהי  $S = \{(x_i, y_i)\}_{i=1}^n$  מדגם בגודל  $n$  אשר כל דגימה  $(x_i, y_i)$  נדגמה  $i.i.d$  מההתפלגות המשותפת  $\mathbb{P}(x, y)$ . נסמן ב- $n_+$  את מספר הדוגמאות החיוביות במדגם, וב- $n_-$  את מספר הדוגמאות השליליות, כך ש- $n = n_+ + n_-$ . נגיד שהמדגם "לא מאוזן" כאשר  $n_+$  גדול משמעותית מ- $n_-$  (או להפך).

(2) להלן איור של מדגם מסוים  $S$  שאינו מאוזן. נניח שלמדנו בעזרת מסווג  $\text{threshold}$  עם אלגוריתם SVM. סמנו על האיור את המקום בו עובר המפריד הנלמד. נמקו את בחירתכם, והסבירו את ההשלכות של מדגם לא מאוזן על תוצאות הלמידה.



will be right in the middle point of the right most (-) and the leftmost (+) as a result of trying to maximize the margin

as a result it will have a higher generalization error and will even be mistaken on (-) that are between the  $\theta$  and 0

(3) עבור התפלגות נתונה ומדגם בגודל  $n$  ממנה, נגדיר את הגדלים הבאים:

$$M(n) = \mathbb{E}[\max(x_1, \dots, x_n)] \quad \circ$$

$$m(n) = \mathbb{E}[\min(x_1, \dots, x_n)] \quad \circ$$

כאשר מדובר בהתפלגות גאוסיאנית עם תוחלת  $\mu$  (ושונות קבועה  $\sigma^2$ ), נסמן ב- $M(n; \mu)$  ו- $m(n; \mu)$  את הגדלים המתאימים.

$$m(n; 0) \cong -\sigma\sqrt{2\ln(n)}, \quad M(n; 0) \cong \sigma\sqrt{2\ln(n)}$$

השלימו:

$$M(n; -\mu) \cong M(n; 0) - \mu$$

$$M(n; \mu) \cong M(n; 0) + \mu$$

$$m(n; -\mu) \cong m(n; 0) - \mu$$

$$m(n; \mu) \cong m(n; 0) + \mu$$

(4) נתון שדגמתם מדגם לא מאוזן אך פריד. שיטה נפוצה להתמודדות עם בעיה זו נקראת subsampling. במסגרת שיטה זו, מוציאים באקראי מ- $S$  דוגמאות **בעלות התיוג הנפוץ ביותר**, עד שמקבלים קבוצת מדגם מאוזנת, כלומר עד ש- $n_+ = n_-$ .

השתמשו בסעיפים הקודמים, ובפרט בסעיף (3), על מנת להסביר כיצד subsampling יכול לסייע לאימון SVM בתרחיש של שאלה זו.  
הדרכה: שימו לב כי הגדלים  $M(n), m(n)$  תלויים ב- $n$ . בנוסף זכרו מה התכונה של המפריד הלינארי שמחזיר SVM.

this will lower the amount of samples until the Expected Max & Min values will be of the same (seder godel)

therefore the Expected (mafrid) will be better generalized

and we'll get the optimal one (sign(x)) as a result of  $(M(n;0) - \mu + m(n;0) + \mu)^2 = 0$ ;

---

---

---

---

---

---

---

---

## שאלה 2: VC-dimension [25 נק']

א. [2 נק'] להלן ההגדרה של "ניתוח". השמטנו מההגדרה את הפקדים.

השלימו את שלושת הכמתים החסרים. בכל מקום כתבו בבירור האם חסר בהגדרה  $\forall$  או  $\exists$ .

$$\mathcal{H} \text{ shatters } C \Leftrightarrow \underbrace{\text{A}}_{\text{השלימו}} y_1, \dots, y_{|C|} \in \mathcal{Y}: \underbrace{\text{E}}_{\text{השלימו}} h \in \mathcal{H}: \underbrace{\text{A}}_{\text{השלימו}} x_i \in C: h(x_i) = y_i$$

בשאלה זו נתון מרחב דוגמאות  $\mathcal{X}$  כלשהו ומרחב תיוגים  $\mathcal{Y} = \{-1, +1\}$ . אם נדרשת הוכחה, הוכיחו באופן פורמלי. אם נדרשת הפרכה, תנו דוגמה נגדית מנומקת היטב והוכיחו כי היא אכן מפריכה את הטענה.

ב. [5 נק'] הוכיחו/הפריכו.

