

תרגיל 3 - IML - מאי ביבי

שאלה 1:

מסווג Bayes optimal:

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

נרצה להראות

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

מכלל בייס מתקיים:

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

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

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 $\Pr(y = 1 | x) \geq \frac{1}{2}$, ולכן נקבל

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

שאלה 2:

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

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

כאשר:

1 ממונטוניות האינטגרל

2 \ln ממונטוניות

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

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

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

6 $\mathbf{x}^\top \Sigma^{-1} \mu_y = \mu_y^\top \Sigma^{-1} \mathbf{x}$ כי Σ אלכסונית כי המשתנים בלתי תלויים והשוויות המשותפת שלהם 0.

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 $\mathbf{x} | y$ נחשב את ממוצע כל הדגימות המתוגות y :

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

עבור מטריצת השוונות, נשמש באומד המוכר הבא (מויקיפדיה):

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

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 spam, כך שה-false-positive חמור יותר.

שאלה 5:

נכתוב את בעיית ה-Hard-SVM הבאה בצורה הקנונית של QP:

$$\begin{aligned}
 &= \underset{(\mathbf{w}, b)}{\operatorname{argmin}} \|\mathbf{w}\|^2 \quad \text{s.t.} \quad \forall i, y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 \\
 &= \underset{(\mathbf{w}, b)}{\operatorname{argmin}} \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix}^T I \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix} \quad \text{s.t.} \quad \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix} \cdot \begin{pmatrix} -x_1 & -1 \\ \vdots & \vdots \\ -x_m & -1 \end{pmatrix} \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix} \geq \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \\
 &= \underset{(\mathbf{w}, b)}{\operatorname{argmin}} \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix}^T I \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix} \quad \text{s.t.} \quad \begin{pmatrix} y_1 x_1^1 & \cdots & y_1 x_d^1 & 1 \\ \vdots & \ddots & \vdots & \vdots \\ y_m x_1^m & \cdots & y_m x_d^m & 1 \end{pmatrix} \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix} \geq \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \\
 &= \underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix}^T 2 \cdot I \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix} + \vec{0}^T \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix} \quad \text{s.t.} \quad \begin{pmatrix} y_1 x_1^1 & \cdots & y_1 x_d^1 & 1 \\ \vdots & \ddots & \vdots & \vdots \\ y_m x_1^m & \cdots & y_m x_d^m & 1 \end{pmatrix} \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix} \geq \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \\
 &= \underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix}^T 2 \cdot I \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix} + \vec{0}^T \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix} \quad \text{s.t.} \quad - \begin{pmatrix} y_1 x_1^1 & \cdots & y_1 x_d^1 & 1 \\ \vdots & \ddots & \vdots & \vdots \\ y_m x_1^m & \cdots & y_m x_d^m & 1 \end{pmatrix} \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix} \leq \begin{pmatrix} -1 \\ \vdots \\ -1 \end{pmatrix}
 \end{aligned}$$

וקיבלנו שהשורה האחרונה היא בצורת QP קנונית עם

$$\mathbf{v} = \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix}, Q = 2 \cdot I, \mathbf{a} = \vec{0}$$

$$A = - \begin{pmatrix} y_1 x_1^1 & \cdots & y_1 x_d^1 & 1 \\ \vdots & \ddots & \vdots & \vdots \\ y_m x_1^m & \cdots & y_m x_d^m & 1 \end{pmatrix}, \mathbf{d} = \begin{pmatrix} -1 \\ \vdots \\ -1 \end{pmatrix}$$

שאלה 6:

נרצה להראות שהבעיה הבאה

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

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

סקולה לבעיית ה-soft SVM כפי שהגדרנו:

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

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 ξ_i 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

