

IML תרגיל 3, מגיש: אבי כוגן, ת.ז: 205417710

שאלה 1

$$h_D(x) = \underset{y \in \{\pm 1\}}{\operatorname{argmax}} Pr(y|x) \stackrel{\text{Bayes theorem}}{=} \underset{y \in \{\pm 1\}}{\operatorname{argmax}} \frac{Pr(x|y)Pr(y)}{Pr(x)} \stackrel{*}{=} \underset{y \in \{\pm 1\}}{\operatorname{argmax}} Pr(x|y)Pr(y)$$

* מתקיים $Pr(x)$ לא תלוי ב- y וחיובי x יצא לכן ההסתברות עליו אינה 0, לכן מקסימום על המונה בלבד ימקסם את הביטוי.

שאלה 2

$$\begin{aligned} h_D(x) &= Pr(x|y)Pr(y) = \\ &= \underset{y \in \{\pm 1\}}{\operatorname{argmax}} \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left\{-\frac{1}{2}(x - \mu_y)^T \Sigma^{-1}(x - \mu_y)\right\} Pr(y) \\ &\stackrel{*}{=} \underset{y \in \{\pm 1\}}{\operatorname{argmax}} \ln\left(\frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left\{-\frac{1}{2}(x - \mu_y)^T \Sigma^{-1}(x - \mu_y)\right\} Pr(y)\right) = \\ &= \underset{y \in \{\pm 1\}}{\operatorname{argmax}} -\frac{1}{2}(x - \mu_y)^T \Sigma^{-1}(x - \mu_y) + \ln(Pr(y)) = \\ &\stackrel{**}{=} \underset{y \in \{\pm 1\}}{\operatorname{argmax}} -\frac{1}{2}x^T \Sigma^{-1}x + \frac{1}{2}x^T \Sigma^{-1}\mu_y + \frac{1}{2}\mu_y^T \Sigma^{-1}x - \frac{1}{2}\mu_y^T \Sigma^{-1}\mu_y + \ln(Pr(y)) = \\ &\triangleq \underset{y \in \{\pm 1\}}{\operatorname{argmax}} x^T \Sigma^{-1}\mu_y - \frac{1}{2}\mu_y^T \Sigma^{-1}\mu_y + \ln(Pr(y)) = \underset{y \in \{\pm 1\}}{\operatorname{argmax}} \delta_y(x) \end{aligned}$$

* \ln פונק' מונוטונית עולה, לכן מקסימום על הביטוי שווה למקסימום על \ln של הביטוי.
 ** $\ln\left(\frac{1}{\sqrt{(2\pi)^d |\Sigma|}}\right)$ לא תלוי ב- y ולכן ניתן
 $\triangle -\frac{1}{2}x^T \Sigma^{-1}x$ לא תלוי ב- y .

שאלה 3

נשתמש באומדי MLE לתוחלות ולשוונות.

נחשב את $\hat{\mu}_{+1}, \hat{\mu}_{-1}, \hat{\Sigma}$: נחפש $\operatorname{argmax} L(\hat{\mu}_{+1}, \hat{\mu}_{-1}, \hat{\Sigma} | x_1/y_1, \dots, x_m/y_m)$

$$\begin{aligned} L(\hat{\mu}_{+1}, \hat{\mu}_{-1}, \hat{\Sigma} | x_1/y_1, \dots, x_m/y_m) &= \prod_{i=1}^m P((x_i/y_i)/\mu_{+1}, \mu_{-1}, \Sigma) = \\ &= \prod_{i=1}^k P((x_i/y_i = 1)/\mu_{+1}, \Sigma) \cdot \prod_{j=k+1}^m P((x_j/y_j = -1)/\mu_{-1}, \Sigma) = \\ &= \left(\frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \right)^m \exp \left\{ \sum_{i=1}^k -\frac{1}{2} (x_i - \mu_{+1})^T \Sigma^{-1} (x_i - \mu_{+1}) \right\} \cdot \exp \left\{ \sum_{j=k+1}^m -\frac{1}{2} (x_j - \mu_{-1})^T \Sigma^{-1} (x_j - \mu_{-1}) \right\} \end{aligned}$$

מתקיים: $\operatorname{argmax} L(\hat{\mu}_{+1}, \hat{\mu}_{-1}, \hat{\Sigma} | x_1/y_1, \dots, x_m/y_m) = \operatorname{argmax} \ln(L(\hat{\mu}_{+1}, \hat{\mu}_{-1}, \hat{\Sigma} | x_1/y_1, \dots, x_m/y_m))$

