

מערכות לומדות תרגיל 3

עמית בסקין 312259013

1. לכל $x \in \mathcal{X}$:

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

צ.להוכיח:

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

הוכחה:

ראשית נשים לב כי:

$$\Pr(x | y) \Pr(y) = \frac{\Pr(x \wedge y)}{\Pr(y)} \cdot \Pr(y) = \Pr(x \wedge y)$$

כלומר בעצם צריך להוכיח ש-

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

יהא $x \in \mathcal{X}$. נחלק למקרים:

• $\Pr(y = 1 | x) \geq \frac{1}{2}$, קרי:

$$\frac{\Pr(x \wedge y = 1)}{\Pr(x)} \geq \frac{1}{2}$$

$$\Pr(x \wedge y = 1) \geq \frac{1}{2} \Pr(x)$$

ובפרט:

$$\frac{\Pr(x \wedge y = -1)}{\Pr(x)} \leq \frac{1}{2}$$

$$\Pr(x \wedge y = -1) \leq \frac{1}{2} \Pr(x)$$

כלומר:

$$\Pr(x \wedge y = 1) \geq \Pr(x \wedge y = -1)$$

ולכן:

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

ומאחר שכל המעברים הם אסם אזי ש- $\frac{1}{2} \geq \Pr(y = 1 | x)$ אסם $\operatorname{argmax}_{y \in \{\pm 1\}} \Pr(x \wedge y) = 1$.

מש"ל מקרה ראשון.

$$\bullet \Pr(y = 1 | x) \leq \frac{1}{2}, \text{ קרי:}$$

$$\frac{\Pr(x \wedge y = 1)}{\Pr(x)} \leq \frac{1}{2}$$

$$\Pr(x \wedge y = 1) \leq \frac{1}{2} \Pr(x)$$

ובפרט:

$$\frac{\Pr(x \wedge y = -1)}{\Pr(x)} \geq \frac{1}{2}$$

$$\Pr(x \wedge y = -1) \geq \frac{1}{2} \Pr(x)$$

כלומר:

$$\Pr(x \wedge y = 1) \leq \Pr(x \wedge y = -1)$$

ולכן:

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

ומאחר שכל המעברים הם אסם אזי ש- $\frac{1}{2} \leq \Pr(y = 1 | x)$ אסם $\operatorname{argmax}_{y \in \{\pm 1\}} \Pr(x \wedge y) = -1$.

מש"ל מקרה שני.

מש"ל

2. נניח ש- $\mathcal{X} = \mathbb{R}^d$ וש- $x | y \sim \mathcal{N}(\mu_y, \Sigma)$ עבור $\mu_y \in \mathbb{R}^d$ ו- $\Sigma \in \mathbb{R}^d$, קרי פונקציית הצפיפות היא:

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

צ.להראות ש-

$$h_{\mathcal{D}}(x) = \operatorname{argmax} \left\{ -\frac{1}{2}x^t \Sigma^{-1}x - \frac{1}{2}\mu_y^t \Sigma^{-1}\mu_y + \ln(\Pr(y)) \mid y \in \{\pm 1\} \right\}$$

הוכחה:

לפי נוסחת בייס:

$$\Pr(y | x) \cdot \Pr(x) = f(x | y) \cdot \Pr(y)$$

קרי:

$$\Pr(y | x) = \frac{f(x | y) \cdot \Pr(y)}{\Pr(x)}$$

נסמן $c = \Pr(x)$ כקבוע כי מספר זה לא מושפע מ- y . אז:

$$\Pr(y | x) = c \cdot f(x | y) \cdot \Pr(y)$$

ונציב את $f(x | y)$:

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

ונסמן $c' = \frac{c}{\sqrt{(2\pi)^d \det(\Sigma)}}$ כקבוע כי לא תלוי ב- y :

$$\Pr(y | x) = c' \Pr(y) \exp\left(-\frac{1}{2}(x - \mu_y)^t \Sigma^{-1}(x - \mu_y)\right)$$

$$\ln(\Pr(y | x)) = \ln(c') + \ln(\Pr(y)) - \frac{1}{2}(x - \mu_y)^t \Sigma^{-1}(x - \mu_y) =$$

$$= \ln(c') + \ln(\Pr(y)) - \frac{1}{2}x^t \Sigma^{-1}x - \frac{1}{2}\mu_y^t \Sigma^{-1}\mu_y + x^t \Sigma^{-1}\mu_y$$

ונסמן:

$$c'' = \ln(c') - \frac{1}{2}x^t \Sigma^{-1}x$$

כי לא תלויים ב- y :

$$\Pr(y | x) = x^t \Sigma^{-1} \mu_y - \frac{1}{2} \mu_y^t \Sigma^{-1} \mu_y + \ln(\Pr(y)) + c''$$

