



מבוא למערכות לומדות (236756)

סמסטר חורף תשפ"ג – 14 במרץ 2023

מרצה: ד"ר יונתן בלינקוב

## מבחן מסכם מועד ב'

### 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:
- תשובה בכתב יד לא קריא – **לא תיבדק**.
- בשאלות רב-ברירה – הקיפו את התשובות בבירור. סימונים לא ברורים יביאו לפסילת התשובה.
- לא יתקבלו ערעורים בנושא.
- במבחן 14 עמודים ממוספרים סה"כ, כולל עמוד שער זה שמספרו 1 ושלושה עמודי טיוטה בסוף הגליון.
- נא לכתוב רק את המבוקש ולצרף הסברים קצרים עפ"י ההנחיות.
- **בתום המבחן יש להגיש את שאלון זה בלבד.**

בהצלחה!

## חלק א' – שאלות פתוחות [82 נק']

### שאלה 1: השפעה של דוגמה יחידה על פעולת מסווגים [24 נק']

נתון סט אימון עם  $m \geq 10$  דוגמאות דו-ממדיות ותיגים בינאריים, משמע לכל  $i = 1, \dots, m$  מתקיים  $y_i \in \{-1, 1\}$ ,  $x_i \in \mathbb{R}^2$ .

לומדים שני מסווגים:

- בשלב הראשון: לומדים מסווג על סט האימון המקורי ומחשבים את הסיווגים על כל הדוגמאות.
- בשלב השני:
  - מסירים דוגמת אימון אחת שרירותית כלשהי.
  - מאמנים מסווג חדש על סט האימון המעודכן, ומחשבים בעזרתו את הסיווג על כל הדוגמאות הנותרות.

עבור כל אלגוריתם למידה, סמנו האם הסיווגים שהמסווג החדש יחזה על דוגמאות האימון הנותרות זהים בהכרח לאלה של המסווג המקורי על דוגמאות אלה.

### הסבירו בקצרה את תשובותיכם (2-4 משפטים בכל סעיף).

הניחו שאין צעדים אקראיים או שגיאות נומריות בריצת האלגוריתמים (בעיות קמורות מתכנסות לפתרון האנליטי במדויק).

א. Soft-SVM ליניארי לא הומוגני עם  $\lambda = 10^{-1}$  בהנחה שהדאטה המקורי פריד ליניארית.

הסיווגים על דוגמאות האימון הנותרות זהים בהכרח? כן / לא

הסבר:

---



---



---



---

ב. Soft-SVM ליניארי לא הומוגני עם  $\lambda \rightarrow 0$  בהנחה שהדאטה המקורי פריד ליניארית.

הסיווגים על דוגמאות האימון הנותרות זהים בהכרח? כן / לא

הסבר:

---



---



---



---

ג. ID3 המשתמש באנטרופיה ועוצר בעומק מירבי 4. הסיווגים על דוג' האימון הנותרות זהים בהכרח? **כן / לא**

הסבר:

ד. kNN עם  $k = 3$  (דוגמה לא נחשבת שכנה של עצמה), כאשר ידוע שלשלושת השכנים הקרובים ביותר

לדוגמה שהוסרה יש תיוג זהה לתיוג שלה. הסיווגים על דוגמאות האימון הנותרות זהים בהכרח? **כן / לא**

הסבר:

## שאלה 2: Generative models [26 נק']

תזכורת: פונק' הצפיפות של התפלגות  $U[a, b]$ , אחידה ורציפה על הקטע הסגור  $[a, b]$ , היא

$$f(z) = \frac{1}{b-a} \mathbb{I}[a \leq z \leq b] = \begin{cases} \frac{1}{b-a}, & a \leq z \leq b \\ 0, & \text{otherwise} \end{cases}$$

א. [5 נק'] **מתרגיל בית:** נתון משתנה אקראי  $X \sim U[0, \theta]$  עבור  $\theta > 0$  לא ידוע.

נתון מדגם אקראי  $S$  של  $m$  דגימות,  $S = \{x_1, \dots, x_m\} \subset \mathbb{R}_{\geq 0}$ , שנדגמו מהמשתנה האקראי באופן i.i.d.

הוכיחו שמשערך ה-MLE שמוגדר בתור  $\hat{\theta}_{\text{MLE}} \triangleq \underset{\text{likelihood}}{\operatorname{argmax}_{\theta}} \Pr[S; \theta]$  הוא  $\hat{\theta}_{\text{MLE}} = \max_{i \in [m]} x_i$ .

