



מבוא למערכות לומדות (236756)
סמסטר אביב תשפ"ג – 16 ביולי 2023
מרצה: ד"ר ניר רחנפלד

מבחן מסכם מועד א' – פתרון חלקי

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

בהצלחה!

שאלה 1: Feature Selection [8 נק']

סמנו את התשובות המתאימות (לפי ההוראות). אין צורך לכתוב הסברים. סימונים לא ברורים יובילו לפסילת התשובה.
נתון ה-dataset הבא.

sample/feat.	$x[1]$	$x[2]$	$x[3]$	y
#1	1	-1	1	-1
#2	-1	1	1	-1
#3	1	-1	1	-1
#4	-1	1	1	-1
#5	1	1	-1	+1
#6	1	1	1	+1

נפעיל sequential feature selection (כפי שהגדרנו בתרגיל הבית) בעזרת אלגוריתם בסיס \mathcal{A} שפותר ERM ומוצא מסווג ליניארי לא-הומוגני עם שגיאת 0-1 מינימלית על הנתונים. כאן – נמצא שני פיצ'רים.

אילו פיצ'רים ייבחרו בתהליך **Backward Selection**?

אילו פיצ'רים ייבחרו בתהליך **Forward Selection**?

a. $x[1], x[2]$

a. $x[1], x[2]$

b. $x[1], x[3]$

b. $x[1], x[3]$

c. $x[2], x[3]$

c. $x[2], x[3]$

d. מהנתונים הקיימים, ניתן לענות רק לגבי פיצ'ר יחיד.

d. מהנתונים הקיימים, ניתן לענות רק לגבי פיצ'ר יחיד.

e. מהנתונים הקיימים, לא ניתן לענות כלל.

e. מהנתונים הקיימים, לא ניתן לענות כלל.

שאלה 2: Capacity של מחלקות שונות [21 נק']

נתונה התפלגות \mathcal{D} על דאטה עם $d \geq 10$ פיצ'רים בינאריים ($\mathcal{X} = \{0,1\}^d$) ותיוגים בינאריים ($\mathcal{Y} = \{-1, +1\}$). ידוע שההתפלגות \mathcal{D} נותנת לכל דוגמה אפשרית $\mathbf{x} \in \mathcal{X}$ הסתברות כלשהי סופית גדולה ממש מאפס. ידוע שהתיוג האמיתי של דוגמה כלשהי הוא חיובי אם ורק אם בדוגמה יש לפחות שני פיצ'רים כלשהם "דולקים" (שערכם 1). לדוגמה: כאשר $d = 5$, לדוגמאות $(0,1,1,0,1)$ ו- $(1,1,0,0,0)$ יש בהכרח תיוג חיובי, ולדוגמה $(0,0,1,0,0)$ בהכרח תיוג שלילי. עבור כל אחת מהמחלקות הבאות, נבדוק האם ניתן ליצור בעזרתה מסווג שיגיע לשגיאת הכללה אפס. שימו לב: זו שאלה על מחלקות היפותזות, ולא על אלגוריתמי למידה ספציפיים.

א. האם קיים עץ החלטה בינארי (צומת שאינו עלה יכול להתפצל רק לשני צמתים) שמגיע לשגיאת הכללה אפס? סמנו: כן / לא

אם קיים – מהו עומק העץ המינימלי הנדרש (בסדר גודל, למשל $\theta(2^d), \theta(d), \theta(\ln d), \theta(1)$ וכו')? הסבירו בקצרה. אחרת – הסבירו באופן מפורט למה לא.

תשובה (לרשותכם דפי טיוטה בסוף הגיליון):

נדרש עץ עם $\theta(d)$ רמות.

בשורש שואלים לגבי פיצ'ר 1. אם דלוק – ברמה הבאה שואלים אם פיצ'ר 2 דלוק. אם כן – חוזים +1.

אחרת, שואלים אם פיצ'ר 3 דלוק וכן הלאה.

נדרשות לכל היותר d שאלות כדי להבין אם התחזית היא חיובית או שלילית.

$$\text{sign}(\mathbf{w}^T \mathbf{x}) = \begin{cases} -1, & \mathbf{w}^T \mathbf{x} < 0 \\ 0, & \mathbf{w}^T \mathbf{x} = 0 \\ +1, & \mathbf{w}^T \mathbf{x} > 0 \end{cases}$$

משמע, margin של אפס בהכרח גורר שגיאת סיווג.