עתה, בסעיף הקודם ראינו ש-

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

אבל למקסם את $\Pr(x | y) \Pr(y)$ זה כמו למקסם את $\ln(\Pr(x | y) \Pr(y))$ קרי:

$$\ln(\Pr(x | y)) + \ln(\Pr(y))$$

קרי את:

$$x^t \Sigma^{-1} \mu_y - \frac{1}{2} \mu_y^t \Sigma^{-1} \mu_y + 2 \ln(\Pr(y)) + c''$$

ומאחר ש- c'' לא תלוי ב- y אזי שזה כמו למקסם את:

$$x^t \Sigma^{-1} \mu_y - \frac{1}{2} \mu_y^t \Sigma^{-1} \mu_y + 2 \ln(\Pr(y))$$

מש"ל

3. נתונה קבוצת דגימות $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$. צלהעריך את $\mu_{+1}, \mu_{-1}, \Sigma, \Pr(y)$ בהתבסס על S .

פיתרון:

נניח בה"כ שקיים $l \leq m$ כך שלכל $i \in [l]$ מתקיים $y_i = 1$.

נסמן:

$$S_{+1} = \{x_i \mid i \in [l]\}$$

$$S_{-1} = \{x_i \mid i \in [m] \setminus [l]\}$$

נסמן ב- X_{+1} ו- X_{-1} את המטריצות שהשורות שלהן הן הוקטורים ב- S_{+1} ו- S_{-1} בהתאמה.

נסמן ב- μ_{+1} ו- μ_{-1} את וקטורי הממוצעים של הוקטורים ב- S_{+1} ו- S_{-1} בהתאמה.

יהא μ וקטור הממוצעים של X , ונסמן ב- \hat{X} את X ממורכזת לפי μ .

נסמן ב- X_{+1}^0 ו- X_{-1}^0 את מטריצות הוקטורים X_{+1} ו- X_{-1} בהתאמה, ממורכזים לפי μ_{+1} ו- μ_{-1} בהתאמה.

נסמן ב- Σ_{+1} ו- Σ_{-1} את מטריצות השונות המשותפת של S_{+1} ו- S_{-1} בהתאמה, קרי:

$$\Sigma_{+1} = \frac{1}{n_{+1}} (X_{+1}^0)^t X_{+1}^0$$

$$\Sigma_{-1} = \frac{1}{n_{-1}} (X_{-1}^0)^t X_{-1}^0$$

ואז מטריצת השונות המשותפת הכוללת מחושבת על ידי:

$$\Sigma_{ij} = \frac{1}{m} \left(n_{+1} (\Sigma_{+1})_{ij} + n_{-1} (\Sigma_{-1})_{ij} \right)$$

נסמן ב- n_{+1} ו- n_{-1} את $|S_{+1}|$ ו- $|S_{-1}|$ בהתאמה.

נסמן $p_{+1} = \Pr(+1)$ ו- $p_{-1} = \Pr(-1)$ ונעריך:

$$p_{+1} = \frac{n_{+1}}{m}$$

$$p_{-1} = \frac{n_{-1}}{m}$$

מש"ל

4. ישנן שתי טעויות אפשריות בסיווג ספאם:

- לסווג הודעה כספאם למרות שהיא לא.

- לסווג כלא ספאם הודעה שהיא כן ספאם.

הטעות שאנחנו הכי לא רוצים לעשות זה לסווג הודעה כספאם למרות שהיא לא. זו טעות חמורה יותר מאשר לסווג כלא ספאם הודעה שהיא כן ספאם. אכן, אם נסווג הודעה כלא ספאם למרות שהיא ספאם, אז המשתמש יאלץ למרבה הצער להיתקל בספאם למרות שהיה מעדיף שלא. לעומת זאת, אם ישנה הודעה שהיא לא ספאם אבל סווגה כספאם, כנראה שהמשתמש יפספס את ההודעה, וזאת על אף שפוטנציאלית מדובר בהודעה חשובה שהמשתמש צריך לקרוא, והנזק שעלול להיגרם למשתמש כתוצאה מפספוס ההודעה, עשוי להיות גדול בהרבה מלקרוא הודעת ספאם שהיה מעדיף שלא לקרוא.

אם כן, אנחנו רוצים שהשגיאה החמורה תהיה שגיאה מסוג ראשון, קרי FP, ולכן ניתן לסיווג הודעה כספאם את הלייבל positive, כלומר אם נסווג הודעה שהיא לא ספאם ב־ positive, אבל טעינו False, אז קיבלנו שגיאה מסוג ראשון, כפי שרצינו.

