



פיתרון מועד א' – טור 1

פיתרון שאלה 1 [10 נק']

ב. [9 נק'] נפתח את הפונקציה שיש להביא למינימום תוך השמטת גורמים שאינם תלויים בסיווג ושימוש בתכונות נוספות של argmax .

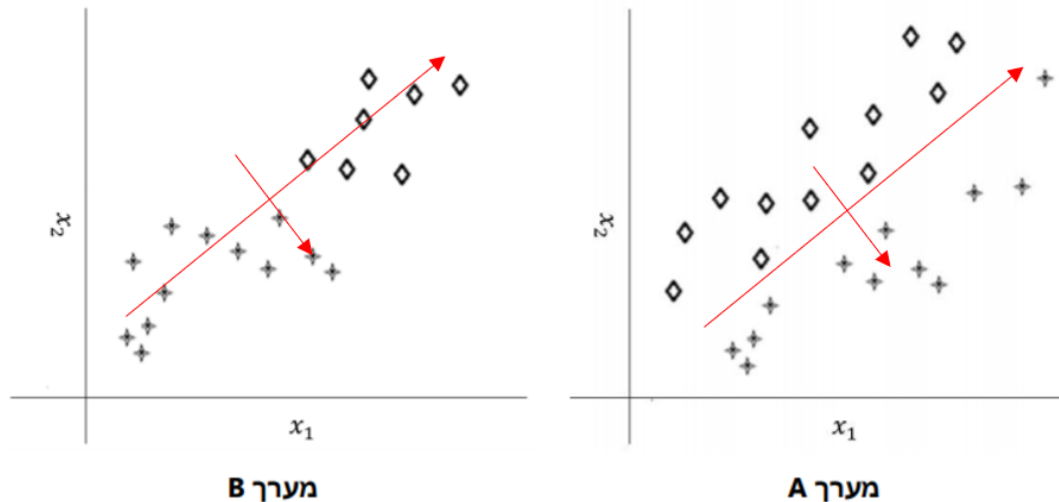
$$\begin{aligned}\hat{y} &= \underset{y=1/3}{\text{argmax}_y} \Pr[y] \cdot \Pr[X_1 = x_1|y] \cdot \Pr[X_2 = x_2|y] \\ &= \underset{y}{\text{argmax}_y} \frac{1}{\sigma_1\sqrt{2\pi}} \exp\left\{-\frac{(x_1 - \mu_{y1})^2}{2\sigma_1^2}\right\} \cdot \frac{1}{\sigma_2\sqrt{2\pi}} \exp\left\{-\frac{(x_2 - \mu_{y2})^2}{2\sigma_2^2}\right\} \\ &= \underset{y}{\text{argmax}_y} \exp\left\{-\frac{(x_1 - \mu_{y1})^2}{8} - \frac{(x_2 - \mu_{y2})^2}{2}\right\} = \underset{y}{\text{argmin}_y} \left[(x_1 - \mu_{y1})^2 + 4(x_2 - \mu_{y2})^2\right]\end{aligned}$$

נותר להציב את הערכים המתאימים לפי סעיף א' לכל סיווג $y = 1, 2, 3$ ולבחור את הסיווג שמחזיר ערך מיטבי.



פיתרון שאלה 2 [20 נק']

בגרפים הבאים מובאים לפניכם שני מערכי נתונים ממימד 2. כל דוגמה מסווגת "◊" או "+".



א. [4 נק'] ציירו במחברת את כל ה- principal components (PC's) עבור כל אחד מהמערכים (אין צורך להעתיק את מערכי הנתונים, אבל יש לצייר מערכת צירים ברורה ואת הווקטורים המבוקשים באופן שהכיוון שלהם ברור, וכך שיהיה ברור מי ה- PC הראשי מבין אלה שציירתם). **ראו ציור.**

ב. [8 נק'] האם ניתן יהיה לסווג נכונה את התצפיות שבמערכים בעזרת מפריד לינארי הפועל על הנקודות כשהן מוטלות על ה- PC הראשון בלבד? הסבירו בקצרה.

מעריך A: לא ניתן להפריד בין הנק' המוטלות על ה PC הראשון בעזרת מפריד לינארי. ניתן לראות שאילו היו מוטלות הנקודות על ה PC המשני אזי כן היה ניתן להפריד אותן בעזרת מפריד לינארי.
מעריך B: כן, הנקודות המוטלות על ה PC הראשון פרידות לינארית.

