

תרגיל 3 - IML - מאי ביבי

שאלה 1:

$$\forall x \in \mathcal{X} \quad h_{\mathcal{D}}(x) = \begin{cases} +1 & \Pr(y = 1 | x) \geq \frac{1}{2} \\ -1 & \text{otherwise} \end{cases}$$

נרצה להראות

$$h_{\mathcal{D}}(x) = \operatorname{argmax}_{y \in \{\pm 1\}} \Pr(x | y) \Pr(y)$$

מכלל בייס מתקיים:

$$\operatorname{argmax}_{y \in \{\pm 1\}} \Pr(x | y) \Pr(y) = \operatorname{argmax}_{y \in \{\pm 1\}} \Pr(y | x) \Pr(x)$$

$$= \operatorname{argmax}_y \{ \Pr(y = 1 | x) \Pr(x), \Pr(y = -1 | x) \Pr(x) \} \quad (*)$$

כיוון שהמאורעות $y = -1, y = 1$ זרים ומהווים חלוקה של מרחב ההסתברות, מנוחסאת ההסתברות השלמה

$$\Pr(y = 1 | x) \Pr(x) + \Pr(y = -1 | x) \Pr(x) = \Pr(x)$$

כיוון ש- $\Pr(x) > 0$ (אחרת $\Pr(y | x)$ לא מוגדר), נחלק בו:

$$\implies \Pr(y = 1 | x) + \Pr(y = -1 | x) = 1$$

ולכן בהכרח או $\Pr(y = 1 | x) \geq \frac{1}{2}$ או $\Pr(y = -1 | x) \geq \frac{1}{2}$. כלומר לעשות \max במקרה הזה שקול לבדיקה $\geq \frac{1}{2}$, ולכן נקבל

$$(*) = \begin{cases} 1 & \Pr(y = 1 | x) \geq \frac{1}{2} \\ -1 & \text{else} \end{cases} = h_{\mathcal{D}}(x)$$

שאלה 2:

פונקצית הצפיפות:

$$f(\mathbf{x} | y) = \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu_y)^\top \Sigma^{-1} (\mathbf{x} - \mu_y) \right\}$$

נרצה להראות שאם היינו יודעים את μ_{-1}, μ_{+1} אז

$$h_{\mathcal{D}}(x) = \operatorname{argmax}_{y \in \{\pm 1\}} \delta_y(\mathbf{x})$$

כאשר

$$\delta_y(\mathbf{x}) = \mathbf{x}^\top \Sigma^{-1} \mu_y - \frac{1}{2} \mu_y^\top \Sigma^{-1} \mu_y + \ln \Pr(y) \quad y \in \{\pm 1\}$$

מהשאלה הקודמת אנחנו יודעים כי $h_{\mathcal{D}}(\mathbf{x}) = \operatorname{argmax}_{y \in \{\pm 1\}} \Pr(\mathbf{x} | y) \Pr(y)$

$$h_{\mathcal{D}}(\mathbf{x}) = \operatorname{argmax}_{y \in \{\pm 1\}} \Pr(\mathbf{x} | y) \Pr(y) \stackrel{1}{=} \operatorname{argmax}_{y \in \{\pm 1\}} f(\mathbf{x} | y) \Pr(y) \stackrel{2}{=} \operatorname{argmax}_{y \in \{\pm 1\}} \ln(f(\mathbf{x} | y) \Pr(y))$$

