

מערכות לומדות - 046195
בחינה סופית - מועד ב'

משך המבחן: 3 שעות

1. אין להשתמש בכל חומר עזר, פרט לדף נוסחאות שיחולק עם הבחינה, ומחשבון.
2. יש לענות על כל השאלות.
3. משקל כל שאלה מצוין בטופס הבחינה. סך הנקודות הוא 100.
4. נא לכתוב בצורה ברורה ומסודרת.
5. יש לכלול פירוט מלא של דרך הפתרון. תשובות לא מנומקות לא יזכו בניקוד.

ב ה צ ל ח ה !!!

שאלה מס' 1 (36%)

שני חלקי השאלה בלתי תלויים זה בזה

חלק א'

נתון: $D = \{x_i, y_i\}_{i=1}^N$, כאשר $x_i \in \mathbb{R}^d$ ו- $y_i \in \mathbb{R}$.
נתון המודל הסטטיסטי הבא: $y_i = w^T x_i + \varepsilon_i$, כאשר $w \in \mathbb{R}^d$ ו- $\{\varepsilon_i\}_{i=1}^N$ הינם משתנים i.i.d. מפולגים נורמלית, $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$.
ידוע מראש כי הפרמטרים לבעיה מתפלגים נורמלית, $w \sim \mathcal{N}(0, \beta^2 I_d)$.
שימו לב: β ו- σ הינם פרמטרים ידועים ואין צורך לשערך אותם.

7% א. הראו כי שערך MAP של הפרמטרים w , שקול לפתרון בעיית הרגרסיה הבאה:

$$(1) w_{\text{linear}} = \arg \min_w \|y - Xw\|^2 + \lambda \|w\|^2$$

רשמו במפורש את $X \in \mathbb{R}^{n \times d}$, $y \in \mathbb{R}^n$ ואת הפרמטר λ כפונקציה של נתוני הבעיה.

6% ב. רשמו את הפתרון לבעיית הרגרסיה (1).

4% ג. הסבירו את המשמעות של הפרמטר λ כתלות בפרמטרים של הבעיה. איזו בעיה הוא נועד לפתור? הסבירו זאת דרך נקודת המבט של בעיית ה-MAP ודרך נקודת המבט של בעיית הרגרסיה.

חלק ב'

הסעיפים בחלק זה בלתי תלויים בסעיפים הקודמים.

נתון: $D = \{(x_i, y_i)\}_{i=1}^N$ כאשר $x \in \mathbb{R}^d$ ו- $y \in \{1, \dots, K\}$.
בחלק זה של השאלה, נדון ברגרסיה של הסתברויות.
במילים אחרות, נרצה למצוא פילוג דיסקרטי $p(y|x)$ רציף ב- x .
לשם כך, נניח מודל פרמטרי, $p_{\text{model}}(y|x, w)$.
בנוסף, נגדיר פונקציית מרחק בין הפילוג האמפירי של ה-Data \hat{p}_{data} ובין ההתפלגות p_{model} :

$$d_{KL}(\hat{p}_{\text{data}} \| p_{\text{model}}) = \mathbb{E}_{y|x \sim \hat{p}_{\text{data}}} \log \frac{\hat{p}_{\text{data}}(y|x)}{p_{\text{model}}(y|x, w)}$$

8% ד. הראו שמזעור ה- KL-divergence, d_{KL} , שקול לשערך MLE. כלומר, הראו בצורה מפורשת שהפתרונות של שתי הבעיות שקולים.

המשך שאלה 1:

הוצע להחסיר את האנטרופיה של p_{model} כאיבר רגולריזציה, באופן הבא:

$$(2) \quad w_{\text{logistic}} = \arg \min_w \sum_{i=1}^N \left(-\log p_{\text{model}}(y_i | x_i, w) + \lambda \sum_{k=1}^K p_{\text{model}}(y = k | x_i, w) \log p_{\text{model}}(y = k | x_i, w) \right)$$

עבור $\lambda > 0$.

4% ה. הסבירו במילים את המשמעות של האיבר שנוסף. מה מעודדת הרגולריזציה במקרה זה?

$$7\% \quad 1. \quad p_{\text{model}}(y = k | x, w) = \frac{\exp(w_k^T x)}{\sum_{j=1}^K \exp(w_j^T x)}, \quad \text{כעת נניח מודל פרמטרי לוגיסטי,}$$

רשמו אלגוריתם גרדיאנט בעל עדכון סדרתי לפתרון בעיית האופטימיזציה הנתונה ב-(2).