נסמן $l(\hat{\mu}_{+1}, \hat{\mu}_{-1}, \hat{\Sigma} | x_1/y_1, \dots, x_m/y_m) = \ln(L(\hat{\mu}_{+1}, \hat{\mu}_{-1}, \hat{\Sigma} | x_1/y_1, \dots, x_m/y_m))$

$$\begin{aligned} l(\hat{\mu}_{+1}, \hat{\mu}_{-1}, \hat{\Sigma} | x_1/y_1, \dots, x_m/y_m) &= \\ \ln \left(\left(\frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \right)^m \right) &+ \sum_{i=1}^k -\frac{1}{2} ((x_i/y = 1) - \mu_{+1})^T \Sigma^{-1} ((x_i/y = 1) - \mu_{+1}) + \\ \sum_{j=k+1}^m -\frac{1}{2} ((x_j/y = -1) - \mu_{-1})^T \Sigma^{-1} ((x_j/y = -1) - \mu_{-1}) \end{aligned}$$

נקבל את האומדים באמצעות גזירה והשוואה ל-0:

$$\begin{aligned} \hat{\mu}_{+1} &= \frac{\partial l}{\partial \mu_{+1}} = \frac{\partial}{\partial \mu_{+1}} \left(-\frac{1}{2} \sum_{i=1}^k ((x_i/y = 1)^T \Sigma^{-1} (x_i/y = 1) - \right. \\ &\quad \left. (x_i/y = 1)^T \Sigma^{-1} \mu_{+1} - \mu_{+1}^T \Sigma^{-1} (x_i/y = 1) + \mu_{+1}^T \Sigma^{-1} \mu_{+1}) \right) = \\ &= \frac{\partial}{\partial \mu_{+1}} \sum_{i=1}^k \left(-\frac{1}{2} (x_i/y = 1)^T \Sigma^{-1} (x_i/y = 1) + (x_i/y = 1)^T \Sigma^{-1} \mu_{+1} - \frac{1}{2} \mu_{+1}^T \Sigma^{-1} \mu_{+1} \right) = \\ &= \sum_{i=1}^k ((x_i/y = 1)^T \Sigma^{-1} - \mu_{+1}^T \Sigma^{-1}) \end{aligned}$$

נשווה ל-0 למציאת האומד:

$$\begin{aligned}\sum_{i=1}^k ((x_i/y = 1)^T \Sigma^{-1} - \mu_{+1}^T \Sigma^{-1}) &= 0 \Leftrightarrow \\ \sum_{i=1}^k ((x_i/y = 1)^T - \mu_{+1}^T) &= 0 \Leftrightarrow \\ \sum_{i=1}^k (x_i/y = 1) &= k\mu_{+1} \Leftrightarrow \\ \frac{1}{k} \sum_{i=1}^k (x_i/y = 1) &= \mu_{+1}\end{aligned}$$

לכן האומד $\hat{\mu}_{+1} = \frac{1}{k} \sum_{i=1}^k (x_i/y = 1) = \overline{(X/y = 1)} \in R^d$ הוא MLE של μ_{+1} . כאשר d הוא מימד הפיצ'רים.

באופן דומה נקבל ש- $\hat{\mu}_{-1} = \frac{1}{m-k} \sum_{i=k}^m (x_i/y = -1) = \overline{(X/y = -1)}$

ניתן לסכם: $\hat{\mu}_{y \in \{\pm 1\}} = \frac{\sum_{i=1}^m x_i \mathbb{I}\{y_i=y\}}{\sum_{i=1}^m \mathbb{I}\{y_i=y\}}$ עבור \mathbb{I} פונק' אינדיקטור.

נמצא את $\hat{\Sigma}$ באמצעות גזירה לפי Σ^{-1} נציב את האומדים שממקסמים את μ_{+1}, μ_{-1} :

