

IML ex3 3/3/10 4/5/11

① If we know D , our best predictor would have been assigning the class with higher prob:

$$\forall x \in X \quad h_D(x) = \begin{cases} +1 & \text{if } \Pr(Y=+1|x) \geq \frac{1}{2} \\ -1 & \text{else} \end{cases}$$

where the probability is over D . This classifier knows as the Bayes Optimal classifier show that:

$$h_D = \arg \max_{y \in \{-1, 1\}} \Pr(x|y) \Pr(y)$$

$$\arg \max_{y \in \{-1, 1\}} \Pr(x|y) \cdot \Pr(y) = \arg \max_{y \in \{-1, 1\}} \Pr(y|x) \cdot \Pr(x)$$

Since $\Pr(x)$ is the same for all y , we can ignore it and just maximize $\Pr(y|x)$.

$$= \arg \max_{y \in \{-1, 1\}} \Pr(y|x) = \arg \max_{y \in \{-1, 1\}} (\Pr(y=+1|x), \Pr(y=-1|x))$$

Let $x \in X$ be a point. We want to show that $h_D(x) = \arg \max_{y \in \{-1, 1\}} \Pr(y|x)$.
 If $\Pr(y=+1|x) \geq \frac{1}{2}$, then $h_D(x) = +1$.
 If $\Pr(y=+1|x) < \frac{1}{2}$, then $h_D(x) = -1$.
 Since $\Pr(y=+1|x) + \Pr(y=-1|x) = 1$, we have $\Pr(y=-1|x) = 1 - \Pr(y=+1|x)$.
 If $\Pr(y=+1|x) < \frac{1}{2}$, then $\Pr(y=-1|x) > \frac{1}{2}$.

Therefore, $h_D(x) = \arg \max_{y \in \{-1, 1\}} \Pr(y|x)$.
 If $\Pr(y=+1|x) \geq \frac{1}{2}$, then $h_D(x) = +1$.
 If $\Pr(y=+1|x) < \frac{1}{2}$, then $h_D(x) = -1$.
 This is exactly the definition of h_D .

(2) mean vector $\mu_y \in \mathbb{R}^d$, $\forall y \sim N(\mu_y, \Sigma)$, $X = \mathbb{R}^d$ and
 covar matrix $\Sigma \in \mathbb{R}^{d \times d}$

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

הפונקציה הזו היא פונקציית הצפיפות של המשתנה x בהינתן y

Bayes optimal classifier \hat{y} הוא המינימום של הסיכון

$$h_0(y) = \arg \max_{y \in \mathcal{Y}} \delta_y(y)$$

הסיכון $\mathbb{R}^d \rightarrow \mathbb{R}$

$$\delta_y(y) = x^T \Sigma^{-1} \mu_y - \frac{1}{2} \mu_y^T \Sigma^{-1} \mu_y - \ln P(y)$$

הפונקציה $f(x|y) = P(x|y)$ היא פונקציית הצפיפות של x בהינתן y .
 הפונקציה $P(y)$ היא פונקציית הצפיפות של y .
 הפונקציה $\delta_y(y)$ היא פונקציית הסיכון של y .

הפונקציה $\delta_y(y)$ היא פונקציית הסיכון של y .
 הפונקציה $P(y)$ היא פונקציית הצפיפות של y .

$$h_0(y) = \arg \max_{y \in \mathcal{Y}} P(x|y) P(y) = \arg \max_{y \in \mathcal{Y}} f(x|y) P(y) =$$

$$= \arg \max_{y \in \mathcal{Y}} \ln(f(x|y) \cdot P(y)) = \arg \max_{y \in \mathcal{Y}} \ln(f(x|y)) + \ln(P(y)) =$$

$$= \arg \max_{y \in \mathcal{Y}} \ln\left[\frac{1}{(2\pi)^d \det(\Sigma)} \exp\left\{-\frac{1}{2}(x-\mu_y)^T \Sigma^{-1}(x-\mu_y)\right\}\right] + \ln(P(y)) =$$

$$\arg \max_{y \in \mathcal{Y}} \left[\ln\left(\frac{1}{(2\pi)^d \det(\Sigma)}\right) + \ln\left(\exp\left\{-\frac{1}{2}(x-\mu_y)^T \Sigma^{-1}(x-\mu_y)\right\}\right) \right] + \ln(P(y)) =$$