הנחיה: ראשית, רשמו את כלל העדכון כתלות ב- $\nabla_w p_{\text{model}}(y | x, w)$.

לאחר מכן, רשמו במפורש את למה שווה הביטוי $\nabla_w p_{\text{model}}(y | x, w)$

פתרון:

(א)

$$\begin{aligned} \arg \max_w p(w | D) &= \arg \max_w p(D | w) p(w) \\ &= \arg \max_w \log p(D | w) p(w) \\ &= \arg \max_w \sum_{i=1}^N \log p(x_i, y_i | w) + \log p(w) \\ &= \arg \max_w -\frac{N}{2} \log 2\pi\sigma^2 - \sum_{i=1}^N \frac{1}{2\sigma^2} (y_i - w^T x_i)^2 - \sum_{i=1}^d \left(\frac{1}{2} \log 2\pi\beta - \frac{1}{2\beta^2} w_i^2 \right) \\ &= \arg \min_w \sum_{i=1}^N \frac{1}{2\sigma^2} (y_i - w^T x_i)^2 + \frac{1}{2\beta^2} w^T w \\ &= \arg \min_w \sum_{i=1}^N (y_i - w^T x_i)^2 + \frac{\sigma^2}{\beta^2} w^T w \\ &= \arg \min_w \|y - Xw\|^2 + \frac{\sigma^2}{\beta^2} \|w\|^2 \end{aligned}$$

כאשר:

$$\begin{aligned} \lambda &= \frac{\sigma^2}{\beta^2} \\ y &= (y_1, \dots, y_N)^T \\ X &= (x_1, \dots, x_N)^T \end{aligned}$$

(ב)
נגזור ונשווה ל-0:

$$\begin{aligned}\frac{d}{dw}(\|y - Xw\|^2 + \lambda \|w\|^2) &= 0 \\ \Leftrightarrow X^T(y - Xw) + \lambda w &= 0 \\ \Leftrightarrow (X^T X + \lambda I)w &= X^T y \\ \Leftrightarrow w &= (X^T X + \lambda I)^{-1} X^T y\end{aligned}$$

(ג)

λ הינו מקדם אשר קובע כמה חשיבות נותנים לאיבר הרגולריזציה $\|w\|^2$. $\lambda = \frac{\sigma^2}{\beta^2}$ מכמת את מידת הוודאות שלנו ב-Prior אל מול הוודאות ה-Data. ככל שהיחס $\frac{\sigma}{\beta}$ קטן יותר, אנחנו מתקרבים לבעיה המקורית של שערך MLE ומציאת Fit ל-Data. ככל שהיחס $\frac{\sigma}{\beta}$ גדל, נחפש פתרון בעל משקולות קטנים יותר, על אף שנוסיף Bias לפתרון בעיית הרגרסיה $\|y - Xw\|^2$. $\arg \min_w \|y - Xw\|^2$ חיפוש פתרון בעל משקולות קטנים יותר מקטין באופן אפקטיבי את מחלקת ההשערות שבה אנו מחפשים פתרון לבעיית הרגרסיה ובכך מקטין את ה-Variance של הבעיה. בנוסף, הוא מאפשר לעיתים יציבות נומרית גדולה יותר.

(ד)

נרשום במפורש את תוצאת בעיית מזעור ה-KL-divergence ונראה שהוא זהה לפתרון בעיית ה-MLE:

$$\begin{aligned}\arg \min_w d_{KL}(\hat{p}_{data} \parallel p_{model}) &= \arg \min_w \mathbb{E}_{y|x \sim \hat{p}_{data}} \log \frac{\hat{p}_{data}(y|x)}{p_{model}(y|x, w)} \\ &= \arg \min_w \mathbb{E}_{y|x \sim \hat{p}_{data}} [\log \hat{p}_{data}(y|x) - \log p_{model}(y|x, w)] \\ &= \arg \min_w - \mathbb{E}_{y|x \sim \hat{p}_{data}} \log p_{model}(y|x, w) \\ &= \arg \max_w \mathbb{E}_{y|x \sim \hat{p}_{data}} \log p_{model}(y|x, w) \\ &= \arg \max_w \frac{1}{N} \sum_{i=1}^N \log p_{model}(y_i | x_i, w) \\ &= \arg \max_w p_{model}(y|x, w) \\ &= w_{MLE}\end{aligned}$$

(ה)

ניזכר כי פונקציית האנטרופיה הינה:

$$h(p_{model} | x_i, w) = - \sum_{k=1}^K p_{model}(y = k | x_i, w) \log p_{model}(y = k | x_i, w)$$

כלומר במקרה שלנו אנחנו מנסים למזער את מינוס האנטרופיה, או במילים אחרות, למקסם את האנטרופיה. כזכור מהתרגול על עצים, אנטרופיה היא מדד לחוסר האחידות של משתנה אקראי. לכן מקסום האנטרופיה שקול להעדפת פתרונות שבהם ההתפלגות $p_{model}(y = k | x_i, w)$ רחוקה מהתפלגות דטרמיניסטית.

