



מבוא למערכות לומדות (236756)

סמסטר אביב תשפ"ב – 12 ביולי 2022

מרצה: ד"ר ניר רחנפלד

## מבחן מסכם מועד א' – פיתרון חלקי

שימו לב: הפתרונות המופיעים כאן הם חלקיים בלבד ומובאים בשביל לעזור לכם בתהליך הלמידה. ייתכנו כאן חוסרים / ליקויים / טעויות של ממש.

### הנחיות הבחינה:

- **משך הבחינה:** 3 שעות.
- **חומר עזר:** המבחן בחומר סגור (ללא ספרים, מחברות, דפי נוסחאות).
- מותר להשתמש במחשבון.
- יש לכתוב בעט **בלבד**.
- יש לכתוב את התשובות **על גבי שאלון זה**.
- מותר לענות בעברית או באנגלית.
- קריאות:
  - תשובה בכתב יד לא קריא – **לא תיבדק**.
  - בשאלות רב-ברירה – הקפידו להקיף את התשובות בבירור. סימונים לא ברורים יביאו לפסילת התשובה.
  - לא יתקבלו ערעורים בנושא.
- במבחן 16 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.**

### מבנה הבחינה:

- **חלק א' [76 נק']:** 4 שאלות פתוחות.
- **חלק ב' [24 נק']:** 4 שאלות סגורות (אמריקאיות) [כל אחת 6 נק'].

**בהצלחה!**

## חלק א' – שאלות פתוחות [76 נק']

## שאלה 1: Feature selection and Classification models [20 נק']

נתון סט אימון עם  $d \geq 3$  מאפיינים (features) ועם תיוגים בינאריים, משמע  $x_i \in \mathbb{R}^d, y_i \in \{-1, 1\}$ . בסט האימון אין שתי דוגמאות זהות.

נתון שהקורלציה בין שני מאפיינים מסוימים  $A, B$  בסט האימון היא בדיוק  $-1$ .

משמע, לכל דוגמה  $x_i$  מתקיים  $x_i[B] = ax_i[A] + b$  עבור  $a < 0$  ו- $b \in \mathbb{R}$ .

לומדים שני מסווגים:

- בשלב הראשון: לומדים מסווג על סט האימון המקורי ומחשבים עליו את דיוק האימון.
- בשלב השני:
  - מסירים את המאפיין  $B$  ומקבלים סט אימון מעודכן (שבו  $d - 1$  פיצ'רים).
  - מאמנים מסווג חדש על סט האימון המעודכן, ומחשבים עליו את דיוק האימון המעודכן.

עבור כל אלגוריתם למידה, סמנו האם דיוק האימון של המסווג החדש על סט האימון המעודכן זהה בהכרח לזה של המסווג המקורי על סט האימון המקורי.

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

הניחו שאין צעדים סטוכסטיים (אקראיים) בריצת האלגוריתמים ואין שגיאות נומריות (בפרט, בעיות קמורות מתכנסות לפיתרון האנליטי שלהן במדויק).

א.  $kNN$  עם  $k = 1$  (דוגמה לא נחשבת שכנה של עצמה). דיוק האימון זהה בהכרח? כן / לא

הסבר: נניח למשל שהסקאלה של פיצ'ר  $B$  הרבה יותר גדולה מכל האחרות (במקרה קיצוני:  $a \rightarrow -\infty$ ). הפיצ'ר הזה "משתלט" על המרחקים האוקלידיים. לאחר הסרתו נקבל משקלים שונים לגמרי ושכנויות שונות.

ב.  $Hard-SVM$  לינארי לא הומוגני בהנחה שהדאטה המקורי פריד. דיוק האימון זהה בהכרח? כן / לא

הסבר: הדאטה המקורי פריד לינארית. הראינו בתרגול שפיצ'רים ת"ל לא מוסיפים capacity למודלים לינאריים. לכן הדאטה לאחר ההסרה גם הוא פריד לינארית בהכרח. בשני המקרים יתקבל דיוק 100%.