סמנו: **כן** / לא

ב. האם קיים מפרד ליניארי לא הומוגני שמגיע לשגיאת הכללה אפס?

אם קיים – הציעו וקטור $\mathbf{w} \in \mathbb{R}^d$ וסקלר $b \in \mathbb{R}$ שמגיעים לשגיאה אפס (במקרה כזה לא נדרש הסבר נוסף).

אחרת – הסבירו באופן מפורט למה לא.

תשובה:

$$\mathbf{w} = [1, 1, \dots, 1], \quad b = -1.1$$

ג. האם ניתן ליצור מסווג 1-Nearest-Neighbor (עם מטריקת ℓ_2 האוקלידית) שמגיע לשגיאת הכללה אפס?

סמנו: **כן** / לא

(שימו לב, "ליצור מסווג 1NN" דורש למעשה ליצור אוסף מתאים של דוגמאות אימון.)

אם ניתן – מה הגודל המינימלי של סט אימון שנדרש ליצירת מסווג כזה (בסדר גודל, למשל $\Theta(\ln d)$, $\Theta(1)$ וכו')? הסבירו בקצרה.

אחרת – הסבירו באופן מפורט למה לא.

תשובה:

צריך את כל הווקטורים הסטנדרטיים $\{\mathbf{e}_i\}_{i \in [d]}$ שמתוייגים שלילית.

אח"כ צריך את כל הדוגמאות שבהן בדיוק 2 פיצ'רים דולקים (מתוייגות חיובית). יש $\binom{d}{2}$ כאלה.

כל דוגמה עם 2 פיצ'רים דולקים ויותר תהיה יותר קרובה לדוגמה חיובית עם 2 פיצ'רים דולקים מאשר לדוגמאות השליליות. כל דוגמה עם פיצ'ר דולק ומטה, תהיה יותר קרובה לדוגמה שלילית.

בסה"כ צריך $\Theta(d^2)$ דוגמאות באוסף.

הערה: זה לא משנה אם דוגמה "נחשבת שכנה של עצמה" או לא. השאלה עוסקת בשגיאת הכללה ואילו העניין שנק' היא שכנה של עצמה או לא רלוונטי רק ל-training set.

שאלה 3: AdaBoost ופרספטרון [21 נק']

משמאל מופיע אלגוריתם AdaBoost.

מימין סט אימון עם $m = 30$ דוגמאות בדו-ממד ($\mathcal{X} = \mathbb{R}^2$) עם תיוגים בינאריים $+1$ מסומן ב- \blacktriangle ו- -1 מסומן ב- \blacklozenge .Initialize $D^{(1)} = \left(\frac{1}{m}, \dots, \frac{1}{m}\right)$ For $t = 1, \dots, T$:

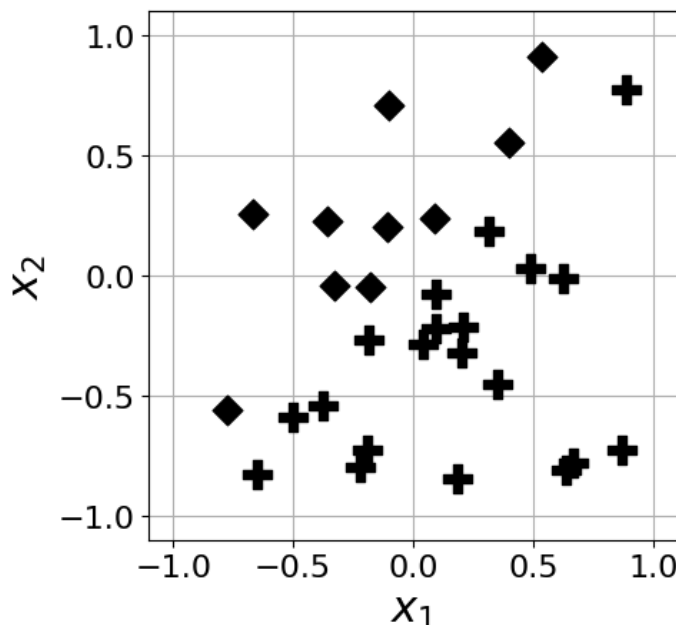
$$h_t = \mathcal{A}(S, D^{(t)})$$

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

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

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

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



א. [5 נק'] כאשר ניתן למקבל הרצות נפרדות של \mathcal{A} (עד T הרצות במקביל), מה פקטור השיפור שניתן לקבל בזמן ההרצה של AdaBoost? הסבירו בקצרה את תשובתכם.