(7)

נסמן את פונקציית האקטיביציה הלוגיסטית:

$$p_{\text{model}}(y=i|x,w) = \frac{e^{w_i^T x}}{\sum_{k=1}^K e^{w_k^T x}}$$

$$\Rightarrow \nabla_{w_j} p_{\text{model}}(y=i|x,w) = \begin{cases} \frac{e^{w_i^T x} \sum_{k=1}^K e^{w_k^T x} x - e^{w_i^T x} e^{w_j^T x} x}{\left(\sum_{k=1}^K e^{w_k^T x}\right)^2} & i=j \\ \frac{-e^{w_i^T x} e^{w_j^T x} x}{\left(\sum_{k=1}^K e^{w_k^T x}\right)^2} & i \neq j \end{cases} = \begin{cases} \frac{e^{w_i^T x} \sum_{k=1}^K e^{w_k^T x} - e^{w_i^T x} e^{w_j^T x}}{\left(\sum_{k=1}^K e^{w_k^T x}\right)^2} x & i=j \\ \frac{-e^{w_i^T x} e^{w_j^T x}}{\left(\sum_{k=1}^K e^{w_k^T x}\right)^2} x & i \neq j \end{cases}$$

כאשר,

$$\nabla_w p_{\text{model}}(y=i|x,w) = (\nabla_{w_1} p_{\text{model}}(y=i|x,w), \dots, \nabla_{w_K} p_{\text{model}}(y=i|x,w))$$

מכאן, נחשב את הגרדיאנט:

$$\begin{aligned} \Delta w &= \nabla_w \left(-\log p_{\text{model}}(y_t|x_t,w) + \lambda \sum_{k=1}^K p_{\text{model}}(y=k|x_t,w) \log p_{\text{model}}(y=k|x_t,w) \right) = \\ &= -\frac{1}{p_{\text{model}}(y_t|x_t,w_t)} \nabla_w p_{\text{model}}(y_t|x_t,w) + \lambda \sum_{k=1}^K (1 + p_{\text{model}}(y=k|x_t,w)) \nabla_w p_{\text{model}}(y=k|x_t,w) \\ &= -\frac{1}{g_{y_t}(x|w)} \nabla_w g_{y_t}(x|w) + \lambda \sum_{k=1}^K (1 + g_k(x|w)) \nabla_w g_k(x|w) \end{aligned}$$

ולבסוף נקבל את כלל העדכון הבא:

$$w_{t+1} = w_t - \eta \Delta w_t$$

שאלה מס' 2 (38%)

בשאלה זו נתמקד ברשת עצבית היזון קדמי (feedforward) עם מטריצות משקולות $\{W_\ell\}_{\ell=1}^L$, כאשר $W_\ell \in \mathbb{R}^{d_\ell \times d_{\ell-1}}$, ללא איברי היסט (Bias).

מוצא כל שכבה עובר דרך פונקציית אקטיבציה לינארית: $\varphi(u) = u$. המופעלות בבעיה עם קלט $X \in \mathbb{R}^{d_0}$ ופלט $O \in \mathbb{R}^{d_L}$, עבור $d_L < d_0$. כלומר, מוצא הרשת הוא: $O = W_L W_{L-1} \dots W_2 W_1 X$. לדוגמא: נסתכל על הנוירון ה- j בשכבה הראשונה אז מוצאו היינו $\varphi_j(\sum_{i=1}^{d_0} w_{j,i} x_i) = \sum_{i=1}^{d_0} w_{j,i} x_i$ ועבור כל המוצאים של השכבה הראשונה נקבל $O_1 = W_1 X$ כאשר:

$$W_1 = \begin{pmatrix} w_{1,1} & \dots & w_{1,d_0} \\ \vdots & & \vdots \\ w_{d_1,1} & \dots & w_{d_1,d_0} \end{pmatrix}, X = \begin{pmatrix} x_1 \\ \vdots \\ x_{d_0} \end{pmatrix}$$

חלק א'

5% א. כמה שאלות לחימום:

- I. כמה שכבות נסתרות יש?
- II. כמה נוירונים יש בשכבה ה- ℓ ?
- III. כמה פרמטרים נלמדים יש ברשת?
- IV. לאחר שהרשת נלמדה, הראו כי מספיק לשמור $d_0 d_L$ פרמטרים ברשת.

