

מבוא ללמידת מכונה - תרגיל 3

מגיש: יואב לוי 314963257

Bayes Optimal and LDA

שאלה 1

$$\operatorname{argmax}_{y \in \{\pm 1\}} \Pr(x|y) \cdot \Pr(y) \stackrel{\text{Bayes Thm.}}{=} \operatorname{argmax}_{y \in \{\pm 1\}} \frac{\Pr(x \cap y)}{\Pr(y)} \cdot \Pr(y) =$$

$$\stackrel{\Pr(x) \neq 0}{=} \operatorname{argmax}_{y \in \{\pm 1\}} \frac{\Pr(x \cap y)}{\Pr(x)} \cdot \Pr(x) = \operatorname{argmax}_{y \in \{\pm 1\}} \Pr(y|x) \cdot \Pr(x) =$$

$$\stackrel{*}{=} \operatorname{argmax}_{y \in \{\pm 1\}} \Pr(y|x) = \begin{cases} +1 & \Pr(y = 1|x) \geq \frac{1}{2} \\ -1 & \text{otherwise} \end{cases} = h_{\mathcal{D}}(x)$$

(*) מכיוון שהסתברות $\Pr(x)$ היא חיובית ואינה תלויה ב- y שעברו argmax מוגדר, ניתן להתעלם מהמכפלה (בו הערה: נשים לב למקרי קצה בהם אם $\Pr(y = 1) = 0$ אז $\Pr(y = -1) = 1$ ואז $\operatorname{Supp}(y) = \{1, -1\}$ מכיוון ש- $\{1, -1\}$ ואז

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

וזה נכון לכל x , בנוסף אם $\Pr(y = -1) = 0$ אז בלי הגבלת הכלליות נקבל אותו דבר כמו מקודם רק הפוך.

שאלה 2

$$\begin{aligned} h_{\mathcal{D}}(x) &= \operatorname{argmax}_{y \in \{\pm 1\}} \{\Pr(\mathbf{x}|y) \Pr(y)\} \\ &= \operatorname{argmax}_{y \in \{\pm 1\}} \left\{ \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu_y)^\top \Sigma^{-1} (\mathbf{x} - \mu_y) \right\} \cdot \Pr(y) \right\} \\ &\stackrel{(\text{removed constants})}{=} \operatorname{argmax}_{y \in \{\pm 1\}} \left\{ \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu_y)^\top \Sigma^{-1} (\mathbf{x} - \mu_y) \right\} \cdot \Pr(y) \right\} \\ &\stackrel{(\text{Log is monotone increasing})}{=} \operatorname{argmax}_{y \in \{\pm 1\}} \left\{ \ln \left(\exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu_y)^\top \Sigma^{-1} (\mathbf{x} - \mu_y) \right\} \cdot \Pr(y) \right) \right\} \\ &= \operatorname{argmax}_{y \in \{\pm 1\}} \left\{ \left(-\frac{1}{2} (\mathbf{x} - \mu_y)^\top \Sigma^{-1} (\mathbf{x} - \mu_y) \right) + \ln(\Pr(y)) \right\} \\ &= \operatorname{argmax}_{y \in \{\pm 1\}} \left\{ \left(\left(-\frac{1}{2} \mathbf{x}^\top \Sigma^{-1} + \frac{1}{2} \mu_y^\top \Sigma^{-1} \right) (\mathbf{x} - \mu_y) \right) + \ln(\Pr(y)) \right\} \\ &= \operatorname{argmax}_{y \in \{\pm 1\}} \left\{ \left(-\frac{1}{2} \mathbf{x}^\top \Sigma^{-1} \mathbf{x} + \frac{1}{2} \mathbf{x}^\top \Sigma^{-1} \mu_y + \frac{1}{2} \mu_y^\top \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_y^\top \Sigma^{-1} \mu_y \right) + \ln(\Pr(y)) \right\} \\ &\stackrel{(-\frac{1}{2} \mathbf{x}^\top \Sigma^{-1} \mathbf{x} \text{ is not depend on } y)}{=} \operatorname{argmax}_{y \in \{\pm 1\}} \left\{ \left(\mu_y^\top \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_y^\top \Sigma^{-1} \mu_y \right) + \ln(\Pr(y)) \right\} \\ &= \operatorname{argmax}_{y \in \{\pm 1\}} \delta_y(x) \end{aligned}$$