יהיו שתי מחלקות  $\mathcal{H}_1, \mathcal{H}_2$ . אזי

$$VCdim(\mathcal{H}_1 \cup \mathcal{H}_2) \geq \max\{VCdim(\mathcal{H}_1), VCdim(\mathcal{H}_2)\}$$

proof:

let  $\max\{VC(\mathcal{H}_1), VC(\mathcal{H}_2)\} = T$

let's say  $VC(\mathcal{H}_1) \geq VC(\mathcal{H}_2)$ ;

therefore there is a set of  $T$  points  $x_1, \dots, x_T$  that for each possible labelling  $y_1, \dots, y_T$  there exists a  $h$  in  $\mathcal{H}_1$  s.t  $h(x_i) = y_i$

let's look at  $T$  and  $\mathcal{H}_1 \cup \mathcal{H}_2$ ;  $\mathcal{H}_1$  shatters  $T$  - therefore all hypotheses  $h$  that help  $\mathcal{H}_1$  shatter  $T$  are also in  $\mathcal{H}_1 \cup \mathcal{H}_2$ ; therefore  $\mathcal{H}_1 \cup \mathcal{H}_2$  shatters  $T$  and  $VC(\mathcal{H}_1 \cup \mathcal{H}_2) \geq T = \max(\dots, \dots)$

ג. [6 נק'] הוכיחו.

נתונה מחלקה  $\mathcal{H}$  סופית. אזי

$$VCdim(|H|) \leq \log_2 |H|$$

הוכחה פורמלית:

let  $VC(H) = d$   
then there exists a set of  $x_1, \dots, x_d$  st for all  $y_1, \dots, y_d$  there exists  $h$  in  $H$  st for all  $i$ :  $h(x_i) = y_i$   
therefore  $|H| \geq 2^d \Rightarrow \log |H| \geq d = VC(H)$



ד. [12 נק'] הוכיחו/הפריכו.

יהיו  $\mathcal{H}_1, \mathcal{H}_2$  מחלקות ונתונה  $\emptyset \neq C \subset \mathcal{X}$ . נניח

(1)  $\mathcal{H}_1$  מנתצת את  $C$ .

(2)  $\mathcal{H}_2$  אינה מנתצת את  $C$ .

אזי בהכרח

$$VCdim(\mathcal{H}_1) > VCdim(\mathcal{H}_2)$$

הוכחה פורמלית:

$$H1 = \{h(x1) = 1 \wedge h(x2) = -1, h'(x1) = -1 \wedge h'(x2) = 1\}$$

$$H2 = \{h''(x1) = 1 \wedge h''(x2) = 1, h'''(x1) = 1 \wedge h'''(x2) = -1\}$$

$VCdim(H1) \leq \log H1 = 1$ , let's take  $x = x1$ ;  $H1$  shatters it  $\Rightarrow VC(H1) = 1$   
same for  $H2$  except  $x = x2$ ;  $VC(H2) = 1$ ;

let  $C$  be  $\{x1\}$ ,  $H1$  shatters it while  $H2$  doesn't, yet  $VC(H1) = VC(H2)$   
therefore we disproved the claim

## שאלה 3: Perceptron [20 נק']

לזכאים להערכה חלופית בלבד (כפי שהוגדרו באתר הקורס): סמנו את התיבה הזו אם ברצונכם לדלג על שאלה זו. המשקל של יתר השאלות יתפזר באופן יחסי על פני 100 נקודות. ניתן לדלג רק על שאלה אחת מתוך שאלות 3,4.

☐

בשאלה זו נניח  $\mathcal{Y} = \{+1, -1\}$  ומרחב דוגמאות  $\mathcal{X} = \mathbb{R}^d$ . לפניכם אלגוריתם הפרספטרון כפי שהוא הוצג בתרגול.

**input:** training set  $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$ , step size  $\eta = 1$

$\mathbf{w} = \mathbf{0}_d$

**while** did not separate the training set:

**for**  $i = 1$  to  $m$ :

$\hat{y}_i = \text{sign}(\mathbf{w}^\top \mathbf{x}_i)$

**if**  $y_i \neq \hat{y}_i$ :

$\mathbf{w} = \mathbf{w} + y_i \mathbf{x}_i$

**return**  $\mathbf{w}$

לפניכם משפט התכנסות הפרספטרון:

תהי קבוצת אימון  $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$ . נניח כי קיימים וקטור משקולות  $\mathbf{w}_*$  כך ש-  $\|\mathbf{w}_*\|_2 = 1$  ו-  $\gamma > 0$  כך שלכל  $i = 1, \dots, m$  מתקיים,

$$y_i(\mathbf{w}_*^\top \mathbf{x}_i) \geq \gamma$$

נניח בנוסף כי לכל  $i = 1, \dots, m$  מתקיים  $\|\mathbf{x}_i\|_2 \leq R$ .