ג.  $Soft-SVM$  לינארי לא הומוגני עם  $\lambda = 1$  בהנחה שהדאטה המקורי פריד. דיוק האימון זהה בהכרח? כן / לא

הסבר: אמנם מבחינת Hinge ניתן להגיע במרחב החדש לאותו loss כמו של המפריד האופטימלי המקורי, אבל רכיב הרגולריזציה עשוי להשתנות. לכן ייתכן שהפיתרון האופטימלי ישתנה ונקבל מפריד שונה ודיוק שונה.

ד.  $ID3$  המשתמש באנטרופיה ובונה עץ בעומק מירבי 4. דיוק האימון זהה בהכרח? כן / לא

הסבר: בכל שלב בריצת  $ID3$ , כאשר נבחר פיצול לפי סף של  $B$  במרחב המקורי, ניתן לבחור פיצול לפי סף שקול של  $A$  במרחב החדש. לאחר הפיצול יתקבלו צמתים זהים ולכן ה- $IG$  של שני הפיצולים המדוברים זהה ומקסימלי (לפי אופן פעולת  $ID3$  במרחב המקורי). לכן בהכרח  $ID3$  יבחר את הפיצול השקול לפי סף של  $A$  ויתקבל עץ זהה ודיוק זהה.

## שאלה 2: Kernel SVM [18 נק']

תזכורת: פונקציה מהווה קרנל חוקי אם ניתן לכתוב אותה בתור  $K(u, v) = \langle \phi(u), \phi(v) \rangle = \sum_{i=1}^n \phi_i(u) \phi_i(v)$  עבור מיפוי כלשהו  $\phi: \mathcal{X} \rightarrow \mathbb{R}^n$ .

נתונות שתי פונקציות קרנל חוקיות  $K_1, K_2: (\mathcal{X} \times \mathcal{X}) \rightarrow \mathbb{R}$  והמיפויים המתאימים להן  $\phi, \psi: \mathcal{X} \rightarrow \mathbb{R}^n$ . משמע, מתקיים:

$$K_1(u, v) = \langle \phi(u), \phi(v) \rangle, \quad K_2(u, v) = \langle \psi(u), \psi(v) \rangle$$

הוכיחו שהפונקציות הבאות מהוות קרנלים חוקיים (כפי שמוסבר בתזכורת לעיל).

בכל סעיף, הגדירו בבירור פונקציית מיפוי  $\phi': \mathcal{X} \rightarrow \mathbb{R}^{n'}$  מתאימה והראו שמתקיים  $K'(u, v) = \langle \phi'(u), \phi'(v) \rangle$  כנדרש.

א. [5 נק'] הפונקציה  $K'(u, v) = (K_1(u, v))^2$ .

הוכחה:

נגדיר  $\phi'$  עבורו  $n' = n^2$ . באופן אינטואיטיבי נתאים כניסה לכל צירוף  $i, j \in [n]$ :

$$\phi'_{ij}(u) = \phi_i(u) \phi_j(u)$$

צריך להראות שמתקיים  $K'(u, v) = \langle \phi'(u), \phi'(v) \rangle = \sum_{ij} \phi'_{ij}(u) \phi'_{ij}(v) = (K_1(u, v))^2$ .

ב. [5 נק'] הפונקציה  $K'(u, v) = K_1(u, v) + 3 \cdot K_2(u, v) + 1$ .

הוכחה:

נגדיר מיפוי בעזרת שרשור וקטורים:  $\phi'(u) = \begin{bmatrix} \phi(u) \\ \sqrt{3}\psi(u) \\ 1 \end{bmatrix}$ . צריך להוכיח את השוויון המבוקש.

ג. [8 נק'] אילו מהפעולות הבאות אמורות להפחית overfitting ב-Kernel SVM?