ג. [8 נק'] לכל אחת מהטענות הבאות, כתבו במחברת התשובות האם היא נכונה או לא נכונה (אין צורך להסביר).

a. המטרה של PCA היא לפרש את המבנה הבסיסי של הנתונים במונחים של הרכיבים העיקריים הטובים ביותר לחיזוי

משתנה הפלט (label). **לא נכון, PCA מקטין את שגיאת ה-Reconstruction והוא בלתי תלוי ב-Labeling.**

b. PC's ששונים מ-0 תמיד יהיו מאונכים זה לזה. **נכון.**

c. למעריך נתונים d -מימדי (dataset המיוצג כמטריצה) תמיד יהיו בדיוק d רכיבים עיקריים (PC's) שונים מ-0.

לא נכון, דוגמא נגדית: מערך נתונים עם מאפיינים קורלטיביים.

d. ניתן לחשב התמרת k-PCA (הטלה ל- k הרכיבים הרכיבים העיקריים הראשונים) באופן שקול על-ידי אימון רשת

עמוקה מסוג autoencoder עם שיכבה אחת נסתרת ברוחב k ופונקציית אקטיבציה טריוויאלית (פונק' הזהות).

נכון, הרשת בעצם טרנספורמציה לינארית עם שמקטינה את ה-Reconstruction error בדיוק כמו ה-PCA.



פיתרון שאלה 3 [25 נק']

היזכרו בבעיית ה-(LS) Least squares: $\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^m (\mathbf{w}^T \mathbf{x}_i - y_i)^2$

א. [4 נק'] יש כמה דרכים לראות את הפיתרון. נדגים אחת.

ראינו בקורס (ומהנוסחה לעיל) שהנורמה מתפרקת על שורות המטריצה \mathbf{X} ועל שורות \mathbf{y} .

לכן, על ידי הוספת שורות, נבנה $\mathbf{X}' = \begin{bmatrix} \mathbf{X} \\ \sqrt{\lambda} \mathbf{I}_{d \times d} \end{bmatrix}$, $\mathbf{y}' = \begin{bmatrix} \mathbf{y} \\ \mathbf{0}_{d \times 1} \end{bmatrix}$

מתקבל:



$$\left\| \begin{bmatrix} \mathbf{X} \\ \sqrt{\lambda} \mathbf{I}_{d \times d} \end{bmatrix} \mathbf{w} - \begin{bmatrix} \mathbf{y} \\ \mathbf{0}_{d \times 1} \end{bmatrix} \right\|_2^2 = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \|\sqrt{\lambda} \mathbf{I}_{d \times d} \mathbf{w} - \mathbf{0}_{d \times 1}\|_2^2 = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

ב. [8 נק'] נציג פיתרון אחד (התקבלו גם פתרונות נוספים, למשל שעמודות \mathbf{X} לא מנורמלות וה-*features* בסדרי גודל שונים).

נניח ש-100 העמודות הראשונות של \mathbf{X} כמעט זהות, למשל עמודות רעש או עמודות זהות ממש לצורך הפשטות.

נניח שאחריהן יש עוד 2 עמודות שמסבירות את וקטור התיגים \mathbf{y} באופן מושלם: $y(x) = 5x_{101} + 2x_{102}$.

שאר העמודות לא משנות.

הקופסה השחורה עשויה להחזיר את וקטור המשקלים $\mathbf{w} = \left[\underbrace{10^8, 10^8, \dots}_{50 \text{ weights}}, \underbrace{-10^8, -10^8, \dots}_{50 \text{ weights}}, 5, 2, \underbrace{0, 0, \dots}_{10^5 - 102 \text{ weights}}, 0 \right]$

שעבורו שגיאת האימון היא 0.

השיטה שהציע הצוות תחזיר את וקטור המשקלים ה"דליל" $\mathbf{w} = \left[\underbrace{10^8, 10^8, \dots}_{50 \text{ weights}}, \underbrace{-10^8, -10^8, \dots}_{50 \text{ weights}}, \underbrace{0, 0, \dots}_{10^5 - 100 \text{ weights}}, 0 \right]$

שעבורו שגיאת האימון גבוהה.

עם זאת, קיים פיתרון דליל $\mathbf{w} = \left[\underbrace{0, 0, \dots}_{100 \text{ weights}}, 0, 5, 2, \underbrace{0, 0, \dots}_{10^5 - 102 \text{ weights}}, 0 \right]$ שעבורו שגיאת האימון היא 0.

ג. [5 נק'] בכלליות רגולריזציה יכולה לעזור כי היא לא תאפשר משקלים כמו 10^8 עבור מאפיינים שלא תורמים להפרדה.