$$\begin{aligned}
&= \operatorname{argmax}_{y \in \{\pm 1\}} \ln f(\mathbf{x} | y) + \ln \Pr(y) \\
&= \operatorname{argmax}_{y \in \{\pm 1\}} \ln \left(\frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu_y)^\top \Sigma^{-1} (\mathbf{x} - \mu_y) \right\} \right) + \ln \Pr(y) \\
&= \operatorname{argmax}_{y \in \{\pm 1\}} \ln \left(\frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \right) - \frac{1}{2} (\mathbf{x} - \mu_y)^\top \Sigma^{-1} (\mathbf{x} - \mu_y) + \ln \Pr(y) \\
&= \operatorname{argmax}_{y \in \{\pm 1\}} -\ln \sqrt{(2\pi)^d \det(\Sigma)} - \frac{1}{2} (\mathbf{x} - \mu_y)^\top \Sigma^{-1} (\mathbf{x} - \mu_y) + \ln \Pr(y) \\
&\stackrel{3}{=} \operatorname{argmax}_{y \in \{\pm 1\}} -\frac{1}{2} (\mathbf{x} - \mu_y)^\top \Sigma^{-1} (\mathbf{x} - \mu_y) + \ln \Pr(y) \\
&\stackrel{4}{=} \operatorname{argmax}_{y \in \{\pm 1\}} -\frac{1}{2} (\mathbf{x}^\top - \mu_y^\top) \Sigma^{-1} (\mathbf{x} - \mu_y) + \ln \Pr(y) \\
&= \operatorname{argmax}_{y \in \{\pm 1\}} -\frac{1}{2} \mathbf{x}^\top \Sigma^{-1} \mathbf{x} + \frac{1}{2} \mathbf{x}^\top \Sigma^{-1} \mu_y + \frac{1}{2} \mu_y^\top \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_y^\top \Sigma^{-1} \mu_y + \ln \Pr(y) \\
&\stackrel{5}{=} \operatorname{argmax}_{y \in \{\pm 1\}} \frac{1}{2} \mathbf{x}^\top \Sigma^{-1} \mu_y + \frac{1}{2} \mu_y^\top \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_y^\top \Sigma^{-1} \mu_y + \ln \Pr(y) \\
&\stackrel{6}{=} \operatorname{argmax}_{y \in \{\pm 1\}} \mathbf{x}^\top \Sigma^{-1} \mu_y - \frac{1}{2} \mu_y^\top \Sigma^{-1} \mu_y + \ln \Pr(y) = \operatorname{argmax}_{y \in \{\pm 1\}} \delta_y(\mathbf{x})
\end{aligned}$$

כאשר:

1 ממונוטוניות האינטגרל

2 ממונוטוניות \ln

3 $-\ln \sqrt{(2\pi)^d \det(\Sigma)}$ קבוע ולא משפיע על argmax

4 $(\mathbf{x} - \mu_y)^\top = (\mathbf{x}^\top - \mu_y^\top)$ מתכונות transpose

5 $-\frac{1}{2} \mathbf{x}^\top \Sigma^{-1} \mathbf{x}$ קבוע ולא משפיע על argmax

6 $\mathbf{x}^\top \Sigma^{-1} \mu_y = \mu_y^\top \Sigma^{-1} \mathbf{x}$ - לכתוב הסבר אחר כך

שאלה 3:

נשתמש באומדים על מנת להעריך את ההסתברויות: (כאשר $\mathbf{1}_{y_i=y}$ הוא המציין של המאורע $\{y_i = y\}$ את $\Pr(y)$ נקבע כאחוז הפעמים שהוא מופיע ב- y_1, \dots, y_m)

$$\Pr(y) = \frac{1}{m} \sum_{i=1}^m \mathbf{1}_{y_i=y}$$

עבור התוחלות של $y | \mathbf{x}$ נחשב את ממוצע כל הדגימות המתויגות y (אומד מוכר מקורס הסתברות):

$$\mu_y = \frac{\sum_{i=1}^m \mathbf{x}_i \mathbf{1}_{y_i=y}}{\sum_{i=1}^m \mathbf{1}_{y_i=y}}$$

עבור מטריצת השונות, נשמש באומד המוכר הבא:

$$\Sigma = \frac{1}{m} \sum_{y \in \{\pm 1\}} \sum_{i \in [m]: y_i=y} (\mathbf{x}_i - \mu_y)^\top (\mathbf{x}_i - \mu_y)$$

שאלה 4:

יש לנו שתי שגיאות: סיווג אימייל ספאם כלא ספאם וסיווג אימייל לא ספאם כספאם. השגיאה היותר חמורה היא סיווג אימייל לא ספאם כספאם - כי אז האימייל החשוב לא יקרא. סיווג אימייל ספאם כלא ספאם הוא פחות חמור כי בסך הכל יגיע אימייל ספאם ליעד. לכן ה-negative יהיה not-spam, ואז ה-false-negative תהיה השגיאה הפחות חמורה, וה-positive יהיה ה-spam, כך שה-false-positive חמור יותר.