סמנו את כל התשובות המתאימות (השאלה אינה עוסקת במקרי קצה אלא במקרה הסביר).

a. לעבור מפונקציית הקרנל  $K(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^T \mathbf{v})^p$  לפונקציית הקרנל  $K(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^T \mathbf{v} + 1)^p$ , כאשר  $p \in \mathbb{N}_{\geq 2}$ .

b. להקטין את השונות  $\sigma^2$  של קרנל RBF (שמוגדר בתור  $K(\mathbf{u}, \mathbf{v}) = \exp\left\{-\frac{1}{2\sigma^2} \|\mathbf{u} - \mathbf{v}\|_2^2\right\}$ ).

c. לפתור (במדויק) את הבעיה ה-primal במקום לפתור (במדויק) את הבעיה ה-dual.

d. להגדיל את מקדם הרגולריזציה  $\lambda$  (בהתאמה: להקטין את  $C$ ).

e. להגדיל את סט האימון (באופן *i. i. d.* מאותה התפלגות של הדאטה המקורי).

f. להגדיל את סט המבחן (באופן *i. i. d.* מאותה התפלגות של הדאטה המקורי).

הערות בדיקה: בסעיף ג' ירדו 2 נקודות על כל טענה מיותרת שסומנה או טענה נכונה שחסרה.

## שאלה 3: VC-dimension [16 נק']

נגדיר את מחלקת ההיפותזות של מסווגים ליניאריים הומוגניים ב- $d$  ממדים:

$$\mathcal{X} = \mathbb{R}^d, \mathcal{H}_{\text{lin}}^d = \{x \mapsto \text{sign}(\mathbf{w}^\top x) : \mathbf{w} \in \mathbb{R}^d\}$$

כדי למנוע טעויות, נקבע שלאורך כל השאלה  $\text{sign}(0) = 0$  (ואחרת הפונקציה מחזירה  $\pm 1$  בהתאם לסימן).

א. [7 נק'] הוכיחו שמתקיים  $\text{VCdim}(\mathcal{H}_{\text{lin}}^d) \geq d$ .

הוכחה:

כמו בתרגול.

ב. [9 נק'] הוכיחו שממד ה-VC הוא בדיוק  $d$  על ידי כך שתוכיחו שמתקיים  $\text{VCdim}(\mathcal{H}_{\text{lin}}^d) < d + 1$ .

רמז: כל אוסף  $x_1, \dots, x_d, x_{d+1}$  של  $d + 1$  וקטורים כלשהם ב- $\mathbb{R}^d$  הינו תלוי ליניארית, ובהכרח אחד הווקטורים באוסף (בה"כ  $x_{d+1}$ ) מקיים  $x_{d+1} = \sum_{i=1}^d z_i x_i$  עבור אוסף סקלרים  $z_1, \dots, z_d \in \mathbb{R}$  שלפחות אחד מהם שונה מאפס.

הוכחה:

יהי אוסף  $x_1, \dots, x_d, x_{d+1}$  של  $d + 1$  וקטורים כלשהם ב- $\mathbb{R}^d$ . לפי הרמז מתקיים  $x_{d+1} = \sum_{i=1}^d z_i x_i$ .

$$y_i = \begin{cases} \text{sign}(z_i), & i \leq d, z_i \neq 0 \\ +1, & i \leq d, z_i = 0 \\ -1, & i = d + 1 \end{cases}$$

נגדיר תיוג

נניח בשלילה שקיים  $\mathbf{w} \in \mathbb{R}^{d+1}$  שמתייג נכונה את כל הנקודות. משמע,  $\forall i: \text{sign}(\mathbf{w}^\top x_i) = y_i$ .

לכן הפרדיקציה עבור  $x_{d+1}$  תהיה (נשתמש בכך שיש לפחות מקדם  $z_i$  שונה מאפס):

$$\text{sign}(\mathbf{w}^\top x_{d+1}) = \text{sign}\left(\mathbf{w}^\top \sum_{i=1}^d z_i x_i\right) = \text{sign}\left(\sum_{i=1}^d \underbrace{z_i \mathbf{w}^\top x_i}_{\geq 0}\right) = 1 \neq -1 = y_{d+1}$$

כי  $\text{sign}(\mathbf{w}^\top x_i) = y_i = \text{sign}(z_i)$

קיבלנו סתירה. קיים תיוג שלא ניתן ללמוד ולכן  $\text{VCdim}(\mathcal{H}_{\text{lin}}^d) < d + 1$ .