7% ב. נניח שניתן לבחור את L ו- $\{W_\ell\}_{\ell=1}^L$ כרצוננו. האם ניתן להשתמש במשפט הקירוב האוניברסלי? אם כן, נמקו. אם לא, הגדירו מה היא משפחת הפונקציות אותה ניתן לממש באמצעות רשת זו.

10% ג. חשבו את $\frac{\partial O}{\partial W_\ell}$, לפי אלגוריתם Backpropagation בהינתן קלט X ומטריצות $\{W_\ell\}_{\ell=1}^L$.

חלק ב':

מכאן ואילך נניח כי הפלט הוא סקלר $d_L = 1$, וכי הרשת מאותחלת ע"י דגימה אקראית של כל משקל בצורה בלתי תלויה מפילוג תלוי שכבה עם שונות σ_ℓ^2 וממוצע אפס.

8% ד. נניח שרכיבי הקלט X נדגמו מפילוג עם ממוצע אפס ושונות 1. מצאו אתחול אקראי למשקולות עבורו גם הכניסה לכל נוירון ברשת היא עם ממוצע אפס ושונות 1.

8% ה. נסמן את יציאת שכבת הנוירונים ה- ℓ ב- $V_\ell = W_\ell W_{\ell-1} \dots W_1 X$, ואת הגרדיאנט $G_\ell = \frac{\partial O}{\partial V_\ell}$.

. מצאו אתחול אקראי למשקולות עבורו רכיבי הגרדיאנט $G_\ell : \forall \ell < L$ כולם עם ממוצע אפס ושונות 1.

פתרון:

(א)

ניתן לראות מיידית:
ברשת יש $L-1$ שכבות נסתרות.
בשכבה הנסתרת ה- ℓ יש d_ℓ נויורונים.

בסה"כ יש $\sum_{\ell=1}^L d_\ell d_{\ell-1}$ פרמטרים, מכיוון שזו סכום גדלי מטריצות המשקולות.

אם נסמן $W = W_L W_{L-1} \dots W_2 W_1 \in \mathbb{R}^{d_0 \times d_L}$, נוכל לראות שהחזאי שנלמד שקול לחזאי הלינארי $O = WX$,
ולכן נוכל לשמור בסה"כ $d_0 d_L$ פרמטרים – מספר הערכים ב- W .

(ב)

מכיוון שהחזאי הנלמד הוא לינארי, ניתן לממש רק פונקציות לינאריות מהצורה: $f(X) = WX$. בנוסף ברור שגם ניתן לממש כל פונקציה לינארית מצורה זו: פשוט נקבע $W_1 = W$ ולכל שאר המטריצות $W_\ell = I$.
מכיוון שכך לא ניתן לקרב כל פונקציה רציפה עם רשת מסוג זה. מכאן, משפט הקירוב האוניברסלי לא תקף במקרה זה. הסיבה שהמשפט לא תקף היא שבמשפט יש תנאי שפונקציית האקטיבציה היא לא פולינום, בעוד שמקרה זה היא כן פולינום: $\varphi(u) = u$.

(ג)

נשים לב כי המוצא $O \in \mathbb{R}^{d_L}$, כלומר ישנן d_L יציאות לרשת הנוירונים.

לכן, הנגזרת ביחס למטריצת המשקולות ה- l הינה $\frac{\partial O}{\partial W_\ell} = \left(\frac{\partial O_1}{\partial W_\ell}, \dots, \frac{\partial O_{d_L}}{\partial W_\ell} \right)$.

מעתה, נתייחס לגרדיאנט של היציאה O_i (היציאה ה- i) בלבד.

נסמן את יציאת (וכניסת) שכבת הנוירונים ה- ℓ ב- $V_\ell = W_\ell W_{\ell-1} \dots W_1 X$, ואת הגרדיאנט $G_{i,\ell} = \frac{\partial O_i}{\partial V_\ell}$.

המוצא של השכבה האחרונה באינדקס ה- i הינו למעשה היציאה ה- i , ולכן, $O_i = v_{i,L}$, ולכן,

$$g_{i,L} = \frac{\partial O_i}{\partial v_{i,L}} = \frac{\partial O_i}{\partial O_i} = 1$$

כאשר סימנו $g_{i,L}$ באותיות קטנות כדי לציין גודל סקלרי.