תשובה:

---

---

---

---

---

---

---

---

---

---

ב. [6 נק'] בנוסף על הנתונים שבסעיף הקודם, בסעיף זה בלבד נתון שמתקיים  $\theta \sim U[10, 20]$ .

מצאו (והוכיחו) את משערך ה-MAP לפי כלל הנתונים:  $\hat{\theta}_{\text{MAP}} \triangleq \arg\max_{\theta} \Pr[\theta|S] = \arg\max_{\theta} \underbrace{\Pr[S|\theta]}_{\text{likelihood}} \underbrace{\Pr[\theta]}_{\text{prior}}$

תשובה:

---

---

---

---

---

---

---

---

---

---

בסעיף הבא מרחב הדוגמאות הוא  $\mathcal{X} = \mathbb{R}_{\geq 0}^2$  (הרביע החיובי) ומרחב התיוגים הוא  $\mathcal{Y} = \{-1, +1\}$ . בהינתן מדגם אימון  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\} \subset (\mathcal{X} \times \mathcal{Y})$ , נרצה ללמוד מסווג בינארי.

### תהליך הלמידה:

i. נניח את הנחת Naïve Bayes (NB).

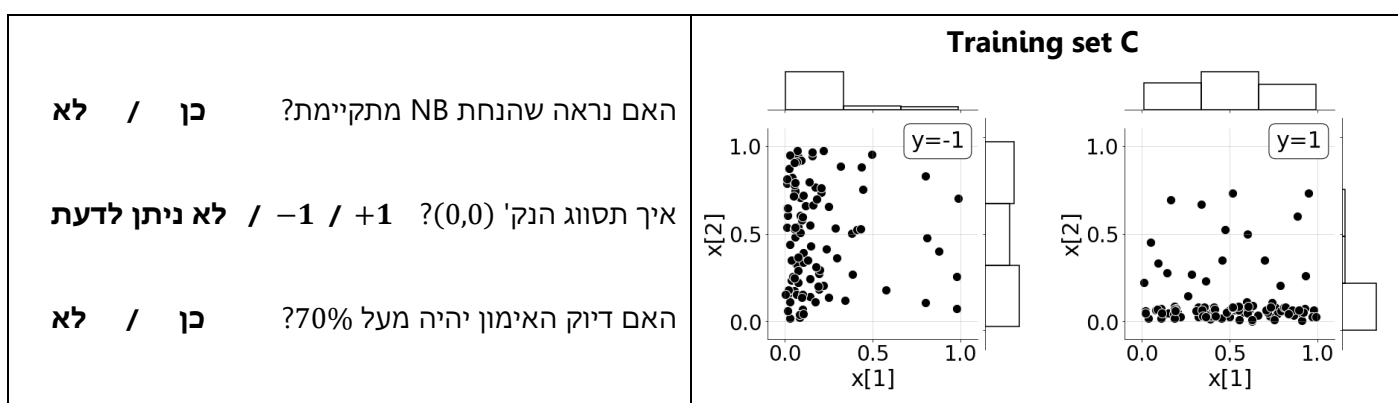
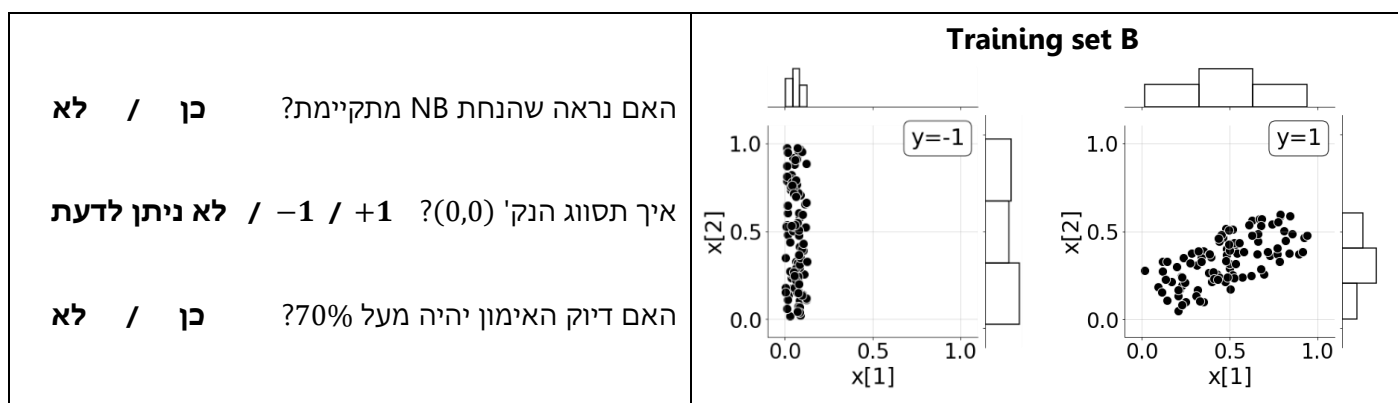
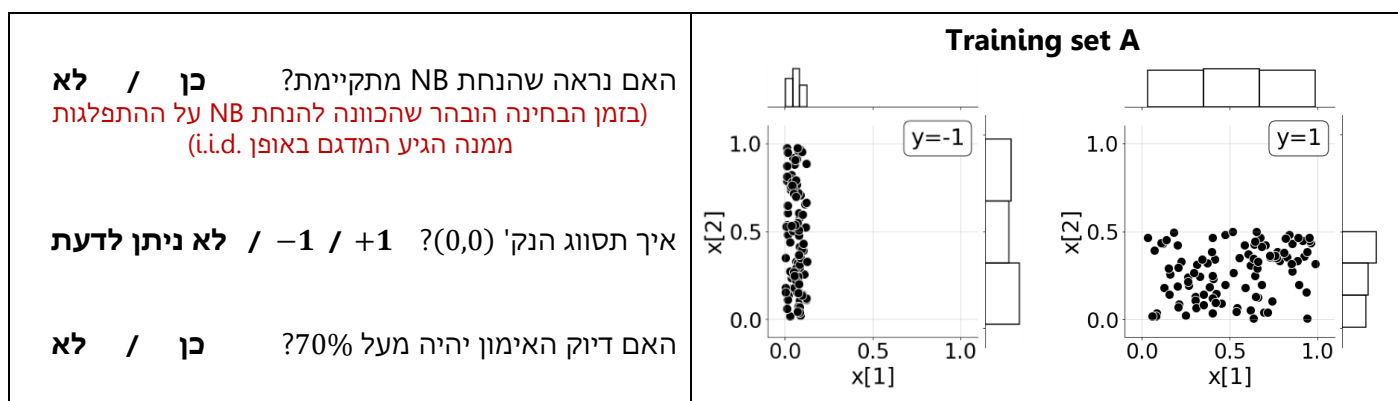
ii. נמדל את בעיות הסיווג בעזרת Uniform NB, משמע,  $X[j]|Y=k) \sim U[0, \theta_k[j]]$ , כאשר  $j \in \{1, 2\}, k \in \{-1, +1\}$ .