## שאלה 4: Bagging and Boosting [22 נק']

משמאל מופיע אלגוריתם AdaBoost.

Initialize  $D^{(1)} = \left(\frac{1}{m}, \dots, \frac{1}{m}\right)$ For  $t=1, \dots, T$ :

$$h_t = \mathcal{A}(S, D^{(t)})$$

$$\epsilon_t = \sum_i D_i^{(t)} \cdot \mathbf{1}_{h_t(x_i) \neq y_i}$$

$$\alpha_t = \frac{1}{2} \log\left(\frac{1}{\epsilon_t} - 1\right)$$

$$D_i^{(t+1)} = \frac{1}{Z_t} D_i^{(t)} \exp(-\alpha_t y_i h_t(x_i))$$

$$h_s(x) = \text{sign}\left(\sum_{t=1}^T \alpha_t h_t(x)\right)$$

א. [5 נק'] שימו לב שבכל איטרציה האלגוריתם מעדכן את ההתפלגות באופן  $D_i^{(t+1)} = \frac{1}{Z_t} D_i^{(t)} \exp(-\alpha_t y_i h_t(x_i))$

כאשר  $Z_t \triangleq \sum_{j=1}^m D_j^{(t)} \exp(-\alpha_t y_j h_t(x_j))$  הוא גורם נירמול.

הוכיחו שמתקיים  $D_i^{(t+1)} = c \cdot \exp(-y_i \sum_{k=1}^t \alpha_k h_k(x_i))$  עבור קבוע כלשהו  $c > 0$  (שתלוי ב- $t$ ).

הוכחה:

כמו בתרגול.

בסעיף הבא נוכיח שהבחירה של האלגוריתם ב- $\alpha_t = \frac{1}{2} \log\left(\frac{1}{\epsilon_t} - 1\right)$  הינה אופטימלית.

תזכורת: הגדרנו בתרגול את ההיפוטזה ה-unthresholded "חזקה" באיטרציה  $t$  בתור:

$$h_s^{(t)} = \sum_{k=1}^t \alpha_k h_k(x_i)$$

והראינו שקיים קבוע חיובי  $C$  כך שערך ה-exp. loss של ההיפוטזה הוא:

$$\mathcal{L}_{\exp}(h_s^{(t)}) \triangleq \frac{1}{m} \sum_{i=1}^m \exp\{-y_i \sum_{k=1}^t \alpha_k h_k(x_i)\} = \underbrace{C}_{>0} \cdot \left(e^{-\alpha_t} + (e^{\alpha_t} - e^{-\alpha_t}) \sum_i D_i^{(t)} \mathbf{1}_{h_t(x_i) \neq y_i}\right)$$

כאשר  $C$  אינו תלוי בבחירה של המקדם  $\alpha_t$  שנבחר בזמן  $t$ .

ב. [5 נק'] הוכיחו שבהינתן  $h_t$  ו- $\{D_i^{(t)}\}_{i \in [m]}$  הבחירה ב- $\alpha_t = \frac{1}{2} \log\left(\frac{1}{\epsilon_t} - 1\right)$  ממזער את  $\mathcal{L}_{\exp}(h_s^{(t)})$ .

הערה: שימו לב ש- $\epsilon_t \triangleq \sum_i D_i^{(t)} \cdot \mathbf{1}_{h_t(x_i) \neq y_i}$  והוא קבוע בהינתן  $h_t$  ו- $\{D_i^{(t)}\}_{i \in [m]}$ .

הוכחה:

כמו בתרגול.

הסעיף הבא לא קשור לאלגוריתם AdaBoost אלא לשיטת Ensemble שממשקלת היפותזות באופן כללי.