שאלה 3

$$\mu_{+1} = \frac{1}{m} \sum_{i=1}^m x_i \cdot (1_{y_i}) \quad \text{where } 1_{y_i} = \begin{cases} 1 & \text{if } y_i = 1 \\ 0 & \text{otherwise} \end{cases}$$

$$\mu_{-1} = \frac{1}{m} \sum_{i=1}^m x_i \cdot (-1_{y_i}) \quad \text{where } -1_{y_i} = \begin{cases} 1 & \text{if } y_i = -1 \\ 0 & \text{otherwise} \end{cases}$$

$$\Sigma = \frac{X_{centered} X_{centered}^T}{m-1} \quad \text{where } X_{centered}^T = \begin{pmatrix} - & x_1 & - \\ & \vdots & \\ - & x_m & - \end{pmatrix} - \begin{pmatrix} \frac{1}{m} \sum_{i=1}^m x_i^1 & \dots & \frac{1}{m} \sum_{i=1}^m x_i^d \\ \vdots & \ddots & \vdots \\ \frac{1}{m} \sum_{i=1}^m x_i^1 & \dots & \frac{1}{m} \sum_{i=1}^m x_i^d \end{pmatrix}$$

$$\Pr(y = 1) = \frac{\sum_{i=1}^m 1_{y_i}}{m} \quad \text{where } 1_{y_i} = \begin{cases} 1 & \text{if } y_i = 1 \\ 0 & \text{otherwise} \end{cases}$$

$$\Pr(y = -1) = \frac{\sum_{i=1}^m -1_{y_i}}{m} \quad \text{where } -1_{y_i} = \begin{cases} 1 & \text{if } y_i = -1 \\ 0 & \text{otherwise} \end{cases}$$

הערה: נשים לב שמטריצת ה- COV מוגדרת בעזרת $X_{centered}^T$ כאשר זאת מוגדרת בעזרת מטריצת של m דגימות על d פיצ'רים, פחות מטריצה שבכל עמודה יש ממוצע של פיצ'ר לפי מספר העמודה (כלומר בעמודה הראשונה יופיעו בכל השורות הממוצע על הפיצ'ר הראשון).

Spam

שאלה 4

- הטעויות שהמסווג שלי עלול לעשות הן, [לסווג אימייל כספאם כאשר הוא אינו ספאם](#), [וההפך לסווג אימייל כלא ספאם כאשר הוא למעשה ספאם](#).

- הטעות שלא נרצה לעשות היא [לסווג אימייל כספאם כאשר הוא אינו ספאם](#).

$$\text{spam} = 1$$

$$\text{not-spam} = -1$$

SVM-Formulation

שאלה 5

$$\operatorname{argmin}_{v \in \mathbb{R}^n} \left(\frac{1}{2} v^\top \cdot 2 \cdot I \cdot v + \overbrace{\vec{0}^\top}^{\vec{a}^\top} v \right)$$

$$s.t. \quad Av \leq d$$

כאשר,

$$Q = 2 \cdot I = \begin{pmatrix} 2 & 0 & \dots & 0 \\ \vdots & 2 & & \vdots \\ 0 & 0 & \ddots & 0 \\ 0 & \dots & 0 & 2 \end{pmatrix} \in \mathbb{R}^{n \times n}$$

$$A = \begin{pmatrix} \overbrace{\begin{pmatrix} -y_1 & \overbrace{(-y_1 x_1)^\top}^{n-1} \\ \vdots & \vdots \\ y_m & (-y_m x_m)^\top \end{pmatrix}}^n \end{pmatrix} \in \mathbb{R}^{m \times n}$$

$$\mathbf{a} = \vec{0} \in \mathbb{R}^n$$

$$\mathbf{d} = -\vec{1} \in \mathbb{R}^m$$

הערות: שורה i במטריצה A היא $row_i = [-y_1, -y_1 x_i^1, -y_1 x_i^2 - y_1 x_i^3, \dots, -y_1 x_i^{n-1}] \in \mathbb{R}_{row}^n$ כאשר x_i^j הוא ה-feature ה- j בדגימה ה- i . בנוסף $I \in \mathbb{R}^{n \times n}$ היא מטריצת הזהות.