אזי אלגוריתם הפרספטרון עושה לכל היותר  $\frac{R^2}{\gamma^2}$  טעויות.

בשאלה זו נוכיח משפט זה בשלבים. תהי  $S$  קבוצת אימון ונניח את קיומם של  $\gamma, R, \mathbf{w}_*$  כמו במשפט.

נגדיר את  $\mathbf{w}_k$  להיות וקטור המשקולות שנלמד עד הטעות ה-  $k$  (לא כולל).

שימו לב: עבור הגדרה זו מתקיים  $\mathbf{w}_1 = \mathbf{0}_d$ .

א. [2 נק'] האם  $S$  פרידה לינארית? נמקו בקצרה.

yes because there is a gamma > 0 st :  $y(\mathbf{w}_*) \geq \gamma > 0$  therefore it (mesaveg) each point correctly therefore it is seperable

---



---



---



---



---



---



---



---

ב. [5 נק'] נתחיל בלהוכיח באינדוקציה את הטענה הבאה:

$$\mathbf{w}_{k+1}^T \mathbf{w}_* \geq k\gamma$$

הדרכה:

- בדקו את בסיס האינדוקציה עבור  $k = 0$ .
- הניחו את נכונות הטענה עבור  $k$  כלשהו, כלומר:  $\mathbf{w}_k^T \mathbf{w}_* \geq (k-1)\gamma$ .
- הוכיחו באמצעות כלל העדכון של הפרספטרון כי  $\gamma + \mathbf{w}_k^T \mathbf{w}_* \geq \mathbf{w}_{k+1}^T \mathbf{w}_*$ .
- השלימו את צעד האינדוקציה.

$$w_{(k+1)} * w^* = (w_{(k)} + \gamma x)w^* \geq (k-1)\gamma + \gamma = k\gamma \text{ as wanted}$$

ג. [2 נק'] לרשותכם אי-שוויון קושי שזורץ

$$\mathbf{u} \in \mathbb{R}^d, \mathbf{v} \in \mathbb{R}^d \quad \|\mathbf{u}\|_2 \|\mathbf{v}\|_2 \geq |\mathbf{u}^T \mathbf{v}|$$

הסיקו

$$\|\mathbf{w}_{k+1}\|_2 \geq \mathbf{w}_{k+1}^T \mathbf{w}_*$$

ד. [9 נק'] הוכיחו באינדוקציה את הטענה

$$\|\mathbf{w}_{k+1}\|_2^2 \leq kR^2$$

○ בדקו את בסיס האינדוקציה עבור  $k = 0$ .

○ הניחו את נכונות הטענה עבור  $k$  כלשהו, כלומר:  $\|\mathbf{w}_k\|_2^2 \leq (k-1)R^2$

○ הוכיחו באמצעות כלל העדכון של הפרספטורן כי  $\|\mathbf{w}_{k+1}\|_2^2 \leq \|\mathbf{w}_k\|_2^2 + R^2$ .

- השלימו את צעד האינדוקציה.

זכרו: עבור  $\mathbf{v} \in \mathbb{R}^d$  מתקיים  $\|\mathbf{v}\|_2^2 = \mathbf{v}^\top \mathbf{v}$ .

[illegible]

$$k^2 \gamma^2 \leq \|\mathbf{w}_{k+1}\|_2^2 \leq kR^2$$

ה. [2 נק'] הסיקו כעת את המשפט. כלומר:

$$k \leq \frac{R^2}{\gamma^2}$$

כאשר  $k$  הוא מספר הטעויות.

## שאלה 4: Deep Learning [20 נק']

לזכאים להערכה חלופית בלבד (כפי שהוגדרו באתר הקורס): סמנו את התיבה הזו אם ברצונכם **לדלג** על שאלה זו. המשקל של יתר השאלות יתפזר באופן יחסי על פני 100 נקודות. ניתן לדלג רק על שאלה אחת מתוך שאלות 3,4. ☐

בשאלה זו נעבוד מעל מרחב דוגמאות  $\mathcal{X} = \mathbb{R}$  ומרחב תיוגים  $\mathcal{Y} = \{0,1\}$ . בסעיפים בהם נדרש חישוב, עגלו את תשובותיכם לדיוק של עד שלוש ספרות אחרי הנקודה.