הפונקציה $\ln\left(\frac{1}{(2\pi)^d \det(\Sigma)}\right)$ היא קבוע, ולכן היא לא משפיעה על המינימום.

$$\arg \max_{y \in \mathcal{Y}} \left[-\frac{1}{2}(x-\mu_y)^T \Sigma^{-1}(x-\mu_y) \right] + \ln(P(y)) =$$

$$= \arg \max_{y \in \mathcal{Y}} \left[-\frac{1}{2} x^T \Sigma^{-1} x + \frac{1}{2} x^T \Sigma^{-1} \mu_y - \frac{1}{2} \mu_y^T \Sigma^{-1} \mu_y \right] + \ln(P(y)) =$$

$$\arg \max_{y \in \mathcal{Y}} \left(x^T \Sigma^{-1} \mu_y - \frac{1}{2} \mu_y^T \Sigma^{-1} \mu_y \right) + \ln(P(y))$$

$$y \in \mathbb{R}^d, y_i \in \mathbb{R}^{\pm 13}$$

זכרתי את שאלה 988 קודם, האנחנו לא יודעים את $\mu_{n+1}, \Sigma, P(y)$ (3)
 צריך להגדיר $S = (y_1, y_2, \dots, y_m, y)$ ϵ training set
 האם קיימת האומדן של $\mu_{n+1}, \Sigma, P(y)$ זכרנו $S \subset$
 (העבר הישן)

$$Pr(y) = \frac{1}{m} \sum_{i \in \{1\}} \mathbb{1}\{y_i = y\} = \begin{cases} Pr(y=1) = \frac{1}{m} \sum_{i \in \{1\}} \mathbb{1}\{y_i = 1\} \\ Pr(y=-1) = \frac{1}{m} \sum_{i \in \{-1\}} \mathbb{1}\{y_i = -1\} \end{cases}$$

$$q_i^n = \frac{1}{\#y_i} \sum_{j=1}^n \{ (x_j) | \{y_i\} \}$$

$$\hat{\mu}_1 = \frac{1}{\#Y=1} \sum_{i=1}^n \mathbb{I}(Y_i=1) X_i$$

(32 no. L. 14211) over matrix (106)

$$\Sigma = \frac{1}{m} \sum_{y=\pm 1} \sum_{i=1}^m (x_i - \mu_y)(x_i - \mu_y)^T$$

(4) Sprim : Not a BIP but IR classifier

[illegible]

SVM - Formulation

⑤ Quadratic Program (QP) $\arg \min_{v \in \mathbb{R}^n} \left(\frac{1}{2} v^T Q v - c^T v \right)$
 subject to $v \in \mathbb{R}^n$

$Q \in \mathbb{R}^{n \times n}$, $A \in \mathbb{R}^{m \times n}$, $a \in \mathbb{R}^m$, $b \in \mathbb{R}^m$

QP is a hard-SVM problem. We want to find the maximum margin hyperplane.

$\arg \min_{w,b} \|w\|^2$ s.t. $\forall i, y_i (\langle w, x_i \rangle + b) \geq 1$

w is the weight vector, v, a are given.