ג. [12 נק'] נתון סט אימון  $\{(x_i, y_i)\}_{i=1}^m$  עם סיווגים בינאריים  $(y_i \in \{-1, +1\})$ .

צוות של 50 מומחים ומומחיות הגדיר ידנית 50 היפותזות בינאריות  $h_1, \dots, h_{50}: \mathcal{X} \rightarrow \{-1, +1\}$ .

המטרה שלנו היא למשקל אותן, משמע, ליצור היפותזה  $h_\alpha(x) = \text{sign}(\sum_{k=1}^{50} \alpha_k h_k(x))$ .

כדי למנוע טעויות, נקבע שלאורך כל השאלה  $\text{sign}(0) = 0$  (ואחרת הפונקציה מחזירה  $\pm 1$  בהתאם לסימן).

נתון: קיימים משקלים  $\alpha_1^*, \dots, \alpha_{50}^* \in \mathbb{R}$  כך שהשגיאה האמפירית (שגיאת האימון) היא אפס.

כעת, נרצה ללמוד אוסף משקלים  $\alpha_1, \dots, \alpha_{50} \in \mathbb{R}$ .

(i) נסחו בעיית אופטימיזציה מתאימה וקמורה (ביחס למשקלים הנלמדים  $\alpha_1, \dots, \alpha_{50}$ ) כך שהפיתרון שלה ממזער את השגיאה האמפירית.

תשובה (לרשותכם טיוטה בסוף הגיליון):

נסתכל על  $\alpha = [\alpha_1, \dots, \alpha_{50}]^T \in \mathbb{R}^{50}$  ועל  $\mathbf{h}(x) = [h_1(x), \dots, h_{50}(x)]^T \in \{\pm 1\}^{50}$  כוקטורים.

$$\operatorname{argmin}_{\alpha \in \mathbb{R}^{50}} \sum_{i=1}^m \max \left\{ 0, 1 - y_i \sum_{k=1}^{50} \alpha_k h_k(x_i) \right\} = \operatorname{argmin}_{\alpha \in \mathbb{R}^{50}} \sum_{i=1}^m \max \{ 0, 1 - y_i \alpha^T \mathbf{h}(x_i) \}$$

(ii) נמקו בקצרה מדוע הבעיה שהגדרתם קמורה (מומלץ להשתמש בטענות מההרצאה ומהתרגול).

תשובה (לרשותכם טיוטה בסוף הגיליון):

נשים לב שקיבלנו בדיוק בעיית Soft-SVM בלי רגולריזציה במרחב החדש.

כפי שהראינו בתרגיל הבית – הפונקציה  $1 - y_i \alpha^T \mathbf{h}(x_i)$  לינארית ב- $\alpha$  ולכן קמורה ביחס אליו.

מקסימום בין שתי פונקציות קמורות הוא קמור.

ולסיום – סכום של פונקציות קמורות הוא קמור.

(iii) יהי אוסף המשקלים  $\alpha_1, \dots, \alpha_{50} \in \mathbb{R}$  פיתרון (אופטימלי) של הבעיה הקמורה שהגדרתם לעיל. הוכיחו שהשגיאה האמפירית של פיתרון זה היא מינימלית.

הוכחה:

לפי הנתון – קיים  $\alpha^*$  עבורו לכל דוגמה מתקיים  $\text{sign}(\alpha^{*\top} \mathbf{h}(x_i)) = y_i \in \{\pm 1\}$ .

ולכן  $y_i \alpha^{*\top} \mathbf{h}(x_i) > 0$  לכל דוגמה.

מכאן מובטח שקיים פיתרון  $\beta \alpha^*$  עבור סקאלר כלשהו  $\beta > 0$  גדול מספיק,

כך שמתקיים  $\forall i: y_i (\beta \alpha^*)^\top \mathbf{h}(x_i) > 1$  ושהוא מאפס את ה-objective שהגדרנו.

מובטח שיימצא פיתרון שמאפס את ה-objective (כי קיימים פתרונות כאלה וה-objective חיובי).