תשובה: המיקבול לא יעזור.

זה אלגוריתם סדרתי. בכל איטרציה צריך להשתמש בהתפלגות שחושבה באיטרציה הקודמת.

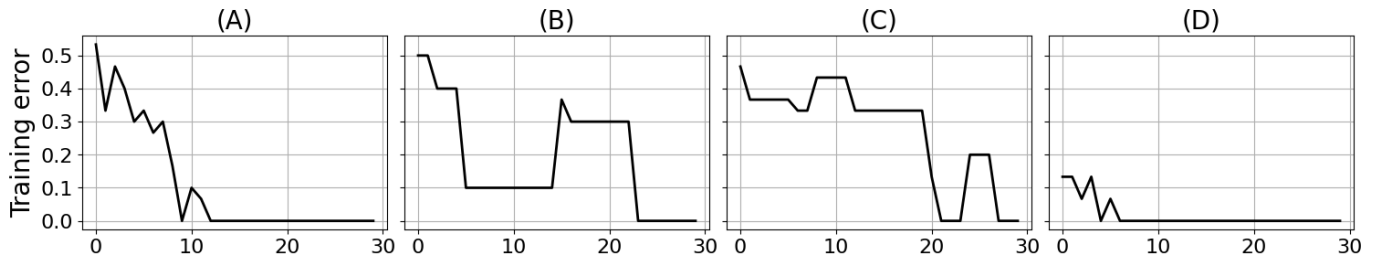
ב. [5 נק'] מצאו $\mathbf{w} \in \mathbb{R}^2$ שיוצר מפריד ליניארי הומוגני עם שגיאת אימון 0 על הדאטה (משמע, $\text{sign}(\mathbf{w}^T \mathbf{x}_i) = y_i$).

$$\mathbf{w} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

כתבו את שני הערכים במסגרת:

(ניתן להבין את התשובה על סמך התרשים; בבדיקת התשובה יינתן מרווח שגיאה סביר בערכים המספריים עצמם.)

ג. [11 נק'] לפניכם ארבע עקומות התכנסות שונות של שגיאת ה-0-1 על הנתונים שבתרשים הקודם.



צירי y בתרשימים מראים את שגיאת האימון על כל 30 דוגמאות האימון.

נתאים **שניים** מהתרשימים לאלגוריתמי הלמידה הבאים:

i. אלגוריתם פרספטרון עם גודל צעד $\eta = 1$ ואתחול $\mathbf{w}_0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$.

ציר x בתרשים המתאים מראה 30 איטרציות של האלגוריתם (בכל איטרציה האלגוריתם רואה דוגמה בודדת).

הקיפו את האות של התרשים ששייך לאלגוריתם זה. (A) (B) (C) (D)

הסבירו בקצרה את תשובתכם.

הסבר: ברגע שפרספטרון מגיע לשגיאת אימון אפס הוא לא מעדכן יותר את הפיתרון.

כמו כן אנחנו מצפים שיהיו הרבה איטרציות שבהן אין עדכון (ויתקבלו איזורים שטוחים בעקומה).

ii. אלגוריתם AdaBoost כאשר \mathcal{A} הוא אלגוריתם ERM ללמידת Decision Stump (עץ החלטה בעומק 1) ביחס לדוגמאות ולמשקלים הנתונים.

ציר x בתרשים המתאים מראה 30 איטרציות של האלגוריתם (כפי שמופיע בתחילת השאלה).

הקיפו את האות של התרשים ששייך לאלגוריתם זה. (A) (B) (C) (D)

הסבירו בקצרה את תשובתכם.

הסבר: רואים על הדאטה שבאיטרציה הראשונה ייבחר dec. stump שעושה 4-5 טעויות מתוך 30.