$$\begin{aligned}\hat{\Sigma}^{-1} &= \frac{\partial l}{\partial \Sigma^{-1}} \left(-\frac{dm}{2} \ln(2\pi) - \frac{m}{2} \ln(|\Sigma|) - \frac{1}{2} \sum_{i=1}^k ((x_i/y = 1) - \overline{(X/y = 1)})^T \Sigma^{-1} ((x_i/y = 1) - \overline{(X/y = 1)}) + \right. \\ &\quad \left. \sum_{j=k+1}^m -\frac{1}{2} ((x_j/y = -1) - \overline{(X/y = -1)})^T \Sigma^{-1} ((x_j/y = -1) - \overline{(X/y = -1)}) \right) \\ &= \frac{m}{2} \Sigma + \frac{\partial l}{\partial \Sigma^{-1}} \left[-\frac{1}{2} \sum_{i=1}^k ((x_i/y = 1) - \overline{(X/y = 1)})^T \Sigma^{-1} ((x_i/y = 1) - \overline{(X/y = 1)}) + \right. \\ &\quad \left. - \frac{1}{2} \sum_{j=k+1}^m ((x_j/y = -1) - \overline{(X/y = -1)})^T \Sigma^{-1} ((x_j/y = -1) - \overline{(X/y = -1)}) \right] \\ &\stackrel{*}{=} \frac{m}{2} \Sigma - \frac{1}{2} \sum_{i=1}^k ((x_i/y = 1) - \overline{(X/y = 1)}) ((x_i/y = 1) - \overline{(X/y = 1)})^T + \\ &\quad - \frac{1}{2} \sum_{j=k+1}^m ((x_j/y = -1) - \overline{(X/y = -1)}) ((x_j/y = -1) - \overline{(X/y = -1)})^T\end{aligned}$$

* מתקיים $|\Sigma| = |\Sigma^{-1}|^{-1}$ לכן $\ln(|\Sigma|) = -\ln(|\Sigma^{-1}|)$ בנוסף $\frac{\partial l}{\partial \Sigma^{-1}} \ln(|\Sigma|) = \Sigma^{-T}$ נקבל $\frac{\partial l}{\partial \Sigma^{-1}} \ln(|\Sigma|) = \frac{\partial l}{\partial \Sigma^{-1}} - \ln(|\Sigma^{-1}|) = -\Sigma^T \stackrel{\Sigma \text{ Symmetric}}{=} -\Sigma$
 ** מכיוון ש- $x^T A x$ הוא סקלר נוכל לקחת את ה- tr שלו ומתקיים $\frac{\partial}{\partial A} x^T A x = \frac{\partial}{\partial A} tr[x^T A x] = tr[x^T] = x x^T$
 $\frac{\partial}{\partial A} tr[A x x^T] = (x x^T)^T = x x^T$

נשווה ל-0 למציאת האומד:

$$\begin{aligned}
& \frac{m}{2} \Sigma - \frac{1}{2} \sum_{i=1}^k ((x_i/y = 1) - \overline{(X/y = 1)})((x_i/y = 1) - \overline{(X/y = 1)})^T + \\
& - \frac{1}{2} \sum_{j=k}^m ((x_j/y = -1) - \overline{(X/y = -1)})((x_j/y = -1) - \overline{(X/y = -1)})^T = 0 \Leftrightarrow \\
& \hat{\Sigma} = \frac{1}{m} \sum_{i=1}^k [(x_i/y = 1) - \overline{(X/y = 1)})((x_i/y = 1) - \overline{(X/y = 1)})^T + \\
& + \sum_{j=k}^m ((x_j/y = -1) - \overline{(X/y = -1)})((x_j/y = -1) - \overline{(X/y = -1)})^T] \\
& \hat{\Sigma} = \frac{1}{m} \sum_{y \in \{\pm 1\}} \sum_{i \in [m]} \sum_{s.t. y_i = y} (x_i - \overline{\mu_y})(x_i - \overline{\mu_y})^T \\
& \hat{Pr}_{y \in \{\pm 1\}}(y) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}\{y_i = y\}
\end{aligned}$$

קיבלנו: $\hat{\Sigma} = \frac{1}{m} \sum_{y \in \{\pm 1\}} \sum_{i \in [m]} \sum_{s.t. y_i = y} (x_i - \overline{\mu_y})(x_i - \overline{\mu_y})^T$
ראינו בתרגול שמתקיים $\hat{Pr}_{y \in \{\pm 1\}}(y) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}\{y_i = y\}$

שאלה 4

הייתי מגדיר את הטעות החמורה יותר כטעות בה סיווגתי מייל שאינו ספאם כ-ספאם, מכיוון שאיבדתי מייל. בהתאם נסווג מייל ספאם כ: *Positive 1*, מייל לא ספאם כ: *Negative -1*. טעות מסוג ראשון, חמורה יותר - *FP* - סיווגתי מייל שאינו ספאם כ-ספאם. (כלומר איבדתי מייל אמיתי). טעות מסוג שני, פחות חמורה - *FN* - סיווגתי מייל ספאם כלא ספאם. (קיבלתי בטעות ספאם).