שאלה 6

ראשית נשים לב שהבעיות שקולות פרט לכך שבבעיה הראשונה אנחנו ממעעים על ξ_i לכל $i \in [m]$ כאשר יש לנו את התנאי $\forall i, y_i \langle w, x_i \rangle \geq 1 - \xi_i$ and $\xi_i \geq 0$ (נסמנו λ) ולכן כדי להראות שהבעיות שקולות יש להראות שלכל i הבחירה של ξ_i עבור אותו w_i (כאשר ξ_i יגרור לעמידה בתנאי λ), תהיה שקולה להפעלת פונקציית ה- ℓ^{hinge} על $y_i \langle w, x_i \rangle$. קרי שמתקיים $\forall i \in [m], \ell^{hinge}(y_i \langle w, x_i \rangle) = \xi_i$ וגם שזה גורר לכך שעומדים בתנאי λ . כדי להראות זאת אחלק לשני מקרים: (עבור $i \in [m]$)

- במידה ומתקיים $y_i \langle w, x_i \rangle \geq 1$ נשים לב שלפי הגדרת פונקציית ה- ℓ^{hinge} אנחנו נקבל

$$\ell^{hinge}(y_i \langle w, x_i \rangle) = \max\{0, 1 - y_i \langle w, x_i \rangle\} \stackrel{y_i \langle w, x_i \rangle \geq 1}{=} 0$$

כעת, נתבונן בבעיה הראשונה ונשים לב שבמקרה הנ"ל אם נחבר $\xi_i > 0$ אז $\xi_i < 0$ וזה לא יתקיים. אבל מכיוון שאנחנו מחפשים להגיע לערך המינימלי של $\frac{\lambda}{2} \|w\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i$ אז ברור שעדיף ש- $\xi_i = 0$, בנוסף עבור אותו i גם יתקיים התנאי λ ,

$$y_i \langle w, x_i \rangle \geq 1 - \xi_i = 1 - 0 = 1 \text{ and } \xi_i \geq 0$$

כלומר במקרה זה פונקציית ה- ℓ^{hinge} תהיה שקולה לבחירת $\xi_i = 0$ עבור ה- i הנ"ל, וגם אנחנו עומדים בתנאי λ .

- במידה ומתקיים $y_i \langle w, x_i \rangle < 1$ נשים לב שנקבל

$$\ell^{hinge}(y_i \langle w, x_i \rangle) = \max\{0, 1 - y_i \langle w, x_i \rangle\} = 1 - y_i \langle w, x_i \rangle$$

כעת, אם נבחר ב- $\xi_i = 1 - y_i \langle w, x_i \rangle$ אז תנאי λ יתקיים כי,

$$y_i \langle w, x_i \rangle = y_i \langle w, x_i \rangle = 1 - \overbrace{1 - y_i \langle w, x_i \rangle}^{-\xi_i} = 1 - \xi_i \quad (1)$$

$$\text{and } \xi_i = 1 - y_i \langle w, x_i \rangle \stackrel{y_i \langle w, x_i \rangle < 1}{\iff} \stackrel{0 < 1 - y_i \langle w, x_i \rangle}{>} 0 \quad (2)$$

בנוסף מכיוון שאנחנו רוצים להקטין כמה שיותר את פונקציית המטרה שלנו בבעיית ה-Soft-SVM, אז נרצה לבחור ξ_i קטן ככל שאפשר

אבל נניח בשלילה שיש כזה $\xi_i = 1 - y_i \langle w, x_i \rangle < \xi'$ אז נקבל

$$y_i \langle w, x_i \rangle = y_i \langle w, x_i \rangle = 1 - \overbrace{1 + y_i \langle w, x_i \rangle}^{-\xi_i} = 1 - \xi_i < 1 - \xi'$$

בניגוד לתנאי 1 $(\forall i, y_i \langle w, x_i \rangle \geq 1 - \xi_i)$, ולכן בהכרח מצאנו את ה- ξ הנכון עבור ה- i הנ"ל.

כלומר למעשה הוכחנו ששני הבעיות שקולות.

■

Implementation and simulation-comparison of different classifiers

שאלה 7

This is a coding question, the code is provided in the 'models.py' python file.

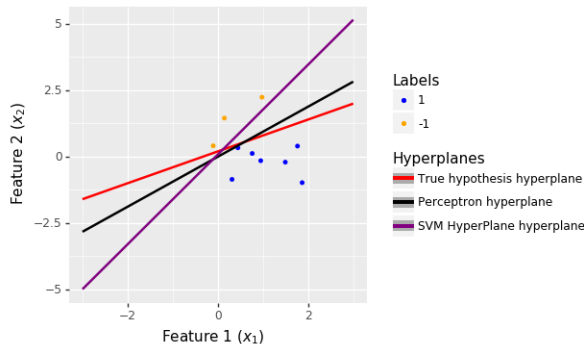
שאלה 8

This is a coding question, the code is provided in the 'comparison.py' python file.