שאלה 4: מסווגים ליניאריים ואופטימיזציה [25 נק']

בכל השאלה נתון מדגם $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ (כאשר $\mathcal{X} = \mathbb{R}^d$ ו- $\mathcal{Y} = \{-1, +1\}$) ונניח שהוא פריד ליניארית הומוגנית.

נגדיר בעיית למידה קמורה בעזרת ה-exponential loss:

$$\argmin_{\mathbf{w} \in \mathbb{R}^d} \underbrace{\sum_{i=1}^m \exp(-y_i \mathbf{w}^T \mathbf{x}_i)}_{\triangleq \mathcal{L}(\mathbf{w})}$$

א. [4 נק'] השתמשו במטריצת ההסיאן (שהיא $\nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w}) = \sum_{i=1}^m \exp(-y_i \mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^T$) כדי להוכיח שהבעיה קמורה.

הוכחה: כל $\mathbf{x}_i \mathbf{x}_i^T$ היא מטריצת גראם ולכן PSD. המקדמים $\exp(-y_i \mathbf{w}^T \mathbf{x}_i)$ הם בהכרח חיוביים.

מכאן שההסיאן הוא צירוף ליניארי עם מקדמים חיוביים של מטריצות PSD ולכן הוא גם PSD.

ב. [8 נק'] למרות הקמירות, הוכיחו שתחת הנחת הפרידות, לבעיה אין מינימום.

הוכחה: הבעיה פרידה ולכן קיים $\bar{\mathbf{w}} \in \mathbb{R}^d$ כך שלכל $i \in [m]$ מתקיים $y_i \bar{\mathbf{w}}^T \mathbf{x}_i > 0$

לכל $\alpha > 0$ הפיתרון $\alpha \bar{\mathbf{w}}$ גם הוא מקיים $y_i (\alpha \bar{\mathbf{w}})^T \mathbf{x}_i > 0$ לכל i .

משמע, לכל $i \in [m]$, הפונק' $y_i (\alpha \bar{\mathbf{w}})^T \mathbf{x}_i$ חיובית ועולה ב- α ולכן הפונק' $\exp(-y_i \mathbf{w}^T \mathbf{x}_i)$ יורדת ב- α .

בסה"כ בעיית האופטימיזציה היא סכום של פונקציות חיוביות יורדות (ממש) ב- α ולכן אין לה מינימום.

הערה: בניגוד לתרגיל הבית, לא הנחנו כאן שהדוגמאות בת"ל. במצב כזה לא ניתן לטעון בפשטות שהגרדיינט $\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = -\sum_{i=1}^m \exp(-y_i \mathbf{w}^T \mathbf{x}_i) y_i \mathbf{x}_i$ לא מתאפס. שימו לב גם שהגרדיינט לא היה נתון בשלב הזה בשאלה. יכול להיות חיובי או שלילי

טענה (אין צורך להוכיח): מתקיים $\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = -\sum_{i=1}^m \exp(-y_i \mathbf{w}^T \mathbf{x}_i) y_i \mathbf{x}_i$.

נסמן: יהיו $\mathbf{w}(1), \dots, \mathbf{w}(t) \in \mathbb{R}^d$ הפתרונות המתקבלים מהרצת GD (מלא, לא סטוכסטי) על ה-exponential loss החל מהווקטור $\mathbf{w}(0) = \mathbf{0}_d$, עם גודל צעד $\eta > 0$.

ג. [5 נק'] בסעיף זה בלבד נניח שיש רק דוגמה אחת (\mathbf{x}, y) במדגם, שהתיג שלה הוא $y = 1$ ושהיא מנורמלת ($\|\mathbf{x}\| = 1$). כמו כן, נניח שגודל הצעד בהרצת GD הוא $\eta = 1$.

קל לראות שתחת תנאי הסעיף מתקיים: $\mathbf{w}(t) = \underbrace{\|\mathbf{w}(t)\|}_{\in \mathbb{R}} \underbrace{\mathbf{x}}_{\in \mathbb{R}^d}$ (ניתן להשתמש בכך מבלי להוכיח זאת).