iii. נשערך את ארבעת הפרמטרים בעזרת MLE, משמע  $\hat{\theta}_{-1} = \left[ \begin{matrix} \max_{i:y_i=-1} x_i[1] \\ \max_{i:y_i=-1} x_i[2] \end{matrix} \right]$  ו-  $\hat{\theta}_{+1} = \left[ \begin{matrix} \max_{i:y_i=+1} x_i[1] \\ \max_{i:y_i=+1} x_i[2] \end{matrix} \right]$ .

iv. בהמשך לכל ההנחות לעיל, נבנה כלל החלטה הסתברותי  $\hat{y}(\mathbf{x}) = \operatorname{argmax}_{y \in \{-1, +1\}} \Pr(\mathbf{x}; y)$ .

כעת נתונים שלושה מדגמי אימון, כל אחד מהתפלגות שונה ומכיל 100 דוגמאות חיוביות ו-100 שליליות. המדגמים מופיעים בתרשימים הבאים (הדוגמאות מכל תיוג מופיעות בנפרד ביחד עם ההיסטוגרמות השוליות המתאימות). לכל מדגם (בנפרד), מבצעים את תהליך הלמידה המתואר לעיל.

ג. [15 נק'] לכל מדגם, ענו על השאלות ביחס לתהליך הלמידה שלו. התשובות אמורות להיות ברורות מהגרפים.



## שאלה 3: Multi-Layer Perceptron (MLP) and VC dimension [23 נק']

קראו היטב את הנתונים הבאים.

בשאלה זו מרחב הנתונים הוא  $\mathcal{X} = \mathbb{R}^2, \mathcal{Y} = \{-1, +1\}$ .נגדיר מחלקה  $\mathcal{H}$  של רשתות MLP עם שכבה חבויה אחת ברוחב  $p \in \mathbb{N}$  (היפר-פרמטר), אקטיביציות ReLU ופלט בינארי יחיד.

בכל הרשתות במחלקה, המשקלים של השכבה השנייה קבועים להיות 1 וללא bias.