לכל שכבה אחרת, נחשב בצורה רקורסיבית:

$$G_{i,\ell-1} = \frac{\partial O_i}{\partial V_{\ell-1}} = \frac{\partial V_\ell}{\partial V_{\ell-1}} \frac{\partial O_i}{\partial V_\ell} = W_\ell^T G_{i,\ell}$$

כך שיתקיים:

$$G_{i,\ell} = W_{\ell+1}^T \dots W_{L-1}^T W_{i,L}^T g_{i,L} = W_{\ell+1}^T \dots W_{L-1}^T W_{i,L}^T$$

לסיום, הנגזרת לפי המשקל היא:

$$\frac{\partial O_i}{\partial W_\ell} = \frac{\partial O_i}{\partial V_\ell} \frac{\partial V_\ell}{\partial W_\ell} = G_{i,\ell} V_{\ell-1}^T$$

כך שנקבל:

$$\frac{\partial O_i}{\partial W_\ell} = \frac{\partial O_i}{\partial V_\ell} \frac{\partial V_\ell}{\partial W_\ell} = W_{\ell+1}^T \dots W_{L-1}^T W_{i,L}^T X^T W_1^T \dots W_{\ell-2}^T W_{\ell-1}^T$$

כדי לוודא את נכונות החישוב, נשים לב למימדיות:

$$W_{\ell+1}^T \cdots W_{L-1}^T W_{i,L}^T = (W_{i,L} W_{L-1} \cdots W_{\ell+1})^T \in \mathbb{R}^{d_i \times 1}$$

ובאותו אופן,

$$X^T W_1^T \cdots W_{\ell-2}^T W_{\ell-1}^T = (W_{\ell-1} \cdots W_1 X)^T \in \mathbb{R}^{1 \times d_{\ell-1}}$$

ולכן,

$$\frac{\partial O_i}{\partial W_\ell} \in \mathbb{R}^{d_i \times d_{\ell-1}}$$

לבסוף, ניזכר שעשינו זאת עבור כל יציאה בנפרד, לכן, הנגזרת הכוללת הינה:

$$\frac{\partial O}{\partial W_\ell} = \left(\frac{\partial O_1}{\partial W_\ell}, \dots, \frac{\partial O_{d_L}}{\partial W_\ell} \right) \in \mathbb{R}^{d_L \times d_\ell \times d_{\ell-1}}$$

כצפוי לנגזרת של וקטור באורך d_L במטריצה בגודל $d_\ell \times d_{\ell-1}$.

(ד)

נדרוש שסכום הכניסות לכל ניורון יהיה בעל ממוצע אפס ושונות 1. נתחיל בכניסה לשכבה הנסתרת הראשונה.

הדרישה לממוצע אפס מתקיימת בכל מקרה כי $E[X] = 0$ ו- $E[W_1] = 0$

$$0 = E[W_1 X] = E[W_1] E[X]$$

מהדרישה לשונות 1 (לשם הפשטות נסמן ב- w_{ij} במקום ב- $w_{ij,1}$ את רכיבי המטריצה W_1):

$$\forall i \in 1, \dots, d_1: 1 = \text{Var}((W_1 X)_i) = \text{Var}\left(\sum_{j=1}^{d_0} w_{ij} x_j\right) = \text{Cov}\left(\sum_{j=1}^{d_0} w_{ij} x_j, \sum_{r=1}^{d_0} w_{ir} x_r\right) = \sum_{r=1}^{d_0} \sum_{j=1}^{d_0} \text{Cov}(w_{ij} x_j, w_{ir} x_r)$$

כאשר:

$$\begin{aligned} \text{Cov}(w_{ij} x_j, w_{ir} x_r) &= E[w_{ij} x_j w_{ir} x_r] - E[w_{ij} x_j] E[w_{ir} x_r] \\ &= E[w_{ij} w_{ir}] E[x_j x_r] - E[w_{ij}] E[w_{ir}] E[x_j] E[x_r] \\ &= E[w_{ij} w_{ir}] E[x_j x_r] \end{aligned}$$

כאשר במעבר האחרון השתמשנו בכך ש- $E[X] = 0$ ו- $E[W_1] = 0$. מכך שרכיבי האתחול בלתי תלויים ועם

שונות σ_1^2 , נקבל

$$\text{Cov}(w_{ij} x_j, w_{ir} x_r) = \delta_{jr} \sigma_1^2 \text{Var}(x_j) = \delta_{jr} \sigma_1^2$$

ולכן:

$$1 = \sum_{r=1}^{d_0} \sum_{j=1}^{d_0} \text{Cov}(w_{ij} x_j, w_{ir} x_r) = \sum_{r=1}^{d_0} \sum_{j=1}^{d_0} \delta_{jr} \sigma_1^2 = \sum_{j=1}^{d_0} \sigma_1^2 = d_0 \sigma_1^2$$

מכאן הבחירה $\sigma_1 = \frac{1}{\sqrt{d_0}}$ תיתן $1 = \text{Var}((W_1 X)_i)$, כנדרש.