הוכיחו שתחת תנאי הסעיף מתקיים: $\|\mathbf{w}(t)\| = \|\mathbf{w}(t-1)\| + e^{-\|\mathbf{w}(t-1)\|}$ (אין צורך להשתמש באינדוקציה).

הוכחה: $\mathbf{w}(t) = \mathbf{w}(t-1) + \exp(-\mathbf{w}_{t-1}^T \mathbf{x}) \mathbf{x} = \|\mathbf{w}(t-1)\| \mathbf{x} - \exp(-\|\mathbf{w}(t-1)\| \mathbf{x}^T \mathbf{x}) \mathbf{x}$

$$= (\|\mathbf{w}(t-1)\| - e^{-\|\mathbf{w}(t-1)\|}) \mathbf{x}$$

ו- \mathbf{x} וקטור מנורמל.

$\{1, -1\}$

$$w_t = w_{t-1} + \sum e^{y_i x_i} = w_{t-1} + e$$

תחת תנאי הסעיף האחרון ובהמשך לנוסחה הרקורסיבית שהוכחתם, ניתן להראות שמתקיים $\|\mathbf{w}(t)\| \approx \ln(t)$.
 כעת נדון בתופעה כללית יותר (לדאטה פריד ליניארית עם $m \geq 1$ דוגמאות).

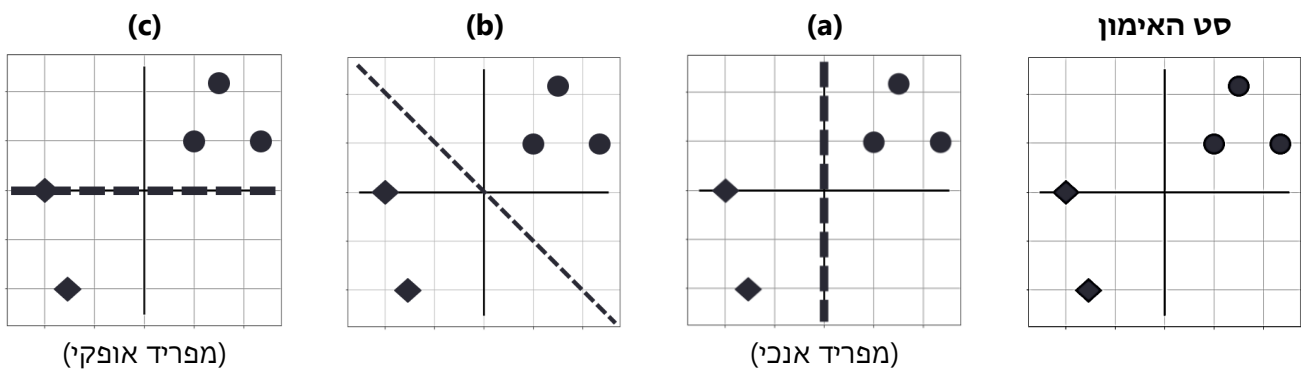
נסמן:
$$\hat{\mathbf{w}} = \operatorname{argmax}_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{\|\mathbf{w}\|} \text{ s.t. } \left(\left(\min_{i \in [m]} |\mathbf{w}^\top \mathbf{x}_i| = 1 \right) \text{ and } (\forall i \in [m]: y_i \mathbf{w}^\top \mathbf{x}_i > 0) \right)$$

משפט: לכל $t \geq 3$, הפתרונות $\mathbf{w}(1), \dots, \mathbf{w}(t) \in \mathbb{R}^d$ שמתקבלים מהרצת GD על ה-exponential loss הם מהצורה:

$$\mathbf{w}(t) = \underbrace{\ln(t)}_{\in \mathbb{R}} \cdot \underbrace{\hat{\mathbf{w}}}_{\in \mathbb{R}^d} + \underbrace{\tilde{\mathbf{r}}(t)}_{\in \mathbb{R}^d}$$

עבור וקטורים $\tilde{\mathbf{r}}(1), \dots, \tilde{\mathbf{r}}(t)$ לא ידועים שמקיימים $\|\tilde{\mathbf{r}}(t)\| = \mathcal{O}(\ln \ln(t))$.

ד. [8 נק'] לפיכך תרשימים של סט אימון פריד ליניארית עם 5 דוגמאות ובנוסף כמה גבולות החלטה (decision boundaries).