Assume $v = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$, $A = \begin{bmatrix} x_1^T \\ \vdots \\ x_m^T \end{bmatrix}$. Then $\arg \min_{w,b} \|w\|^2 \rightarrow \arg \min_{w,b} \left(\frac{1}{2} v^T Q v - c^T v \right)$ s.t. $A w \leq b$

$\arg \min_{w,b} \|w\|^2 = \langle w, w \rangle = w^T w = w^T I w$ (1)

$a=0, Q=I, \arg \min_{w,b} w^T I w = \frac{1}{2} v^T Q v - c^T v$ s.t. $A w \leq b$

$\arg \min_{w,b} \frac{1}{2} \|w\|^2 = \frac{1}{2} w^T I w$ s.t. $\forall i, y_i (\langle w, x_i \rangle + b) \geq 1$

s.t. $\forall i, y_i (\langle w, x_i \rangle + b) \geq 1 \Leftrightarrow \forall i, y_i \langle w, x_i \rangle + y_i b \geq 1$ (2)

$\langle w, x_i \rangle = w^T x_i = [y_1 \ x_1 \ y_2 \ x_2 \ \dots \ y_m \ x_m] \cdot \begin{bmatrix} w \\ b \end{bmatrix} \geq 1$

$\begin{bmatrix} (y_1 x_1) & y_1 \\ \vdots & \vdots \\ (y_m x_m) & y_m \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_m \\ b \end{bmatrix} \geq \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \Leftrightarrow$

$\underbrace{\begin{bmatrix} (y_1 x_1) & y_1 \\ \vdots & \vdots \\ (y_m x_m) & y_m \end{bmatrix}}_A \underbrace{\begin{bmatrix} w_1 \\ \vdots \\ w_m \\ b \end{bmatrix}}_v \leq \underbrace{\begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}}_d$

$\frac{1}{2} w^T I w$ is the squared norm of w . We want to minimize it.

The constraints are linear in w and b . This is a linear programming problem.

We can solve this problem using the Simplex method or other LP solvers.

The optimal solution gives us the maximum margin hyperplane.

$w = (w_1, \dots, w_m, b)$ s.t. $\forall i, y_i \langle w, x_i \rangle + y_i b \geq 1$

⑥ In the soft-SVM we defined the problem:

$$(I) \arg \min_{w, \xi_i \geq 0} \frac{\lambda}{2} \|w\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i \quad \text{s.t. } \forall_i \quad y_i \langle w, x_i \rangle \geq 1 - \xi_i \quad \text{and } \xi_i \geq 0$$

בין קבוצת שברים הנחלקת על ידי קבוצת נקודות.

$$(II) \arg \min_w \frac{\lambda}{2} \|w\|^2 + \frac{1}{m} \sum_{i=1}^m \ell_{\text{hinge}}(y_i \langle w, x_i \rangle) \quad \text{כאשר } \ell_{\text{hinge}}(a) = \max\{0, 1-a\}$$

לכן אנו רוצים להבין כי יש לנו וקטור w שנקראת המישור
שהוא תלוי בפרמטרים w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

$$(I) \arg \min_{w, \xi_i \geq 0} \frac{\lambda}{2} \|w\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i \quad \text{s.t. } \forall_i \quad y_i \langle w, x_i \rangle \geq 1 - \xi_i \quad \text{and } \xi_i \geq 0$$