באופן דומה לשאר הניורונים: $\forall \ell = 1, \dots, L: \sigma_\ell = \frac{1}{\sqrt{d_{\ell-1}}}$.

(ה)

לשם הגיוון, בסעיף זה נבצע את החישוב על הביטויים הסופיים, ולא שכבה-שכבה כמו בסעיף הקודם. מכיוון שכל המשקולות בלתי תלויות ועם ממוצעים שווים לאפס, נקבל:

$$E[G_\ell] = E[W_{\ell+1}^T \cdots W_{L-1}^T W_L^T] = 0$$

עכשיו, נחשב את השונות ושוב, מכיוון שכל המשקולות בלתי תלויות ועם ממוצעים שווים לאפס, נקבל:

$$\begin{aligned}\text{Var}(g_{i,\ell}) &= \text{Var}\left(\sum_{i_{\ell+1}, \dots, i_L} w_{ii_{\ell+1}, \ell+1} \cdots w_{i_L, L}\right) \\ &= \sum_{i_{\ell+1}, \dots, i_L} \text{Var}(w_{ii_{\ell+1}, \ell+1}) \cdots \text{Var}(w_{i_L, L}) \\ &= \sum_{i_{\ell+1}, \dots, i_L} \sigma_{\ell+1}^2 \cdots \sigma_L^2 = d_{\ell+1} \cdots d_L \sigma_{\ell+1}^2 \cdots \sigma_L^2\end{aligned}$$

מכאן, נקבל שבחירה $\sigma_\ell = \frac{1}{\sqrt{d_\ell}}$, $\forall \ell = 1, \dots, L$, תיתן לנו $\text{Var}(g_{i,\ell}) = 1$.

אסופת שאלות בנושאים שונים:

(שערוך)

8% א.

נתונות לנו N מדידות IID $\{x_i\}_{i=1}^N$ כאשר x_i מגיע מהתפלגות הבאה:

$$P_x(x) = \frac{1}{\mu} e^{-\frac{1}{\mu}(x-\theta)}, \quad x \geq \theta, \quad \mu > 0$$

(1) מצאו את משעריך MLE עבור הפרמטר μ בהנחה כי θ פרמטר ידוע.(2) מצאו את משעריך MLE עבור הפרמטר θ בהנחה כי μ פרמטר ידוע.פתרון:

א. כמו בסעיף הקודם, פונקציית ה-Likelihood (הפעם כפונקציה של μ כי הוא המשתנה הלא ידוע בסעיף זה):

$$L(\mu) = \prod_{i=1}^N \frac{1}{\mu} e^{-\frac{1}{\mu}(x_i-\theta)} = \frac{1}{\mu^N} e^{-\frac{1}{\mu} \sum_{i=1}^N (x_i-\theta)} I_{\{\mu > 0\}}$$

נניח כאן כי $\mu > 0$ אחרת האינדקטור מתאפס.

$$l(\mu) = \log L(\mu) = -N \log \mu - \frac{1}{\mu} \sum_{i=1}^N (x_i - \theta)$$

מגזירה והשוואה לאפס נקבל

$$l'(\mu) = -\frac{N}{\mu} + \frac{1}{\mu^2} \sum_{i=1}^N (x_i - \theta) = 0 \Rightarrow \hat{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^N (x_i - \theta)$$

הנגזרת השנייה שלילית ולכן זוהי אכן נקודת מקסימום.

ב. נכתוב את ה-Likelihood:

$$L(\theta) = \prod_{i=1}^N \frac{1}{\mu} e^{-\frac{1}{\mu}(x_i-\theta)} I_{\{x_i \geq \theta\}} = \frac{1}{\mu^N} e^{-\frac{1}{\mu} \sum_{i=1}^N (x_i-\theta)} I_{\{\min_i x_i \geq \theta\}}$$

כאשר I_A היא פונקציית אינדקטור (מקבלת 1 אם המאורע A מתקיים ו-0 אחרת). נשים לב כי $L(\theta)$ היא פונקציה מונוטונית עולה ב- θ בתחום שבו $\theta \leq \min_i x_i$. לכן, משעריך הסבירות המירבית יתקבל

בערך המקסימלי האפשרי עבור θ בתחום זה: $\hat{\theta}_{MLE} = \min_i x_i$.