נאמר שאוסף פרמטרים  $\theta$  הוא חוקי, אם המשקלים בו אי-שליליים (אילוץ זה לא כולל את רכיבי ה-bias).נתאר את הפונקציה המתקבלת  $F_\theta: \mathbb{R}^2 \rightarrow \{-1, +1\}$  בצורה גרפית ובצורה פורמלית:

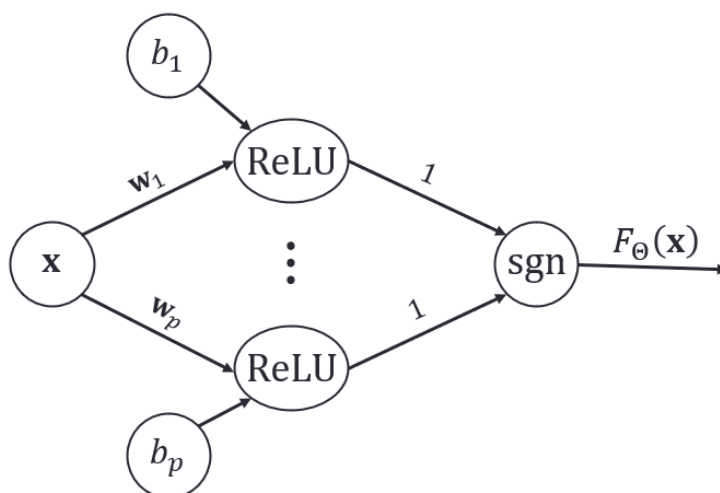
$$F_\theta(\mathbf{x}) = \text{sgn} \left( \sum_{t=1}^p \text{ReLU}(\mathbf{w}_t^\top \mathbf{x} + b_t) \right),$$

where:

$$\theta = (\mathbf{w}_1, \dots, \mathbf{w}_p, b_1, \dots, b_p),$$

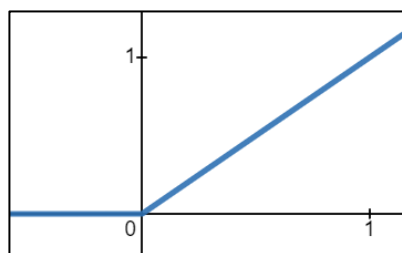
$$\mathbf{w}_1, \dots, \mathbf{w}_p \in \mathbb{R}_{\geq 0}^2,$$

$$b_1, \dots, b_p \in \mathbb{R}.$$



וכמו כן,

$$\text{ReLU}(z) = \begin{cases} 0, & z \leq 0 \\ z, & z > 0 \end{cases} \quad \text{תזכורת:}$$



$$\text{sgn}(z) = \begin{cases} -1, & z \leq 0 \\ +1, & z > 0 \end{cases} \quad \text{נגדיר:}$$

(כך שלא מתקבל אפס לשום קלט)

סימון: יהיו שני קלטים  $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^2$ . נסמן  $\mathbf{x}_i \succcurlyeq \mathbf{x}_j$  אם ורק אם  $\mathbf{x}_i[1] \geq \mathbf{x}_j[1] \wedge \mathbf{x}_i[2] \geq \mathbf{x}_j[2]$ .

א. [8 נק'] הוכיחו: ברשתות שהגדרנו, לכל רוחב  $p$  ולכל  $\theta$  חוקי, אם  $\mathbf{x}_i \succcurlyeq \mathbf{x}_j$  אזי  $F_\theta(\mathbf{x}_i) \geq F_\theta(\mathbf{x}_j)$ .

הוכחה (לרשותכם טיוטה בסוף הגיליון):

---

---

---

---

---

---

---

---

---

---

ב. [8 נק'] נגדיר את ה-XOR dataset:  $S = \left\{ \left( \underbrace{(0,1)}_{\mathbf{x}_1}, \underbrace{+1}_{y_1} \right), \left( \underbrace{(1,0)}_{\mathbf{x}_2}, \underbrace{+1}_{y_2} \right), \left( \underbrace{(0,0)}_{\mathbf{x}_3}, \underbrace{-1}_{y_3} \right), \left( \underbrace{(1,1)}_{\mathbf{x}_4}, \underbrace{-1}_{y_4} \right) \right\}$

הוכיחו שלכל רוחב  $p$  ולכל  $\theta$  חוקי, לא ניתן להגיע לשגיאת אימון אפס על  $S$  ע"י  $F_\theta \in \mathcal{H}$  (כדאי להיעזר בסעיף הקודם).  
**הבהרה:** כל הסעיפים עוסקים ב-capacity של המחלקה ולא במציאת דרכי אימון יעילות.

הוכחה:

---

---

---

---

---

---

---

---

---

---



יהיו  $n \geq 2$  קלטים  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}_{\geq 0}^2$  שונים ומנורמלים ברביע החיובי (משמע  $\forall i: \mathbf{x}_i \neq \mathbf{x}_j$  וגם  $\|\mathbf{x}_i\|_2 = 1$ ).