בגלל שה-Hinge loss הוא surrogate עבור השגיאה האמפירית, גם השגיאה האמפירית תתאפס

עבור כל פיתרון אופטימלי של הבעיה שהוגדרה.



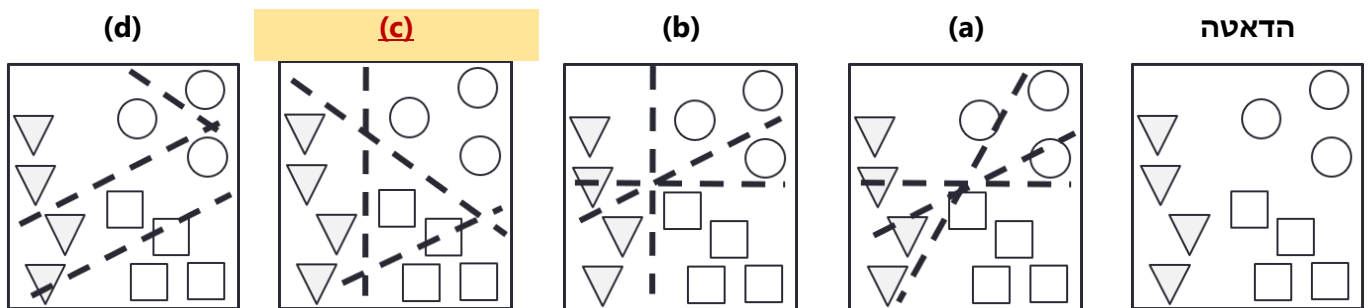
## חלק ב' – שאלות אמריקאיות [24 נק']

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

**הערות בדיקה:** בסעיפים ב'-ד' ירדו 2 נקודות על כל טענה מיותרת שסומנה או טענה נכונה שחסרה.

א. [6 נק'] נתון דאטה דו-ממדי עם שלושה תיוגים אפשריים – משולש, מעגל וריבוע.

מפעילים One-vs-All (OVA) עם פרספטרון בתור מסווג בסיס. לכל בעיה בינארית מריצים את הפרספטרון מספיק איטרציות כך שאם הבעיה פרידה ליניארית – הפרספטרון יתכנס (חישבו מה קורה אחרת).  
בכל תרשים מצוירים שלושה גבולות החלטה שאמורים לתאר את הגבולות שהתקבלו ע"י מסווגי הבסיס (הבינאריים).  
**הקיפו** את האות שמתאימה לתרשים היחיד שמתאר גבולות החלטה שיכולים להתקבל ע"י השיטה שתוארה לעיל.



ב. [6 נק'] נתון דאטה  $d$ -ממדי  $\{(x_i, y_i)\}_{i=1}^m$  עם סיווגים בינאריים  $(\pm 1)$ .

סמנו את כל הטענות הנכונות **בהכרח** ביחס **לשיטות להורדת ממד** לממד נמוך כלשהו.

הבהרה: זו שאלת חשיבה ולא שאלת טריוויה.

a. הדאטה המקורי פריד ליניארית  $\Leftrightarrow$  הדאטה פריד ליניארית לאחר הורדת ממד ליניארית עם PCA.

b. הדאטה המקורי פריד ליניארית  $\Leftrightarrow$  בהסתברות גבוהה, הדאטה פריד ליניארית לאחר הורדת ממד ליניארית עם

Random Projections.

~~c. הדאטה המקורי פריד ליניארית  $\Leftrightarrow$  הדאטה פריד ליניארית לאחר הורדת ממד עם Kernel-PCA (עם קרנל לא ליניארי).~~

d. הדאטה פריד ליניארית לאחר הורדת ממד ליניארית עם PCA  $\Leftrightarrow$  הדאטה המקורי פריד ליניארית.

e. PCA היא שיטה Supervised בעוד ש-Random Projections היא שיטה Unsupervised.