שאלה 5:

להסביר

$$\begin{aligned}
 & \arg \min_{(\mathbf{w}, b)} \|\mathbf{w}\|^2 \\
 \text{s.t. } & \forall i, y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 \\
 & \arg \min_{(\mathbf{w}, b)} \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix}^T \mathbf{I} \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix} \\
 = & \begin{pmatrix} (y_1 x_1) & y_1 \\ \vdots & \vdots \\ (y_m x_m) & y_m \end{pmatrix} \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix} \geq \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \\
 \text{s.t. } & \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix} + \mathbf{0}^T \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix} \\
 & \arg \min_{(\mathbf{w}, b)} \frac{1}{2} \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix}^T [2 \cdot \mathbf{I}] \begin{pmatrix} (y_1 x_1) & y_1 \\ \vdots & \vdots \\ (y_m x_m) & y_m \end{pmatrix} \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix} \leq \begin{pmatrix} -1 \\ \vdots \\ -1 \end{pmatrix}
 \end{aligned}$$

שאלה 6:

נרצה להראות שהבעיה הבאה

$$\arg \min_{\mathbf{w}} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \ell^{\text{hinge}}(y_i \langle \mathbf{w}, \mathbf{x}_i \rangle),$$

$$\text{where } \ell^{\text{hinge}}(a) = \max\{0, 1 - a\}$$

שקולה לבעיית ה-SVM soft כפי שהגדרנו:

$$\arg \min_{\mathbf{w}, \{\xi_i\}} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i \text{ s.t. } \forall i, y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1 - \xi_i \text{ and } \xi_i \geq 0$$

נגדיר את ξ_i עבור $i = 1, \dots, m$ להיות:

$$\xi_i := \begin{cases} 0 & y_i \langle \mathbf{x}_i, \mathbf{w} \rangle > 1 \\ 1 - y_i \langle \mathbf{x}_i, \mathbf{w} \rangle & \text{otherwise} \end{cases}$$

נראה כי עומד באילוצים $\xi_i \geq 0$ ו- $\xi_i \geq 1 - y_i \langle \mathbf{x}_i, \mathbf{w} \rangle \iff y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1 - \xi_i$ אם $y_i \langle \mathbf{x}_i, \mathbf{w} \rangle > 1$ אז $\xi_i = 0$ וכן מתקיים $1 - y_i \langle \mathbf{x}_i, \mathbf{w} \rangle < 0$ ולכן גם האילוץ השני מתקיים. אחרת, $1 - y_i \langle \mathbf{x}_i, \mathbf{w} \rangle \geq 0$ ואנו מגדירים $\xi_i = 1 - y_i \langle \mathbf{x}_i, \mathbf{w} \rangle$ וזה בפרט עומד באילוץ השני. נשים לב שהגדרנו את ξ_i להיות המספר המינימלי שעומד באילוצים - אם $y_i \langle \mathbf{x}_i, \mathbf{w} \rangle > 1$ הוא המינימלי ש- ξ_i יכול להיות מהאילוץ $\xi_i \geq 0$, ואחרת, מהאילוץ השני המספר המינימלי עבור ξ_i הוא $1 - y_i \langle \mathbf{x}_i, \mathbf{w} \rangle$.