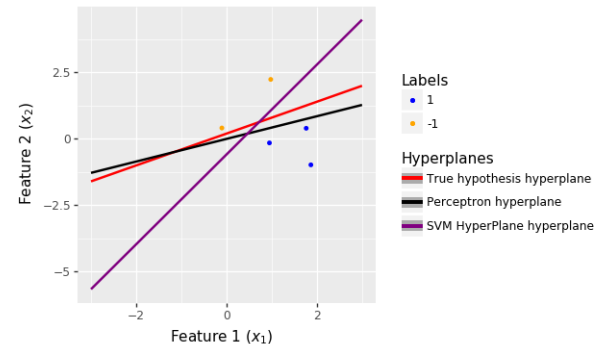
שאלה 9

הגרפים הבאים מציגים את העל-מישורים שנבחרו ע"י האלגוריתמים SVM, Perceptron, כאשר הפונקציה האמיתית שמתארת את העל מישור היא $\left\langle \begin{pmatrix} 0.3 \\ -0.5 \end{pmatrix}, x \right\rangle + 0.1 = 0$, והשוני בכל גרף הוא מספר הדגימות (samples) שעליו כל אלגוריתם קלסיפיקציה אומן.

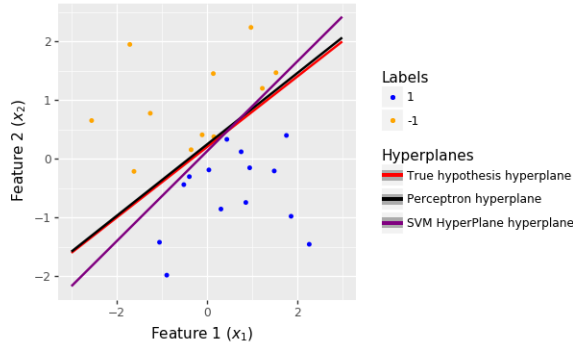
Question 9: Testing SVM algorithm VS Perceptron (number of samples:10)



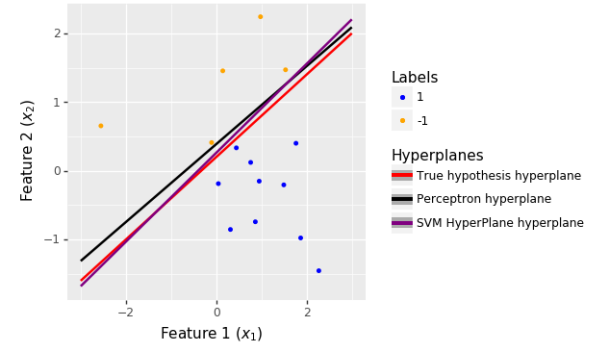
Question 9: Testing SVM algorithm VS Perceptron (number of samples:5)



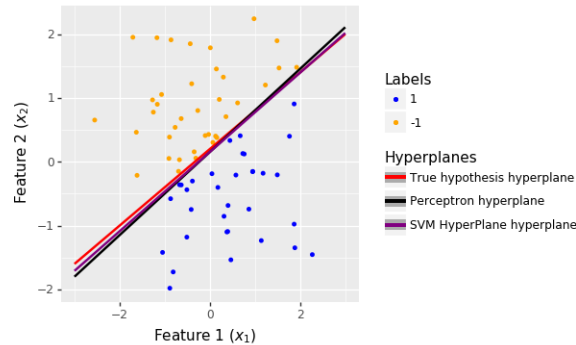
Question 9: Testing SVM algorithm VS Perceptron (number of samples:25)



Question 9: Testing SVM algorithm VS Perceptron (number of samples:15)



Question 9: Testing SVM algorithm VS Perceptron (number of samples:70)

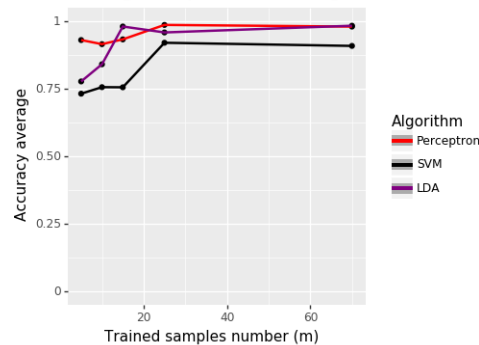


*בנוסף הקוד של הגרפים ויצירתם נמצא בקובץ 'comparison.py'.