נתונה רשת נוירונים שעוצבה לפתירת בעיית סיווג בינארי המוגדרת ע"י הקשרים הבאים:

$$\begin{bmatrix} a_1 = w_{11}x + b_{11} \\ a_2 = w_{12}x + b_{12} \\ z_1 = \text{ReLU}(a_1) \\ z_2 = \text{ReLU}(a_2) \\ a_3 = w_{21}z_1 + w_{22}z_2 + b_{21} \\ \hat{y} = \sigma(a_3) \end{bmatrix}$$

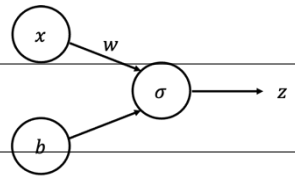
כאשר  $x \in \mathbb{R}$  הוא הקלט ו-  $\hat{y} \in [0,1]$  הוא הפלט של הרשת.

האימון של הרשת מתבצע באמצעות cross-entropy loss. תזכורת:

$$\begin{bmatrix} \text{ReLU}(x) = \max(0, x) \\ \sigma(x) = \frac{1}{1 + e^{-x}} \\ \ell^{CE}(y, \hat{y}) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y}) \end{bmatrix}$$

א. [2 נק'] ציירו תרשים של הרשת (לרשותכם דוגמה לתרשים בגוף התשובה).

ציינו מי הן המשקולות של הרשת שביחס אליהן אנו עושים אופטימיזציה.



$z = \sigma(xw + b)$

---

---

---

---

---

---

---

---

---

---

$$\begin{bmatrix} b_{11} = & 0.04 \\ b_{12} = & 0.01 \\ b_{21} = & 0.08 \\ w_{11} = & 0.2 \\ w_{12} = & -0.1 \\ w_{21} = & 0.7 \\ w_{22} = & 0.7 \end{bmatrix}$$
[illegible]

ג. [1 נק'] באיזה כלל אנו משתמשים בשביל לחשב את הנגזרות החלקיות ביחס למשקולות של רשת נוירונים כלשהי? סמנו את התשובה הנכונה:

1. כלל הפיצה.
2. כלל האצבע.
3. כלל השרשרת.
4. כלל השורש.

בסעיף הבא הניחו כי עבור הדוגמה  $x = 0.3$  התיוג האמיתי שלה הוא  $y = 1$ .

נעסוק כעת באלגוריתם ה-backpropagation.

ד. [12 נק'] בצעו את אלגוריתם backpropagation על המשתנה  $b_{12}$ . עליכם לכתוב את הנגזרת החלקית של פונקציית ההפסד  $\ell(y, \hat{y})$  ביחס למשתנה  $b_{12}$ , דהיינו  $\frac{\partial \ell}{\partial b_{12}}$ , באמצעות הנגזרות החלקיות  $\frac{\partial \alpha}{\partial \beta}$ , כאשר  $\alpha, \beta$  יכולים להיות כל אחד מהבאים:

$$\ell, \hat{y}, z_i, a_i, b_{ij}, w_{ij}, x$$

עבור כל הערכים החוקיים של  $i, j$ . וודאו כי כל נגזרת חלקית  $\frac{\partial \alpha}{\partial \beta}$  לא ניתנת לפירוק לנגזרות חלקיות פשוטות יותר.

לאחר מכן, **חשבו** את  $\frac{\partial \ell}{\partial b_{12}}$  עבור  $x = 0.3$ . תזכורת:  $\frac{d}{dx} \sigma(x) = \sigma(x) \cdot (1 - \sigma(x))$

$$\frac{\partial \ell}{\partial b_{12}} =$$



מסגרת נוספת (יש לציין אם מדובר בטייטה או בהמשך לתשובה אחרת):

This image shows a single page of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are approximately 25 lines in total. The left edge of the paper has rounded corners, while the right edge is straight. The paper appears to be part of a notebook or a set of loose-leaf papers.

מסגרת נוספת (יש לציין אם מדובר בטיוטה או בהמשך לתשובה אחרת):

A large rectangular box with rounded corners, containing 25 horizontal lines for writing. The box is empty, with the lines spaced evenly from top to bottom.

מסגרת נוספת (יש לציין אם מדובר בטיוטה או בהמשך לתשובה אחרת):

A large rectangular box with rounded corners, containing 20 horizontal lines for writing. The lines are evenly spaced and extend across the width of the box. The box is intended for a student to provide a second answer or a clarification to the question above.