כעת נראה כי $\ell^{\text{hinge}}(y_i \langle \mathbf{x}_i, \mathbf{w} \rangle) = \xi_i$

$$\ell^{\text{hinge}}(y_i \langle \mathbf{x}_i, \mathbf{w} \rangle) = \max\{0, 1 - y_i \langle \mathbf{x}_i, \mathbf{w} \rangle\} = \begin{cases} 0 & y_i \langle \mathbf{x}_i, \mathbf{w} \rangle > 1 \\ 1 - y_i \langle \mathbf{x}_i, \mathbf{w} \rangle & \text{otherwise} \end{cases} = \xi_i$$

ולכן

$$\frac{1}{m} \sum_{i=1}^m \ell^{\text{hinge}}(y_i \langle \mathbf{w}, \mathbf{x}_i \rangle) = \frac{1}{m} \sum_{i=1}^m \xi_i$$

הגדרנו ξ_i מינימליים שעומדים באילוצים, וזה מתאים לבעיה המקורית בה אנו רוצים למזער אותם. ממה שהראנו קודם מתקיים

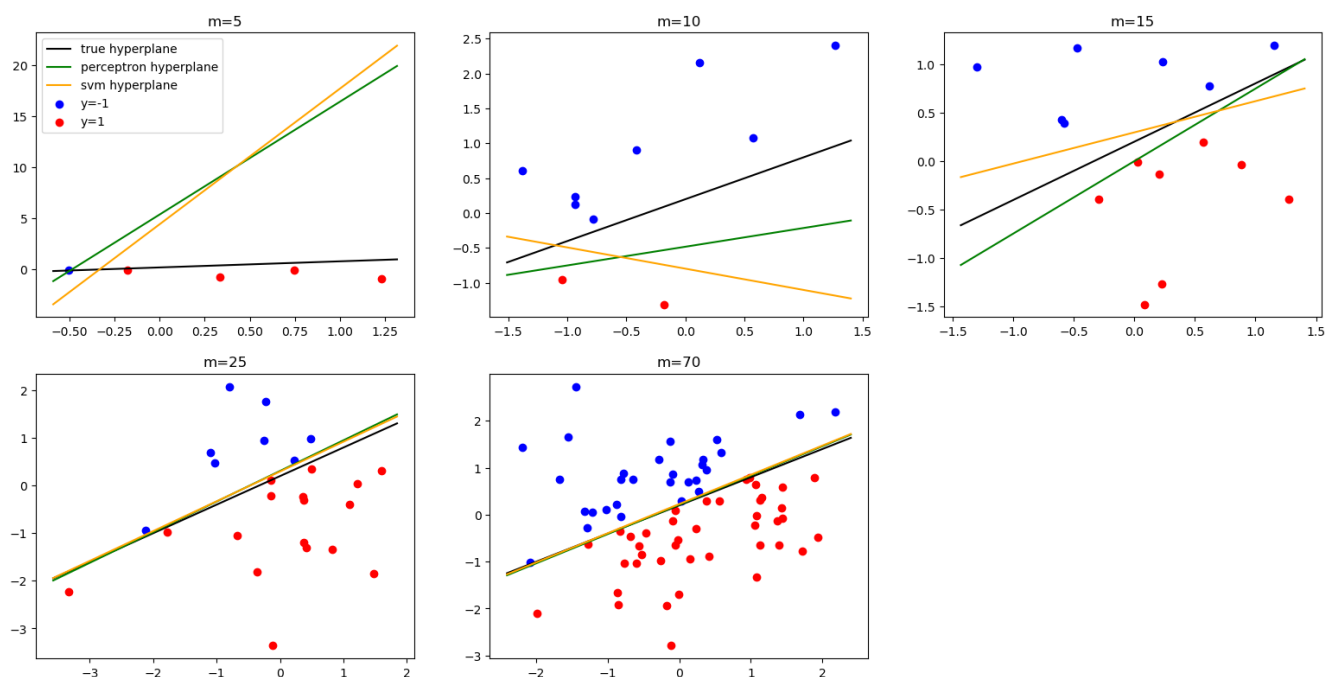
$$\arg \min_{\mathbf{w}, \{\xi_i\}} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i = \arg \min_{\mathbf{w}} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \ell^{\text{hinge}}(y_i \langle \mathbf{w}, \mathbf{x}_i \rangle)$$

ולכן הבעיות שקולות.