ובפרט, למדנו שרגרסיה ליניארית עם רגולריזציית \mathcal{L}^1 (LASSO) נוטה להחזיר פתרונות דלילים.

ד. [8 נק'] מטילים את הנתונים בעזרת PCA על ידי $\tilde{\mathbf{X}} = \underbrace{\mathbf{X}}_{m \times 10^5} \cdot \underbrace{\mathbf{V}^{(100)}}_{10^5 \times 100}$ ואז פותרים LS ומקבלים $\hat{\mathbf{w}}^{(100)} = LS(\tilde{\mathbf{X}}, \mathbf{y})$

1. שיטה זו יכולה למנוע overfitting כי: (1) PCA נוטה להחליק רעשים ולעיתים התאמת-היתר נובעת מרעשים,

ובנוסף (2) המעבר לרגרסיה במימד נמוך יותר, מקטין את מרחב ההיפותזות ומקשה על הגעה להתאמת יתר.

2. מודל החיזוי יהיה $\tilde{\mathbf{X}} \hat{\mathbf{w}}^{(100)} = \mathbf{X} \underbrace{\mathbf{V}^{(100)} \hat{\mathbf{w}}^{(100)}}_{\hat{\mathbf{w}}} = \mathbf{X} \hat{\mathbf{w}}$ והוא ליניארי ביחס לנתונים המקוריים.



פיתרון שאלה 4 [15 נק']

א. [7 נק'] איזה מודל תבחרו? נמקו את בחירתכם.

מודל B

TN = 96	FP = 4
FN = 10	TP = 90

המודל עומד בדרישות עם עלות:

$$cost_B = 10 \cdot x + 4 \cdot 5x = 30x$$

מודל A

TN = 91	FP = 9
FN = 22	TP = 78

המודל לא עומד בדרישה ($TPR = 78\%$)

מודל D

TN = 98	FP = 2
FN = 18	TP = 82

המודל עומד בדרישות עם עלות:

$$cost_D = 18 \cdot x + 2 \cdot 5x = 28x$$

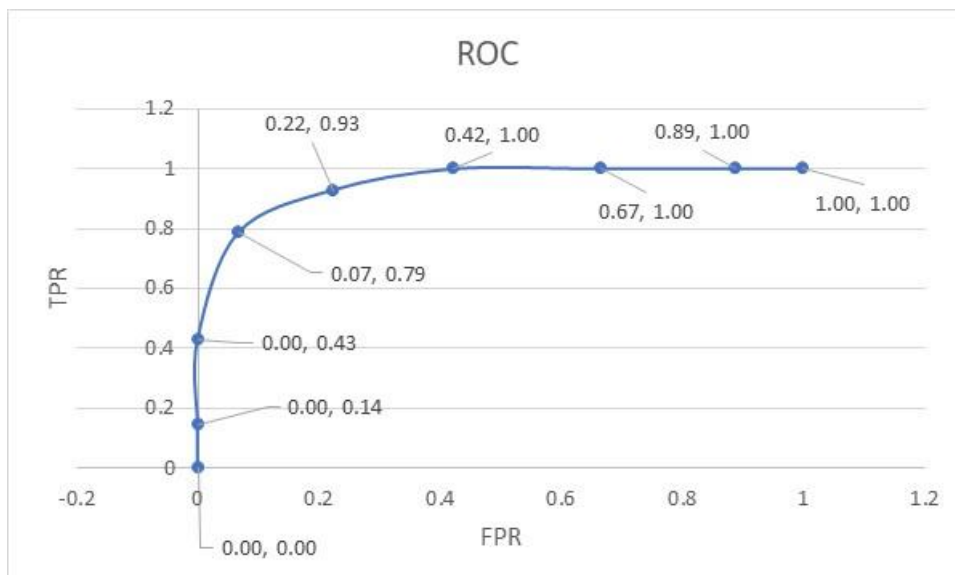
מודל C

TN = 99	FP = 1
FN = 21	TP = 79

המודל לא עומד בדרישה ($TPR = 79\%$)

ולכן התשובה היא מודל D (בטור 2 התשובה הייתה מודל A).

ב. [8 נק'] עקומת ROC מציגה את ביצועי המודל על פי שיעור החיוביים האמיתיים (TPR) מול שיעור החיוביים הכחבים (FPR) כנגד סף החלטה משתנה (decision threshold).