שאלה 10

בשאלה זו נדרשנו לבדוק את ה-Accuracy של האלגוריתמים SVM, LDA, Perceptron, עבור מספר דגימות משתנה שעליהן אימנו את ה-DATA, בנוסף ה-Testing data נעשה על 10,000 דגימות נוספות, ומיצענו על 500 איטרציות עבור כל מספר דגימות שעליו אימנו את ה-DATA.

Question 10: Testing Perceptron VS SVM algorithm VS LDA



*בנוסף הקוד של הגרפים ויצירתם נמצא בקובץ 'comparison.py'.

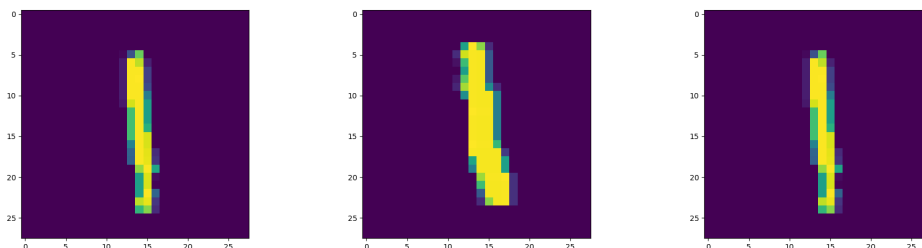
שאלה 11

ניתן לראות בגרף שהמסווג Perceptron היה המוצלח ביותר מבין שלושת המסווגים, וה-SVM גם מאוד קרוב אליו. לדעתי זה בגלל שהפרספטרון מנסה למצוא את ההפרדה הטובה ביותר, בעוד שה-SVM מנסה למצוא את ה-margin הטוב ביותר, נשים לב שכל האלגוריתמים מתכנסים יחסית ככל שאנחנו מאמנים על יותר DATA.

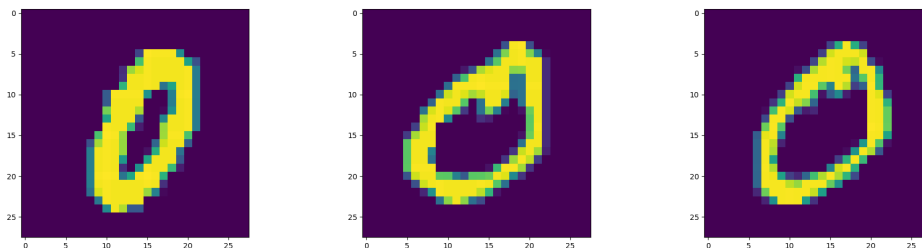
שאלה 12

דגימות מה-DATA SET:

אחדים-



אפסים-



שאלה 13

This is a coding question, the code is provided in the 'mnist_data.py' python file.

שאלה 14

- הגרף הבא מתאר את הדיוק (Accuracy) של כל אחד מהאלגוריתמים למידה שנדרשנו לבדוק. מדדנו את התוצאות על DATA של כתבי יד של מספרים (0 או 1) וציר ה- X הוא מספר הדגימות שדגמנו כדי לאמן את האלגוריתמים. נשים לב שבמקרה של KNN דווקא חלה ירידה ככל שאימנו על יותר וזה מכיוון שהגדרתי את שהאלגוריתם יבדוק על מספר שכנים שהוא שליש ממספר הדגימות, כלומר כשיש יותר דגימות ככה הוא מתחשב ביותר שכנים. דבר זה כנראה מביא לבעיה מסויימת כי הוא מתחיל להחשיב שכנים שכנראה מטעים אותו. בקשר לשאר האלגוריתמים, הם די יציבים וגם ככה רובם קרובים ל-100% רוב הזמן.
- בנוסף זמן הריצה של ה-KNN (כלומר הפרדיקציה\חיזוי) הולך וגדל ככל שמאמנים אותו על יותר דגימות. ובאופן כללי שמתי לב שהזמן של החיזוי גדל ככל שאנחנו עובדים עם אלגוריתמים שאימנו אותם על יותר דאטא (חוץ מה-SVD)

Question 14: Accuracies of LR vs SVM vs DC-Tree vs KNN on MNIST