(הניחו שאין שגיאות נומריות)

איזה תרשים מתאר את גבול ההחלטה שמתקבל ע"י $\lim_{t \rightarrow \infty} \frac{\mathbf{w}(t)}{\|\mathbf{w}(t)\|}$?

(לא ניתן לדעת) (c) (b) (a)

סמנו באופן ברור והסבירו בקצרה.

הסבר:

כפי שלמדנו $\hat{\mathbf{w}}$ הוא הפיתרון של Hard-SVM (אפשר לראות זאת גם ע"י פעולות פשוטות על הגדרתו).

$$\lim_{t \rightarrow \infty} \frac{\mathbf{w}(t)}{\|\mathbf{w}(t)\|} = \lim_{t \rightarrow \infty} \frac{\ln(t) \hat{\mathbf{w}} + \|\tilde{\mathbf{r}}(t)\| \frac{1}{\|\tilde{\mathbf{r}}(t)\|} \tilde{\mathbf{r}}(t)}{\|\mathbf{w}(t)\|} = \lim_{t \rightarrow \infty} \frac{\ln(t) \hat{\mathbf{w}} + c \ln \ln(t) \frac{1}{\|\tilde{\mathbf{r}}(t)\|} \tilde{\mathbf{r}}(t)}{\left\| \ln(t) \hat{\mathbf{w}} + c \ln \ln(t) \frac{1}{\|\tilde{\mathbf{r}}(t)\|} \tilde{\mathbf{r}}(t) \right\|}$$

$$= \lim_{t \rightarrow \infty} \frac{\hat{\mathbf{w}} + \overbrace{\frac{c \ln \ln(t)}{\ln(t)}}^{\rightarrow 0} \overbrace{\frac{1}{\|\tilde{\mathbf{r}}(t)\|} \tilde{\mathbf{r}}(t)}^{=O(1)}}{\left\| \hat{\mathbf{w}} + \frac{c \ln \ln(t)}{\ln(t)} \frac{1}{\|\tilde{\mathbf{r}}(t)\|} \tilde{\mathbf{r}}(t) \right\|} = \frac{\hat{\mathbf{w}}}{\|\hat{\mathbf{w}}\|}$$

ולכן יתקבל מפריד זהה לזה של ה-Hard-SVM.

שאלה 5: רגרסיה ורגולריזציה [25 נק']

עבור $\mathbf{X} \in \mathbb{R}^{m \times d}$, $\mathbf{y} \in \mathbb{R}^m$ ו- $\lambda > 0$, נסמן את הפתרונות של שלוש בעיות רגרסיה ליניארית שלמדנו:

$$\hat{\mathbf{w}}_{\text{LS}} = \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2, \quad \hat{\mathbf{w}}_{\ell_2} = \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} (\|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2), \quad \hat{\mathbf{w}}_{\ell_1} = \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} (\|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_1)$$

הנחה: בשאלה זו נניח שהעמודות של \mathbf{X} אורתונורמליות (משמע $\mathbf{X}^T \mathbf{X} = \mathbf{I}_{d \times d}$).

תזכורת: הנגזרת של התבנית הריבועית היא $\nabla_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = 2\mathbf{X}^T(\mathbf{X}\mathbf{w} - \mathbf{y})$.

א. [3 נק'] תחת ההנחה (אורתונורמליות), הוכיחו שמתקיים $\hat{\mathbf{w}}_{\text{LS}} = \mathbf{X}^T \mathbf{y}$.

$$\nabla_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = 2\mathbf{X}^T(\mathbf{X}\mathbf{w} - \mathbf{y}) = 2\mathbf{X}^T \mathbf{X}\mathbf{w} - 2\mathbf{X}^T \mathbf{y} = 2\mathbf{w} - 2\mathbf{X}^T \mathbf{y} = \mathbf{0}$$

$$\hat{\mathbf{w}}_{\text{LS}} = \mathbf{X}^T \mathbf{y}$$

ב. [3 נק'] תחת ההנחה, מצאו ביטוי סגור ל- $\hat{\mathbf{w}}_{\ell_2}$ כפונקציה של λ ו- $\hat{\mathbf{w}}_{\text{LS}}$. צרפו לתשובתכם פיתוח קצר מתאים.