שאלה 5

נגדיר $v = \begin{pmatrix} w \\ b \end{pmatrix}$ בהתאם לנוסחה ונקבל:

$$\begin{aligned}
& \underset{(w,b)}{\operatorname{argmin}} ||w||^2 \stackrel{*}{=} \\
& s.t. \forall i, y_i (<w, x_i> + b) \geq 1 \\
& \underset{(w,b)}{\operatorname{argmin}} \begin{pmatrix} w \\ b \end{pmatrix}^T I \begin{pmatrix} w \\ b \end{pmatrix} = \\
& s.t. \begin{pmatrix} (y_1 x_1) & y_1 \\ \cdot & \cdot \\ \cdot & \cdot \\ (y_m x_m) & y_m \end{pmatrix} \begin{pmatrix} w \\ b \end{pmatrix} \geq \begin{pmatrix} 1 \\ \cdot \\ \cdot \\ 1 \end{pmatrix} \\
& \underset{(w,b)}{\operatorname{argmin}} \frac{1}{2} \begin{pmatrix} w \\ b \end{pmatrix}^T 2I \begin{pmatrix} w \\ b \end{pmatrix} \\
& s.t. - \begin{pmatrix} (y_1 x_1) & y_1 \\ \cdot & \cdot \\ \cdot & \cdot \\ (y_m x_m) & y_m \end{pmatrix} \begin{pmatrix} w \\ b \end{pmatrix} \leq \begin{pmatrix} 1 \\ \cdot \\ \cdot \\ 1 \end{pmatrix}
\end{aligned}$$

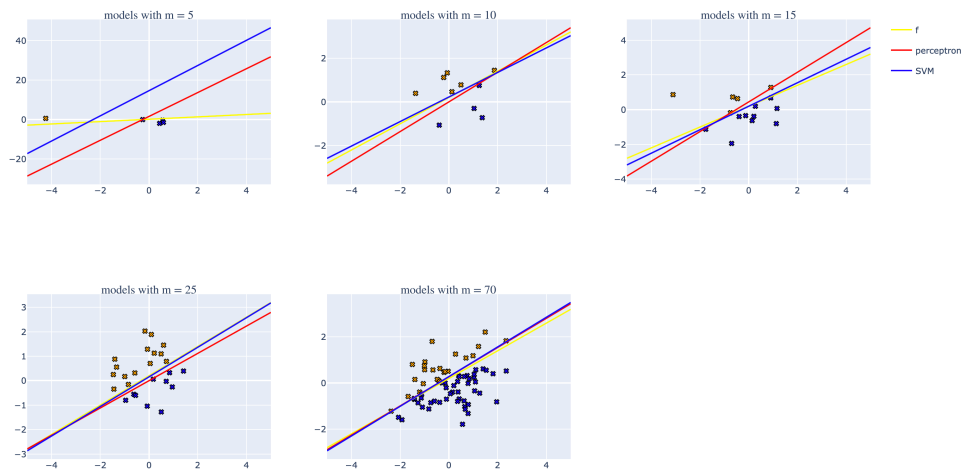
$$A = \begin{pmatrix} (y_1 x_1) & y_1 \\ \vdots & \vdots \\ (y_m x_m) & y_m \end{pmatrix}, Q = 2I, a = 0, d = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \text{ לכן קיבלנו}$$

שאלה 6

ניתן להמיר את הבעיה לצורה האלטרנטיבית מכיוון שמציאת מינ' על פני w, ξ_i שקול, נראה זאת בכך שעבור w מסויים המני' של ξ_i תחת האילוצים שהוא אי שלילי וצריך להתקיים $y_i < w, x_i > \geq 1 - \xi_i$ אבל עבור המקרה נקבל שבמידה ו- $y_i < w, x_i > \geq 1$ ההשמה האופט' עבור ξ_i היא $\xi_i = 0$. $\xi_i = 1 - y_i < w, x_i >$ היא על מנת לשמר את האילוץ בצורה ההדוקה ביותר. לכן נקבל שהצורה האלטרנטיבית היא שקולה לצורה המקורית מכיוון שפונק' l^{hing} מביאה למינ' את ξ_i תחת האילוצים שלו.