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

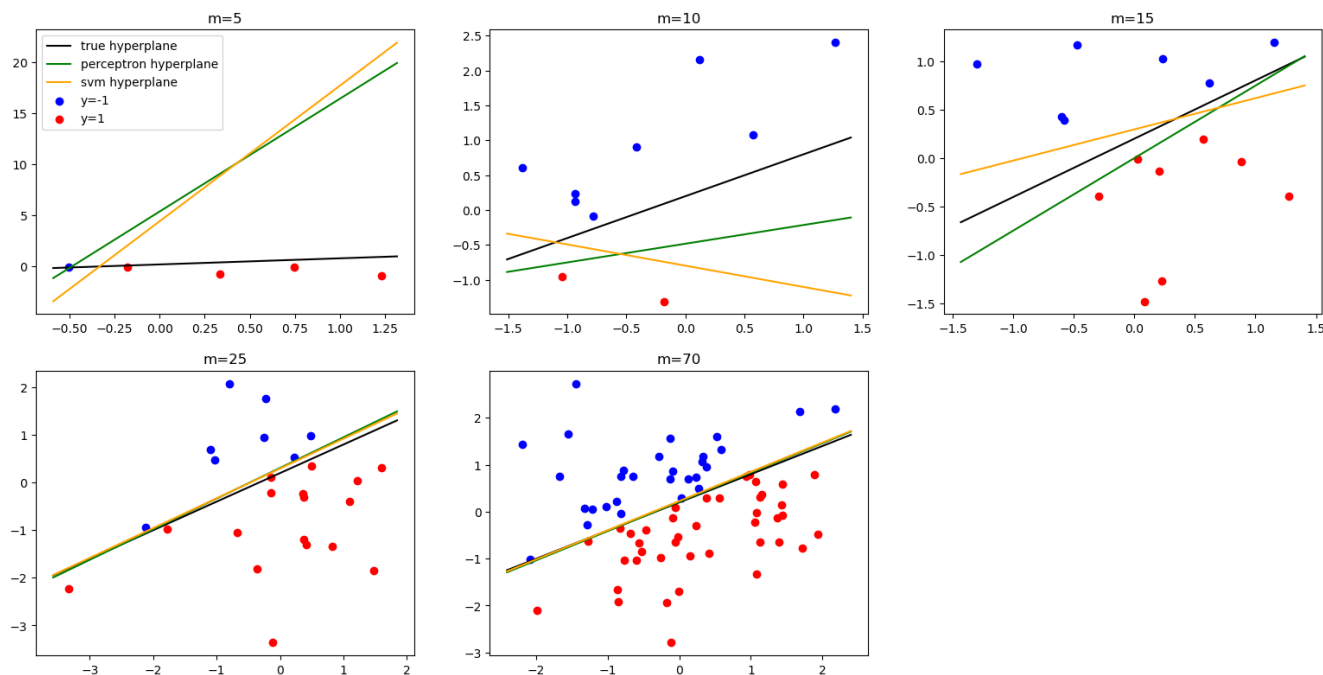
הגדרנו ξ_i מינימליים שעומדים באילוצים, וזה מתאים לבעיה המקורית בה אנו רוצים למזער אותם. ממה שהראנו קודם מתקיים

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

ולכן הבעיות שקולות.

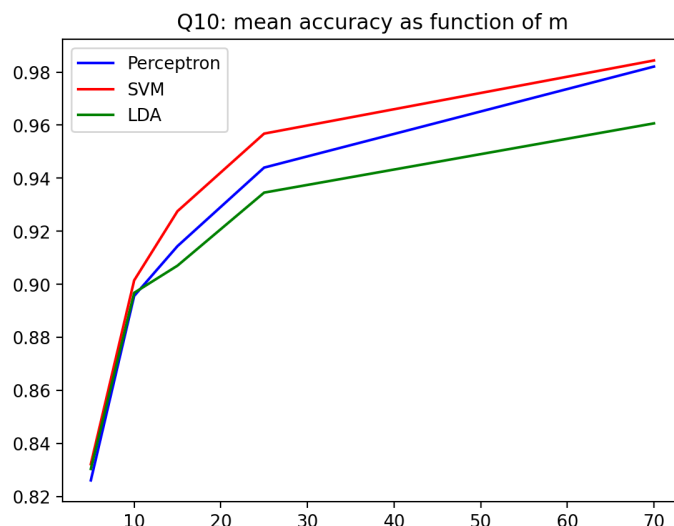
שאלה 9:

Q9: True vs. Perceptron vs. SVM hyperplanes



שאלה 10:

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

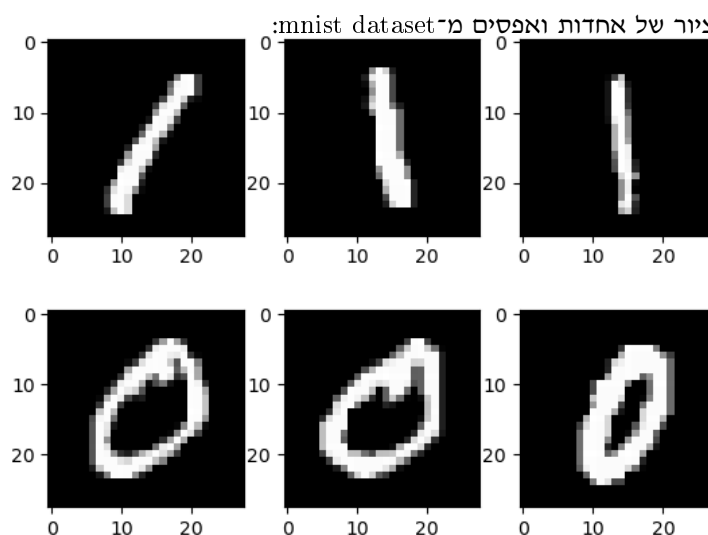


ניתן לראות שה-SVM הוא בעל ה-mean accuracy הכי גבוה, וה-LDA הכי נמוך

שאלה 11:

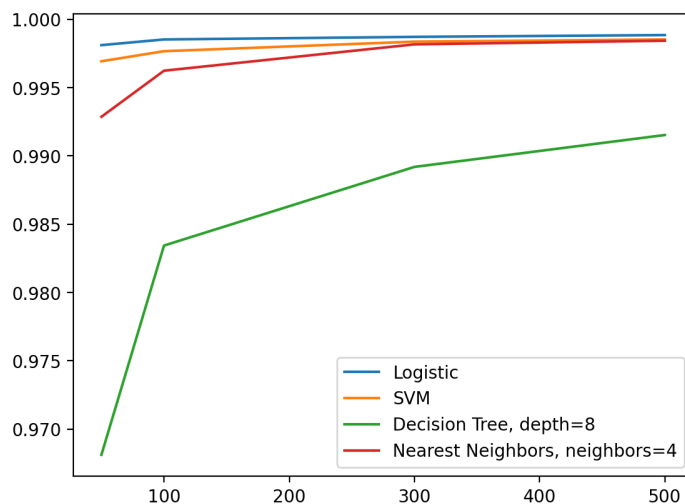
כיוון ש-LDA מניח הנחה שגויה לגבי התפלגות הדגימות (הוא מניח ש- X, y באים מהתפלגות משותפת וכן ש- X מגיע משני גאוסיאניים, בעוד שדגמנו רק את X באופן בלתי תלוי וגאוסייני, ואינו תלוי ב- y) רואים שהוא מסווג הכי פחות טוב (accuracy נמוך יותר). SVM ו-Perceptron לא מניחים הנחה שגויה כזו ולכן מסווגים טוב יותר (accuracy 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:

Q14: mean accuracy as function of m



רואים כי שככל שמספר הדגימות עולה, Logistic, SVM, Nearest Neighbors מסווגים פחות או יותר עם אותו accuracy גבוה, לעומתם Decision tree קצת פחות טוב.

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

ניתן לראות שככל שמספר הדגימות עולה כך גם זמן הריצה. בנוסף ניתן לראות ש-Logistic, Decision Tree לוקחים בערך אותו זמן, SVM לוקח קצת יותר ו-Nearest Neighbors לוקח הרבה יותר זמן. לגבי Nearest Neighbors: ראינו בכיתה שמסווג זה פחות יעיל כי עבור כל דגימה שהוא מסווג (ב-predict) צריך לעבור על כל ה-training set ולמצוא את K השכנים הכי קרובים. לחלופין, גם אם ב-fit הוא מכניס את ה-training set למבנה נתונים שנותן לו יעילות ב-predict (נאמר בהרצאה שניתן לעשות זאת) זה תהליך שלוקח זמן ולכן ההבדל בזמנים.