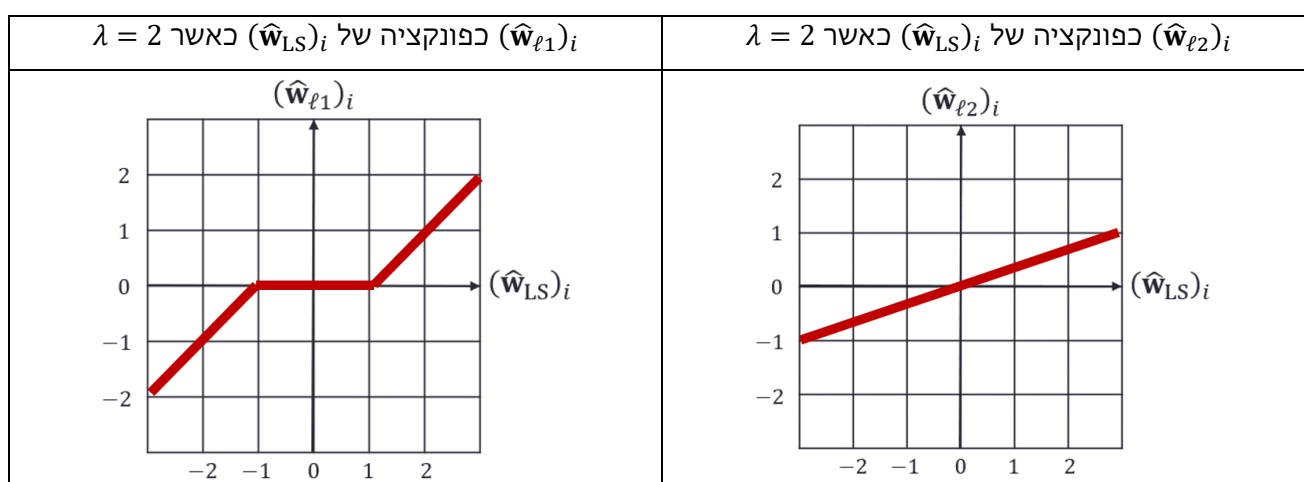
$$\nabla_{\mathbf{w}} (\|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2) = 2\mathbf{X}^T(\mathbf{X}\mathbf{w} - \mathbf{y}) + 2\lambda \mathbf{w} = 2\mathbf{w} - 2\mathbf{X}^T \mathbf{y} + 2\lambda \mathbf{w} = \mathbf{0}$$

$$\hat{\mathbf{w}}_{\ell_2} = \left(\frac{1}{1 + \lambda} \right) \mathbf{X}^T \mathbf{y} = \left(\frac{1}{1 + \lambda} \right) \hat{\mathbf{w}}_{\text{LS}}$$

נתונה טענה 1: תחת ההנחה, מתקיים $(\hat{\mathbf{w}}_{\ell_1})_i = \operatorname{sign}((\hat{\mathbf{w}}_{\text{LS}})_i) \cdot \max\left(0, |(\hat{\mathbf{w}}_{\text{LS}})_i| - \frac{\lambda}{2}\right)$.

ג. [4 נק'] עבור כניסה i שרירותית וערך $\lambda = 2$, השתמשו בביטוי שמצאתם בסעיף הקודם ובטענה 1 וציירו באופן ברור

על גבי התרשימים הבאים את העקומות של $(\hat{\mathbf{w}}_{\ell_2})_i$ ו- $(\hat{\mathbf{w}}_{\ell_1})_i$ כפונקציה של $(\hat{\mathbf{w}}_{\text{LS}})_i$.



ד. [8 נק'] הסבירו כיצד הסעיפים הקודמים ממחישים את התכונות של Ridge regression ו-Lasso והבדלים ביניהם.

Ridge מכווץ את המשקלים ו-LASSO מאפס משקלים שקרובים לאפס.

כעת נוכיח את טענה 1 במקרה חד-ממדי פשוט, משמע: $d = 1$.