שאלה 9:

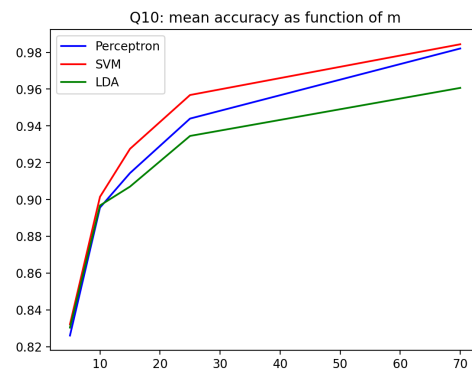
גרף של:

Q9: True vs. Perceptron vs. SVM hyperplanes



שאלה 10:

גרף של mean accuracy כפונקציה של m עבור שלושת המסווגים LDA, SVM, Perceptron:

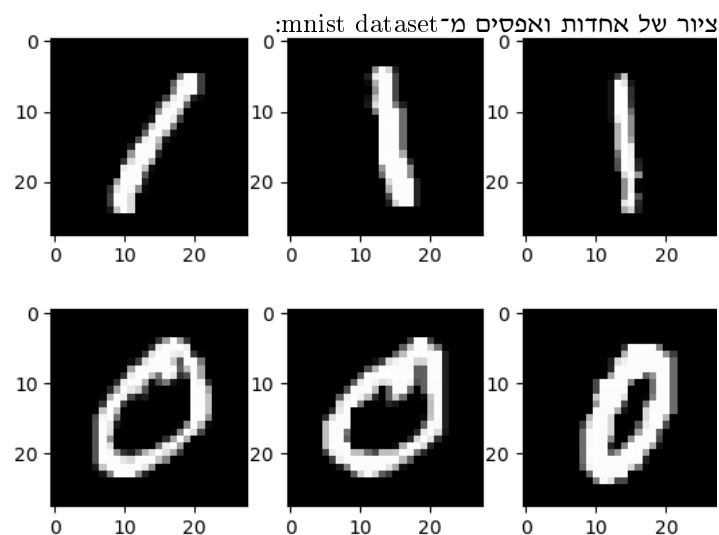


ניתן לראות שה-SVM הוא בעל ה-mean accuracy הכי גבוה, וה-LDA הכי נמוך

שאלה 11:

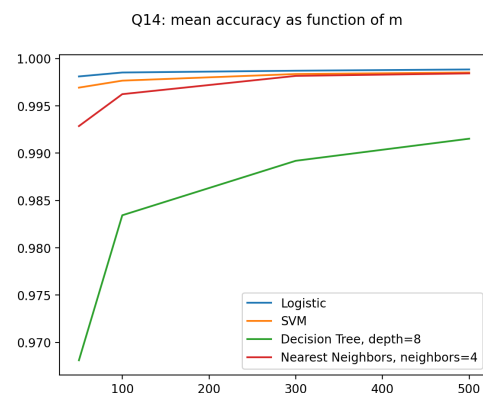
כיוון ש-LDA מניח הנחה שגויה לגבי התפלגות הדגימות (הוא מניח ש- X, y באים מהתפלגות משותפת ותלויים אחד בשני, בעוד שדגמנו רק את X באופן בלתי תלוי וגאוסייני, ואינו תלוי ב- y) רואים שהוא מסווג הכי פחות טוב (accuracy נמוך יותר). SVM ו-Perceptron לא מניחים הנחה שגויה כזו ולכן מסווגים טוב יותר (accuracy גבוה יותר). ניתן לראות ש-SVM סיווג טוב יותר וזאת כי הוא בוחר על מישור (קו ההפרדה) עם שול מקסימלי, בעוד שה-Perceptron בוחר על מישור כלשהו שעובד.

שאלה 12:



שאלה 14:

א:



זמני הריצה:

| | m=50 | m=100 | m=300 | m=500 |
|--------------------------------|----------|----------|----------|----------|
| Logistic | 0.007669 | 0.008070 | 0.011957 | 0.015337 |
| SVM | 0.039158 | 0.047437 | 0.069210 | 0.083393 |
| Decision Tree, depth=8 | 0.003963 | 0.004811 | 0.010244 | 0.017135 |
| Nearest Neighbors, neighbors=4 | 0.180960 | 0.298545 | 0.854777 | 1.470554 |