פיתרון שאלה 5 [30 נק'] + מפתח בדיקה

- א. **יכולות לעזור:** הוספת דוגמאות, גריעה של מאפיינים, הוספה של הנחות מקדימות.
הערה: גם הוספה של מאפיינים יכולה לעזור במקרים מסוימים, לכן לא ירדו נקודות למי שסימן אפשרות זו, אולם מי שלא סימן "גריעה של מאפיינים" ירדו לו נקודות.
- ב. ERM הוא אלגוריתם אשר בוחר, מתוך קבוצת היפותזות, את ההיפותזה שמביאה למינימום את שגיאת האימון.
- ג. דוגמאות להשלמת מאפיינים חסרים: החלפת המאפיינים בממוצע המאפיין על פני הדוגמאות הלא חסרות, אלגוריתם EM (כפי שנלמד בכיתה).
- ד. להלן התשובות:
- מקדם פירסון בתחום הממשי $[-1, 1]$ – **נכון**
 - אם בוחרים היפותזה לפי חצי ראשון של הנתונים, אז השגיאה האמפירית על החצי השני הוא משערך בלתי מוטה של שגיאת ההכללה – **נכון**
 - רגרסיה לוגיסטית היא אלגוריתם לסיווג בינארי באמצעות מפרד לינארי – **נכון**
 - פונקציית הלוגריתם של הסיגמואיד היא לא קמורה ולא קעורה – **לא נכון** (היא קעורה)
 - אלגוריתם EM בכל איטרציה מוריד את פונקציית ה-Log-likelihood או משאיר אותה ללא שינוי – **לא נכון** (הוא מעלה או משאיר ללא שינוי)
- ה. גודל רצוי של test set לבחירה של אחת מתוך n היפותזות - לוגריתמי ב- n
- ו. עבור נתונים $x_1 \dots x_n, k$ פונקציית המטרה היא $f(C_1 \dots C_k) = \sum_{i=1}^k \min_{\mu} \sum_{x \in C_i} \|x - \mu\|^2$ כאשר האשכולות הם $C_1 \dots C_k$.
- הערות: מי שכתב ביטוי שקול (למשל, במקום מינימזציה על μ בתוך הסכום, לקח את μ להיות הממוצע על אברי C_i) קיבל את מלוא הנקודות. מי שבמקום $\|x - \mu\|^2$ רשם משהו כדוגמת $(x - \mu)^2$, גם זכה למלוא הניקוד. במילים אחרות, כל מי שהראה בוודאות שהוא יודע שמדובר במרחק אוקלידי בריבוע קיבל את מלוא הנקודות. לעומת זאת, מי שכתב ביטויים כגון $\|x - \mu\|$ או $(x - \mu)$ (כלי העלאה בריבוע) או פונקציות מרחק כלליות כגון $d(x, \mu)$ או $d^2(x, \mu)$ - ירדו לו נקודות. מי שהוסיף ביטויי נירמול שונים (למשל חלוקה ב $|C_i|$ בתוך הסכום... והיו עוד דוגמאות נירמול שונות) - ירדו לו נקודות.



ז. השלמת האלגוריתם:

Step A: חישוב של מרכזי האשכולות (ממוצע) $\forall i = 1..k: \mu_i = \sum_{j: a[j]=i} x_j / \#\{j: a[j] = i\}$

(הערה 1: יש המון דרכים שקולות לכתוב זאת כפסאודוקוד, כולן התקבלו)

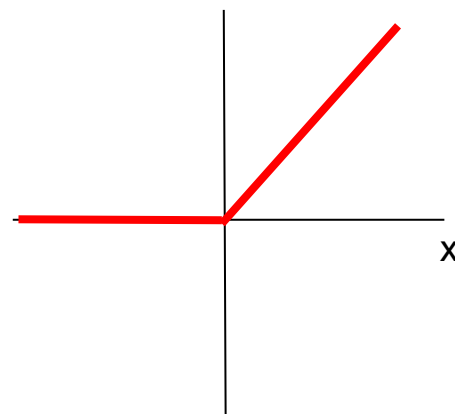
(הערה 2: בעיקרון יש לבדוק שהמכנה ב Step A הוא לא 0, אבל זה "אפקט מסדר שני" ולא נדרש בבחינה)

Step B: שיוך מחדש של נקודות: $\forall j = 1..n: a[j] = \operatorname{argmin}_{i=1..k} \|x_j - \mu_i\|$

(הערה: כאן אפשר אבל לא חובה להעלות את המרחק בריבוע, כי זה שקול, וגם כאן התקבלו ביטויים כגון $(x_j - \mu_i)$,

כלומר בלי סימן של נורמה אוקלידית).

ח. פונקציית ReLU: $\max(0, x)$



ט. פונקציית hinge loss (עבור $y=+1$): $\max(0, 1 - z)$