המשפט הזה הוא הממשלה של w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

אם w הוא וקטור שנקרא w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

אם w הוא וקטור שנקרא w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

אם w הוא וקטור שנקרא w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

אם w הוא וקטור שנקרא w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

אם w הוא וקטור שנקרא w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

אם w הוא וקטור שנקרא w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

אם w הוא וקטור שנקרא w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

אם w הוא וקטור שנקרא w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

אם w הוא וקטור שנקרא w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

אם w הוא וקטור שנקרא w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

אם w הוא וקטור שנקרא w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

אם w הוא וקטור שנקרא w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

אם w הוא וקטור שנקרא w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

אם w הוא וקטור שנקרא w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

אם w הוא וקטור שנקרא w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

אם w הוא וקטור שנקרא w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

אם w הוא וקטור שנקרא w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

אם w הוא וקטור שנקרא w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

אם w הוא וקטור שנקרא w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

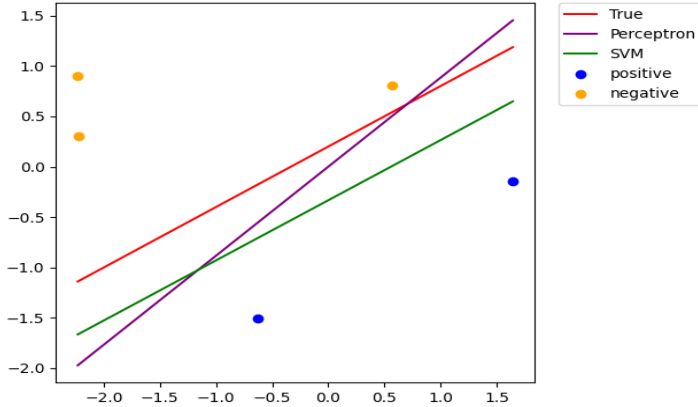
אם w הוא וקטור שנקרא w ו- b ויש לנו $\|w\|^2$
והוא נקרא ℓ_{hinge} ויש לנו $\ell_{\text{hinge}}(a) = \max\{0, 1-a\}$

IML ex3 classifiers

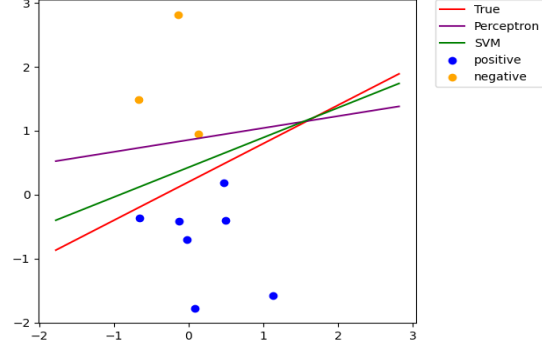
אביב אוהיון 313410458

שאלה 9:

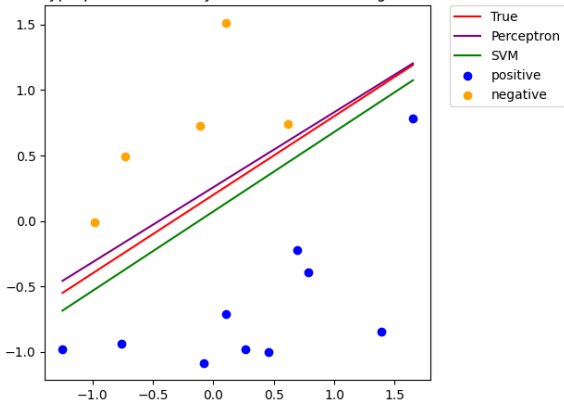
hyperplane created by the each classifier algorithm



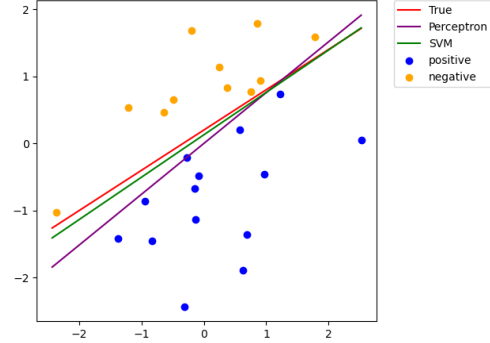
hyperplane created by the each classifier algorithm



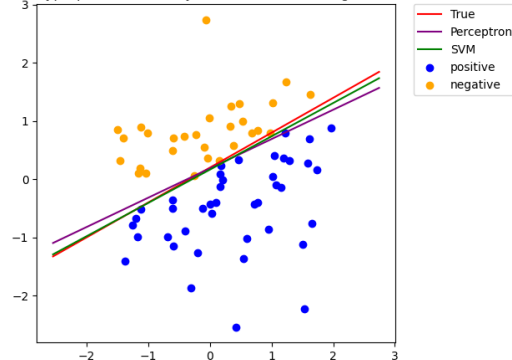
hyperplane created by the each classifier algorithm