(עצים)

6% ב. רותי רוצה לבנות עץ החלטה המבוסס על קריטריון האנטרופיה וקובע על סמך העונה, מצב הלחות והטמפרטורה האם צפוי לרדת גשם או לא. ברשותה 6 מדידות מהעבר:

עונה	מצב הלחות	טמפרטורה	האם ירד גשם?
אביב	גבוהה	נמוכה	כן
אביב	נמוכה	נמוכה	לא
סתיו	גבוהה	נמוכה	לא
סתיו	גבוהה	גבוהה	כן
קיץ	גבוהה	נמוכה	כן
קיץ	נמוכה	גבוהה	לא

מהו הקריטריון שישמש בצומת הראשונה בעץ?

פתרון:

נחשב את האנטרופיה שתושרה לאחר פיצול לפי כל אחד מהמאפיינים האפשריים-
עונה:

עונה	האם ירד גשם?
אביב	כן
אביב	לא
סתיו	לא
סתיו	כן
קיץ	כן
קיץ	לא

אביב: +1/-1.

$$H(\text{rain}|\text{winter}) = -\frac{1}{2}\log_2\left(\frac{1}{2}\right) - \frac{1}{2}\log_2\left(\frac{1}{2}\right) = 1$$

סתיו: +1/-1.

$$H(\text{rain}|\text{winter}) = -\frac{1}{2}\log_2\left(\frac{1}{2}\right) - \frac{1}{2}\log_2\left(\frac{1}{2}\right) = 1$$

קיץ: +1/-1.

$$H(\text{rain}|\text{winter}) = -\frac{1}{2}\log_2\left(\frac{1}{2}\right) - \frac{1}{2}\log_2\left(\frac{1}{2}\right) = 1$$

אנטרופיה כללית בחלוקה על סמך עונה:

$$H(\text{rain}|\text{season}) = \frac{2}{6} \cdot 1 + \frac{2}{6} \cdot 1 + \frac{2}{6} \cdot 1 = 1$$

מצב הלחות:

האם ירד גשם?	מצב הלחות
כן	גבוהה
לא	נמוכה
לא	גבוהה
כן	גבוהה
כן	גבוהה
לא	נמוכה

לחות גבוהה: +3/-1

$$H(\text{rain}|\text{high humidity}) = -\frac{3}{4}\log_2\left(\frac{3}{4}\right) - \frac{1}{4}\log_2\left(\frac{1}{4}\right) = 0.81$$

לחות נמוכה: +0/-2

$$H(\text{rain}|\text{low humidity}) = 0$$

אנטרופיה כללית בחלוקה על סמך מצב הלחות:

$$H(\text{rain}|\text{humidity}) = \frac{4}{6} \cdot 0.81 + \frac{2}{6} \cdot 0 = 0.54$$

טמפרטורה:

האם ירד גשם?	טמפרטורה
כן	נמוכה
לא	נמוכה
לא	נמוכה
כן	גבוהה
כן	נמוכה
לא	גבוהה

טמפרטורה נמוכה: +2/-2

$$H(\text{rain}|\text{low temperature}) = -\frac{2}{4}\log_2\left(\frac{2}{4}\right) - \frac{2}{4}\log_2\left(\frac{2}{4}\right) = 1$$

טמפרטורה גבוהה: +1/-1

$$H(\text{rain}|\text{high temperature}) = -\frac{1}{2}\log_2\left(\frac{1}{2}\right) - \frac{1}{2}\log_2\left(\frac{1}{2}\right) = 1$$

אנטרופיה כללית בחלוקה על סמך גובה הטמפרטורה:

$$H(\text{rain}|\text{humidity}) = 1$$

קריטריון הלחות השיג את האנטרופיה הכללית הכי נמוכה ולכן זהו הולך להיות הפיצול הראשון בעץ.

(SVM)

ג 6% נתונה בעיית האופטימיזציה הבאה:

$$\begin{aligned} \min_{\tilde{w} \in \mathbb{R}^{d+1}, \xi \in \mathbb{R}^n} & \left(\frac{1}{2} \|\tilde{w}\|^2 + C \sum_{i=1}^n \xi_i \right) \\ \text{subject to} & \quad y_i \tilde{w}^T \tilde{x}_i \geq 1 - \xi_i \quad i = 1, \dots, n \\ & \quad \xi_i \geq 0 \quad i = 1, \dots, n \end{aligned} \quad (P)$$

עבור $\tilde{w} = (w, b)$, $\tilde{x}_i = (x_i, 1)$

הוכיחו או הפריכו את הטענה הבאה:

בעיית האופטימיזציה הנתונה שקולה לבעיית ה-Soft-SVM שנלמדה בכתה.