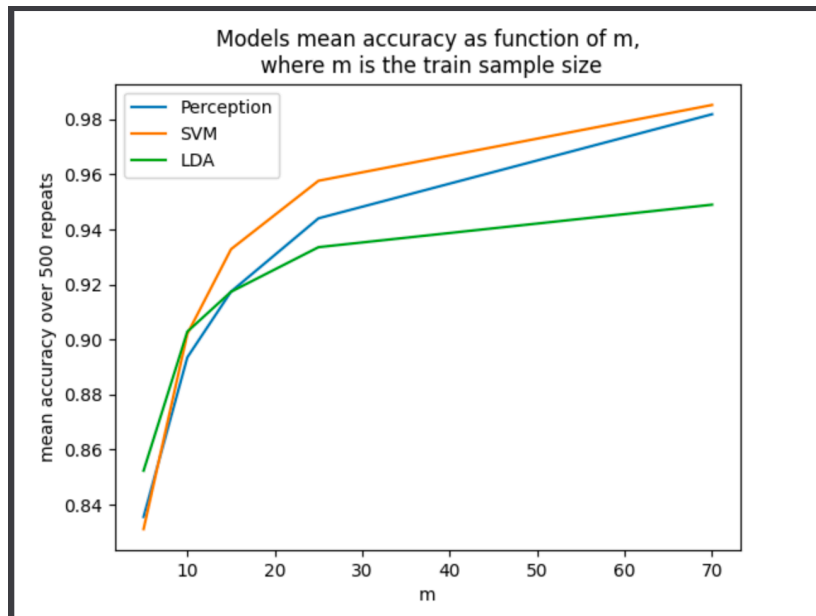
שאלה 9

תרשים פיזור נקודות מוגרלות עבור מדגם אימון בגודל m , יחד איתו על-המישור המתאים לכל מודל.



שאלה 10

התרשים:

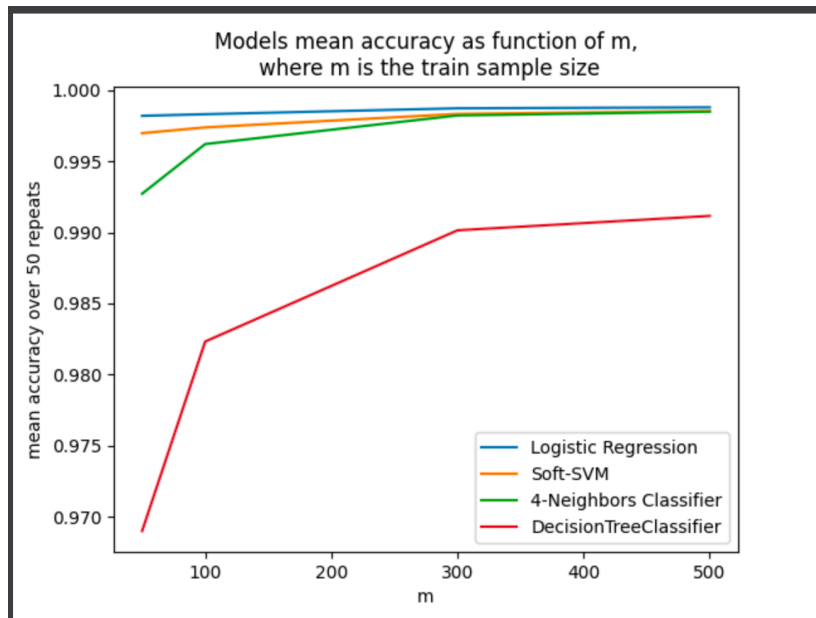


שאלה 11

ניתן לראות ש- SVM הצליח יותר, זאת מכיוון שהוא מתאים את עצמו בצורה טובה יותר כאשר ממקסם את ה- $margin$ בשונה מה- $Perceptron$ שמתאים את הקו הראשון שמצליח לחלק את הדאטה. ה- LDA הצליח פחות משניהם מכיוון שמניח שהדאטה שמייצר את 2 הלייבל מגיע מהתפלגות נורמלית שונה כל אחד, אך במקרה זה הוא מגיע מהתפלגות יחידה.

שאלה 14

התרשים:



כאשר ארבעת המודלים מומשו בעזרת ספריית *sklearn*,
 עבור *Soft – SVM* בחרתי פרמטר רגולציה 1,
 עבור *kNeighborsClassifier* בחרתי ב-4 שכנים,
 עבור *DecisionTreeClassifier* בחרתי בעומק 4 מקסימלי.
 זמני הריצה הם:

עבור *LogisticRegression*: 5.23 שניות

עבור *SVC*: 9.51 שניות

עבור *NeighborsClassifier* – 4: 17.60 שניות

עבור *DecisionTreeClassifier*: 1.87 שניות

ניתן לראות שלעץ עם עומק 4 לקח הכי מעט זמן מכיוון שהאימון שלו מהיר וגם הסיווג בעומק שכזה.
 לעומת זאת מודל *NeighborsClassifier* – 4 לקח הכי הרבה זמן וזאת מכיוון שתהליך סיווג כל תצפית אינו פשוט בהשוואה לאחרים.
 ל-*SVC* לקח יותר זמן בהשוואה ל-*LogisticRegression*, יתכן שבעקבות תהליך האימון שדורש מציאת על-מישור.