משמע, הציעו השמה כזאת (שתלויה ב- $\mathbf{x}_1, \dots, \mathbf{x}_n$ ) והוכיחו שהיא מקיימת את הנדרש.

הוכחה (לרשותכם טיוטה בסוף הגיליון):

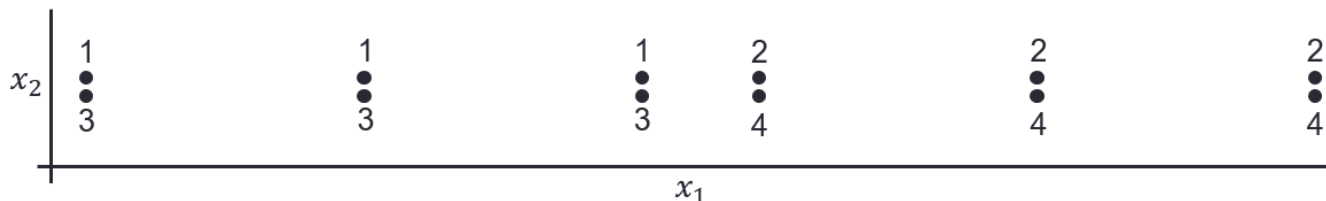
[illegible]

## הוכחה:

## חלק ב' – שאלות רב-ברירה [18 נק']

בשאלות הבאות סמנו את התשובות המתאימות (לפי ההוראות). בחלק זה אין צורך לכתוב הסברים.

א. לפניכם סט אימון דו-ממדי עם 4 מחלקות ו-3 דוגמאות מכל מחלקה (התיוג כתוב מעל/מתחת הדוגמאות).



מבין מודלי ה-multiclass הבאים, סמנו את כָּל אלה שצפויים להגיע לדיוק אימון של 100% על הדאטה לעיל.

- 1-nearest-neighbor (חצה את התיוג של השכן הקרוב ביותר לפי מרחק אוקלידי, דוג' לא נחשבת שכנה של עצמה).
- עץ החלטה בעומק מירבי 3 (הפרדיקציה של כל עלה נקבעת לפי רוב דוגמאות האימון שבתוכו).
- מודל one-vs-one עם decision stump (עץ בעומק 1) כמודל בסיס.
- מודל one-vs-all עם decision stump (עץ בעומק 1) כמודל בסיס.

ב. נגדיר אלגוריתם Random Forest פשוט:

**Random Forest**( $S, k, \text{max\_depth}, \text{min\_samples\_split}$ ):

For  $i=1$  to  $k$ :

$S' = \text{Sample } \sqrt{d} \text{ features out of the original } d \text{ features in } S \text{ (keeping all samples)}$

$h_i = \text{ID3}(S', \text{max\_depth}, \text{min\_samples\_split}, \text{criterion}=\text{"entropy"})$

Return  $H(x) = \frac{1}{k} \sum_{i=1}^k h_i(x)$

אילו מבין הבחירות האלגוריתמיות הבאות צפויות להפחית את ה-Variance של המסווג הכולל שנלמד  $H$ ?

סמנו את כָּל התשובות הנכונות (השאלה אינה עוסקת במקרי קצה אלא במקרה הסביר).

- הגדלת  $k$  (מספר העצים ביער).
- הגדלת  $\text{max\_depth}$  (העומק המירבי המותר).
- הגדלת  $\text{min\_samples\_split}$  (מספר הדוגמאות המינימלי הנדרש לפיצול של צומת).
- נירמול מקדים של הדאטה בשיטת min-max.
- נירמול מקדים של הדאטה בשיטת standardization (Z-score).

(יש שאלה נוספת בעמוד הבא)

נתונה פונקציית מיפוי כלשהי  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{10}$ .

$\underset{\mathbf{w}}{\text{argmin}} \left( \frac{1}{m} \sum_i \max\{0, 1 - y_i \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_i)\} \right)$ 
נגדיר בעיית אופטימיזציה בעזרת hinge loss על  $S$ :

a. כן.

c. רק אם הפונקציה  $\phi$  לא ליניארית.

e. רק אם  $c \geq 1$ .

f. לא, כי חסר גורם רגולריזציה.

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

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are approximately 20 lines visible. The paper has rounded corners on the left side and a straight edge on the right. There is no handwriting or other markings on the paper.

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 across the page.

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 lines are evenly spaced and extend across the width of the box. The box is intended for providing a second answer or clarification to the question above.