כלומר, עבור $\mathbf{X}, \mathbf{y} \in \mathbb{R}^m$ ו- $\lambda > 0$, נגדיר: $\hat{w}_{LS} = \operatorname{argmin}_{w \in \mathbb{R}} \|\mathbf{X}w - \mathbf{y}\|_2^2$, $\hat{w}_{\ell 1} = \operatorname{argmin}_{w \in \mathbb{R}} (\underbrace{\|\mathbf{X}w - \mathbf{y}\|_2^2 + \lambda|w|}_{\triangleq \mathcal{L}(w)})$

נמשיך להניח: מתקיים $\mathbf{X}^T \mathbf{X} = 1$ (ולכן עדיין מתקיים $\hat{w}_{LS} = \mathbf{X}^T \mathbf{y}$).

לשימושכם תזכורת לתתי-נגזרות והכללה של תנאי האופטימליות למקרה של תת-נגזרת:

תזכורת: תהא $f: \mathbb{R} \rightarrow \mathbb{R}$ קמורה. נסמן ב- $\partial f(w)$ את קבוצת תתי-הנגזרות (sub-derivatives) שלה בנקודה $w \in \mathbb{R}$.

למשל, עבור $f(z) = |z|$, מתקיים $\partial f(-2) = \{-1\}$ וגם $\partial f(0) = [-1, 1]$.

תנאי אופטימליות (אין צורך להוכיח): תהא $f: \mathbb{R} \rightarrow \mathbb{R}$ קמורה. $f(w) = \min_{z \in \mathbb{R}} f(z)$ אם ורק אם $0 \in \partial f(w)$.

ה. [8 נק'] הוכיחו שבמקרה החד-ממדי, תחת הנחת האורתונורמליות, מתקיים $\hat{w}_{\ell 1} = \operatorname{sign}(\hat{w}_{LS}) \cdot \max\left(0, |\hat{w}_{LS}| - \frac{\lambda}{2}\right)$.

הוכחה (לרשותכם דפי טיוטה בסוף הגיליון):

$$\mathcal{L}'(w) = 2\mathbf{X}^T(\mathbf{X}w - \mathbf{y}) + \lambda g(w) = 2w - 2\mathbf{X}^T \mathbf{y} + \lambda g(w) = 2w - 2\hat{w}_{LS} + \lambda g(w)$$

כאשר g תת-נגזרת כלשהי $g(w) \in \partial f(w)$ עבור $g(w) = |w|$. נשים לב שמתקיים $\mathcal{L}'(w) \in \partial \mathcal{L}(w)$.

נראה שמתקיים $\mathcal{L}'(\hat{w}_{\ell 1}) = 0$ וסיימנו (לפי תנאי האופטימליות).

מקרה א': מתקיים $\hat{w}_{\ell 1} = |\hat{w}_{LS}| - \frac{\lambda}{2} > 0$. הנגזרת היא בהכרח $g(\hat{w}_{\ell 1}) = \operatorname{sign}(\hat{w}_{LS})$.

$$\mathcal{L}'(\hat{w}_{\ell 1}) = 2\hat{w}_{\ell 1} - 2\hat{w}_{LS} + \lambda g(\hat{w}_{\ell 1}) = 2\left(\operatorname{sign}(\hat{w}_{LS})\left(|\hat{w}_{LS}| - \frac{\lambda}{2}\right)\right) - 2\hat{w}_{LS} + \operatorname{sign}(\hat{w}_{LS})\lambda$$

$$= 2\hat{w}_{LS} - \operatorname{sign}(\hat{w}_{LS})\lambda - 2\hat{w}_{LS} + \operatorname{sign}(\hat{w}_{LS})\lambda = 0$$

מקרה ב': מתקיים $\hat{w}_{\ell 1} = 0 \geq |\hat{w}_{LS}| - \frac{\lambda}{2}$ ולכן $|\hat{w}_{LS}| > 0$ ו- $\lambda \geq 2|\hat{w}_{LS}|$.

נבחר תת-נגזרת $g(0) = \frac{2\hat{w}_{LS}}{\lambda} \in [-1, 1]$.

$$\mathcal{L}'(\hat{w}_{\ell 1}) = 2\hat{w}_{\ell 1} - 2\hat{w}_{LS} + \lambda g(\hat{w}_{\ell 1}) = 0 - 2\hat{w}_{LS} + \lambda \left(\frac{2\hat{w}_{LS}}{\lambda}\right) = 0$$