פתרון:

אין הבדל בנושא זה בין בעיית ה-SVM הרגילה ובעיית ה-Soft-SVM, ולכן נדון בבעיה הרגילה.

בעיית ה-SVM הפרימאלית ניתנת לפי:

$$\begin{aligned} & \begin{cases} \min_{\tilde{w} \in \mathbb{R}^{m+1}} \frac{1}{2} \|\tilde{w}\|^2 \\ \text{s.t. } y_i \tilde{w}^T \tilde{x}_i \geq 1 \quad i = 1, \dots, n \end{cases} = \\ & = \begin{cases} \min_{w \in \mathbb{R}^m, b \in \mathbb{R}} \frac{1}{2} \|w\|^2 + \frac{1}{2} b^2 \\ \text{s.t. } y_i \tilde{w}^T \tilde{x}_i \geq 1 \quad i = 1, \dots, n \end{cases} \end{aligned}$$

כעת יש רגורלריזציה על איבר ה-BIAS בניגוד לבעיה המקורית. לכן מדובר בשתי בעיות שונות עם פתרון שונה.

נפרש בעיה זו – כעת, אנחנו רוצים למצוא מישור מפריד בעל Margin גדול ככל הניתן, אולם תוך כדי מתן העדפה לישרים מפרידים אשר עוברים בקירוב דרך הראשית (כלומר בעלי b קטן).

(PCA)

6% ד. עבור סדרת נקודות נתונות $\{x_1, \dots, x_n\}$ ב- \mathbb{R}^2 (בעלות ממוצע 0), חושבה מטריצת קווריאנס-המדגם הבאה:

$$P_n = \begin{pmatrix} 3 & 2 \\ 2 & 6 \end{pmatrix}$$

(1) איזה מהווקטורים הבאים מייצג (עד כדי נירמול בקבוע) את הכיוון העיקרי (הראשון) של הנקודות הנתונות?

$$w_1 = \begin{pmatrix} -2 \\ 1 \end{pmatrix}, w_2 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, w_3 = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

(2) חשבו את שני הרכיבים הראשיים של $x = (1, 0)^T$

פתרון:

(1)

בחישוב מהיר ניתן לראות:

$$P_n w_1 = \begin{pmatrix} 3 & 2 \\ 2 & 6 \end{pmatrix} \begin{pmatrix} -2 \\ 1 \end{pmatrix} = \begin{pmatrix} -6+2 \\ -4+6 \end{pmatrix} = \begin{pmatrix} -4 \\ 2 \end{pmatrix} = 2 \begin{pmatrix} -2 \\ 1 \end{pmatrix}$$

מכאן, ש w_1 הינו וקטור עצמי בעל ערך עצמי 2.

$$P_n w_2 = \begin{pmatrix} 3 & 2 \\ 2 & 6 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 3+2 \\ 2+6 \end{pmatrix} = \begin{pmatrix} 5 \\ 8 \end{pmatrix} \neq \alpha \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

כלומר, w_2 איננו וקטור עצמי של מטריצת הקווריאנס.

$$P_n w_3 = \begin{pmatrix} 3 & 2 \\ 2 & 6 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} = \begin{pmatrix} 3+4 \\ 2+12 \end{pmatrix} = \begin{pmatrix} 7 \\ 14 \end{pmatrix} = 7 \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

כלומר, w_3 אף הוא וקטור עצמי, בעל ערך עצמי 7.

מאחר שערך עצמי זה המתאים ל w_3 גבוה מהערך העצמי המתאים לקטור העצמי w_1 , הרי שהוא מייצג את הכיוון העיקרי של הנקודות הנתונות.

(2)

כעת, נרצה להטיל את הוקטור $x = (1 \ 0)^T$ על מערכת הצירים של הכיוונים העיקריים (PRINCIPAL COMPONENTS).

$$x_{PCA}(1) = \frac{w_3^T x}{\|w_3\|} = \frac{1}{\sqrt{5}}$$

$$x_{PCA}(2) = \frac{w_1^T x}{\|w_1\|} = \frac{-2}{\sqrt{5}}$$

$$\Rightarrow x_{PCA} = \frac{1}{\sqrt{5}} \begin{pmatrix} 1 & -2 \end{pmatrix}$$