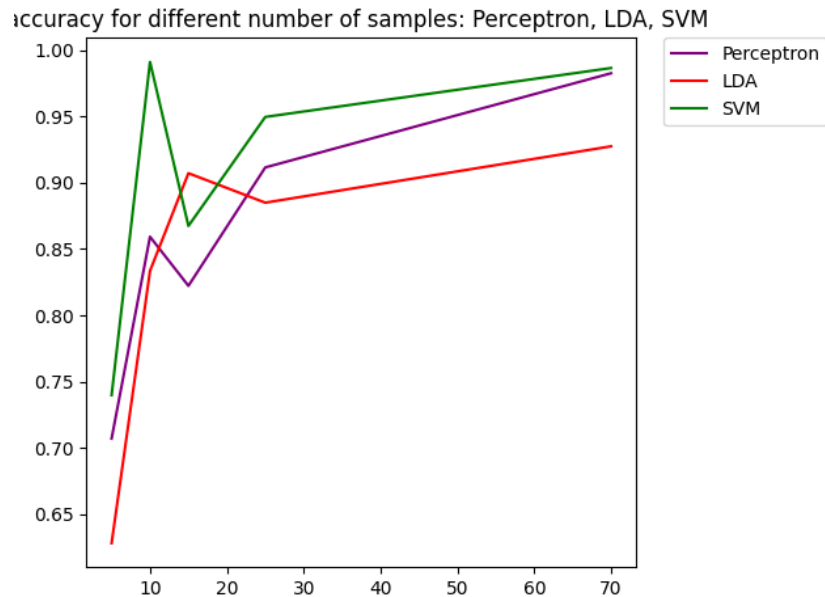
hyperplane created by the each classifier algorithm



hyperplane created by the each classifier algorithm

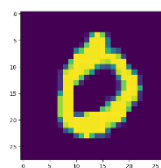
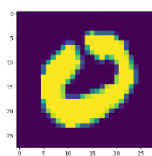
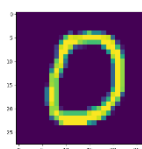
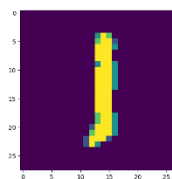
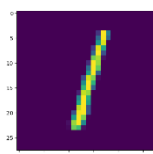
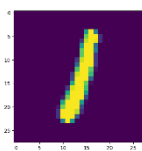


שאלה 11 + 10:



קלספייר ה SVM הוא הטוב ביותר לדעתי. הסיבה לכך היא שהמעקב הגדרתו הוא אינו מניח רליזאביליות על הדטא כלומר הוא "יודע" שאת הדטא לווא דווקא אפשר להפריד לינארית ופשוט מצמצם את טווח הטעות שאנחנו נגדיר לו בניגוד לperceptron שמניח ראלזיביליות (יוצר את המודל באופן איטרטיבי לפי ה training שלו ללא התחשבות בדגימות עתידיות שיגיעו) LDA שמניח שעובד לפי משערכים הסתברותיים (כגון שונות משותפת לדאטא) ותלוי מאוד בסט הדגימות שהוא מקבל לאימון עקב כך בצורה רנדומית מההתפלגות שניתנה לו (במקרה זה נתנו לו את f שהיא דטרמיניסטית ובפרט לא ההתפלגות הנורמאלית שהוא מצפה לה). עקב כך, ככל שכמות הדגימות עולה, ככה גם ההנחות על 2 האלגוריתם הללו הופכות ליותר קשות להבטחה ועל כן הניבוי שהן יניבו עלול להיות שגוי בעוד ש SVM לא "בונה" על ההבטחות האלו ורק מנסה לצמצם את כמות השגיאה לפי הדגימות שניתנו לו לפי עקרון הלמידה של ה max margin. זו הסיבה שאנחנו רואים בנוסף לכך כי LDA מכמות מסוימת של דגימות נישאר יציב ונמוך יותר בעוד ש SVM מאותה כמות דגימות רק עולה ביציבות.