**הערות בדיקה:** על אף שטענה c לא נכונה, הוחלט שלא להוריד ניקוד על סימון שלה בגלל שלא התעכבנו על זה בהרצאות.

ג. [6 נק'] נתונה בעיית רגרסיה לינארית (עבור  $\lambda \geq 0, p \in \{0.5, 1, 2\}$ ):

$$\operatorname{argmin}_{\mathbf{w}} \left( \frac{1}{m} \sum_{i=1}^m (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|_p^p \right) = \operatorname{argmin}_{\mathbf{w}} \left( \frac{1}{m} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_p^p \right)$$

סמנו את כל הטענות הנכונות.

הערה: בכל הסעיפים, שואלים על קמירות ביחס ל- $\mathbf{w}$ .

$$\|\mathbf{w}\|_p^p \triangleq \sum_{k=1}^d |w_k|^p$$

a. כאשר  $\lambda = 0$ : ההסיאן הוא  $\frac{2}{m} \mathbf{X}^\top \mathbf{X}$  שהיא מטריצה PSD אבל לא בהכרח PD, ולכן הבעיה לא בהכרח קמורה.

b. כאשר  $\lambda > 0$  וגם  $p = 2$ : הבעיה קמורה בהכרח.

c. כאשר  $\lambda > 0$  וגם  $p = 1$ : הבעיה קמורה בהכרח.

d. כאשר  $\lambda > 0$  וגם  $p = 0.5$ : הבעיה קמורה בהכרח.

e. כאשר  $\lambda > 0$  וגם  $p = 1$ : הבעיה אינה גזירה ביחס ל- $\mathbf{w}$  ולכן לא ניתן להפעיל שיטות גרדיינט (או סאב-גרדיינט).

ד. [6 נק'] רוצים לאמן רשת נוירונים בתצורת Autoencoder.

נסמן את אוסף כל הפרמטרים של הרשת בתור  $\theta$ . לכל דוגמה  $\mathbf{x}_i \in \mathbb{R}^d$  נסמן את הפלט של הרשת בתור  $\hat{\mathbf{x}}_i \in \mathbb{R}^d$ .

בהרצאה ראינו שניתן לאמן את הרשת ע"י פיתרון הבעיה הבאה:  $\operatorname{argmin}_{\theta} \left( \lambda \|\theta\|_2^2 + \frac{1}{m} \sum_{i=1}^m \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|_2^2 \right)$ .

אילו מבין הבעיות הבאות עשויות גם כן להתאים לאימון רשת Autoencoder (עבור  $\lambda > 0$ )?

סמנו את כל התשובות המתאימות.

$$\operatorname{argmin}_{\theta} \left( \lambda \|\theta\|_2^2 + \frac{1}{m} \sum_{i=1}^m \|\mathbf{x}_i + \hat{\mathbf{x}}_i\|_2^2 \right) \quad a.$$

$$\operatorname{argmin}_{\theta} \left( \lambda \|\theta\|_2^2 + \frac{1}{m} \sum_{i=1}^m \exp\{\|\mathbf{x}_i - \hat{\mathbf{x}}_i\|_2^2\} \right) \quad b.$$

$$\operatorname{argmin}_{\theta} \left( \lambda \|\theta\|_2^2 - \frac{1}{m} \sum_{i=1}^m \|\mathbf{x}_i\|_2^2 + \frac{1}{m} \sum_{i=1}^m \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|_2^2 \right) \quad c.$$

$$\operatorname{argmin}_{\theta} \left( \lambda \|\theta\|_2^2 + \frac{1}{m} \sum_{i=1}^m \max\{0, 1 - \operatorname{sign}(\mathbf{x}_i^\top \hat{\mathbf{x}}_i)\} \right) \quad d.$$

$$\operatorname{argmin}_{\theta} \left( -\lambda \|\theta\|_2^2 + \frac{1}{m} \sum_{i=1}^m \frac{\mathbf{x}_i^\top \hat{\mathbf{x}}_i}{\|\mathbf{x}_i\|_2 \|\hat{\mathbf{x}}_i\|_2} \right) \quad e.$$