שאלה 12:



שאלה 14:

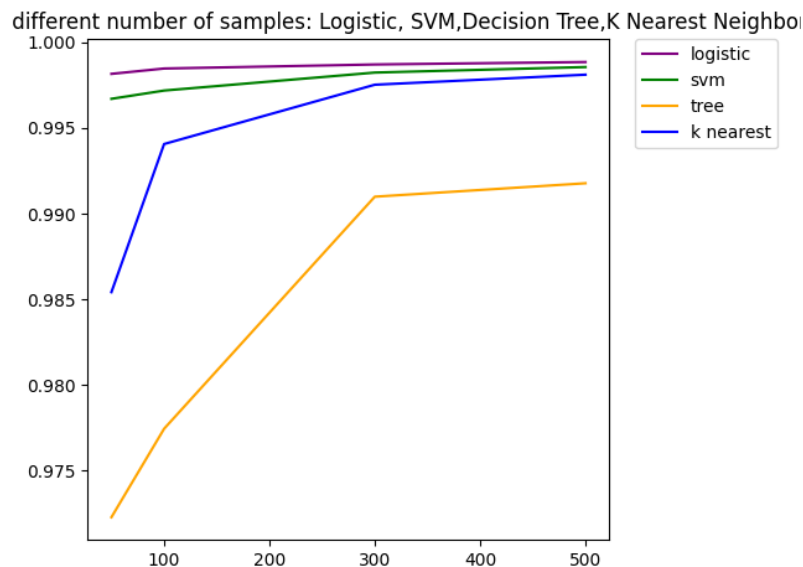
זמני הריצה של כל קלסיפייר (שמאל: קטן, ימין: גדול) הינם:

Logistic mean elapsed time: [0.06364116 0.0705077 0.09678286 0.12411702]

Svm mean elapsed time: [0.07811794 0.10009616 0.15179591 0.18879328]

tree_mean_elapsed_time: [0.01249707 0.01499523 0.02599104 0.04071514]

k nearest mean elapsed time: [0.15015065 0.15144804 0.1791475 0.21961463]



הדברים שהבחנתי בהם בהבדלים בין זמני הריצה של האלגוריתם הם: ראשית, k nearest הינו הקלסיפר שלוקח הכי הרבה זמן, דבר זה אינו מפתיע עקב העובדה שהוא על כל דגימה, עובר על כל סט הדגימות על מנת להכריע מי השכנים הקרובים אליו ביותר. נשין לב שבנוסף עץ ההחלטה עובד על פיצור המרחב וצמצומו עד עומק קבוע ולפי כך הוא יוצר את הקלפסיפקציה, עקב כך אומנם החישוב עצמו יהיה מהיר כי ההכרע במחינה חישובית היא די קלה (סך הכל להגיד אם משהו גדול ממהשו) לעומת השאר אבל כיוון שעומק העץ הוא זה שיכריע את הקלסיפקיה בהינתן עומק לא גדול במיוחד הכרעה של העץ צפויה להיות לא תמיד הכי מדויקת לעומת השאר. ה-SVM משתמש בכך שהוא פותר את בעייה שהיא אופטימוזציה ריבועית שיש לה מקסימום והיא חלקה, עקב ובשכל העובדה ש-SVM הוא איטרטיבי על סט הדגימות (להבדיל מ-k nearest לדוגמא) הוא גם עובד יחסית מהר ושילוב של 2 הסיבות מסביר למה הוא גם נותן את התוצאות הכי טובות.

הערה: בסעיף זה השתמשתי בפוקיציות הפסרייה של sklearn על מנת לא לשנות את הלייבלים ל 0 1 בסעיף זה.