5. הצורה הקאנונית של תכנית ריבועית היא:

$$\operatorname{argmin}_{v \in \mathbb{R}^n} \left(\frac{1}{2} v^t Q v + a^t v \right)$$

כך ש־

$$A v \leq d$$

עבור:

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

נתבונן בבעיית ה־Hard-SVM:

$$\operatorname{argmin}_{(w,b)} \|w\|^2$$

כך שלכל i :

$$y_i (\langle w, x_i \rangle + b) \geq 1$$

צ.לכתוב אותה כתכנית ריבועית קאנונית.

פיתרון:

נניח שיש לנו m דגימות (x_i, y_i) כאשר כל x_i הוא מסדר $n - 1$.

נסמן ב־ A' מטריצה $m \times n - 1$ שהשורות שלה הם הוקטורים $y_i x_i$.

נסמן ב־ A'' את המטריצה שמתקבלת מ־ A' על ידי הוספת עמודה שהכניסות בה הן y_i , וניקח $A = -A''$.

נסמן $Q' = I_{(n-1) \times (n-1)}$, ונתאר את Q כמטריצת בלוקים שנתונים באלכסון הראשי: הבלוק הראשון הוא Q' והבלוק השני הוא הסקלר 0.

את a ניקח להיות 0_n ואת d ניקח להיות $(-1)_m$

אז התכנית הריבועית שלנו היא:

$$\operatorname{argmin}_{v \in \mathbb{R}^n} \left(\frac{1}{2} v^t Q v + a^t v \right)$$

קרי:

$$Q = 2 \begin{pmatrix} 1 & & & \\ & \ddots & & \\ & & 1 & \\ & & & 0 \end{pmatrix}, \mathbf{v} = \mathbf{w} = \begin{pmatrix} w_1 \\ \vdots \\ w_n \\ b \end{pmatrix}, \mathbf{a} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}, A = - \begin{pmatrix} y_1 & y_1 x_{11} & \cdots & y_1 x_{1(n-1)} \\ \vdots & \vdots & \ddots & \vdots \\ y_m & y_m x_{m1} & \cdots & y_m x_{m(n-1)} \end{pmatrix}, \mathbf{d} = - \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

כלומר:

$$\begin{aligned} & \frac{1}{2} \mathbf{v}^t Q \mathbf{v} + \mathbf{a}^t \mathbf{v} = \\ & = \frac{1}{2} \begin{pmatrix} w_1 & \cdots & w_n \end{pmatrix} 2 \begin{pmatrix} 1 & & & \\ & \ddots & & \\ & & 1 & \\ & & & 0 \end{pmatrix} \begin{pmatrix} b \\ w_1 \\ \vdots \\ w_n \end{pmatrix} + \begin{pmatrix} 0 & \cdots & 0 \end{pmatrix} \begin{pmatrix} w_1 \\ \vdots \\ w_n \\ b \end{pmatrix} = \\ & = \begin{pmatrix} w_1 & \cdots & w_n \end{pmatrix} \begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix} = \|\mathbf{w}\|^2 \end{aligned}$$

כאשר:

$$A\mathbf{v} \leq \mathbf{d}$$

קרי:

$$\begin{aligned} & - \begin{pmatrix} y_1 & y_1 x_{11} & \cdots & y_1 x_{1(n-1)} \\ \vdots & \vdots & \ddots & \vdots \\ y_m & y_m x_{m1} & \cdots & y_m x_{m(n-1)} \end{pmatrix} \begin{pmatrix} b \\ w_1 \\ \vdots \\ w_n \end{pmatrix} \leq - \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \\ & \begin{pmatrix} \langle y_1 \mathbf{x}_1, \mathbf{w} \rangle + y_1 b \\ \vdots \\ \langle y_m \mathbf{x}_m, \mathbf{w} \rangle + y_m b \end{pmatrix} \geq \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \\ & \begin{pmatrix} y_1 \langle \mathbf{w}, \mathbf{x}_1 \rangle + y_1 b \\ \vdots \\ y_m \langle \mathbf{w}, \mathbf{x}_m \rangle + y_m b \end{pmatrix} \geq \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \\ & \begin{pmatrix} y_1 (\langle \mathbf{w}, \mathbf{x}_1 \rangle + b) \\ \vdots \\ y_m (\langle \mathbf{w}, \mathbf{x}_m \rangle + b) \end{pmatrix} \geq \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \end{aligned}$$

וקיבלנו שלכל i צריך להתקיים:

$$y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1$$

כפי שרצינו.

מש"ל.

6. נתבונן בבעיית ה-Soft-SVM:

$$\operatorname{argmin}_{\mathbf{w}, \{\xi_i\}} \left(\frac{1}{2} \lambda \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i \right)$$

כך שלכל i :

$$y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1 - \xi_i$$

$$\xi_i \geq 0$$

צ. להוכיח שהיא שקולה לבעייה:

$$\operatorname{argmin}_{\mathbf{w}} \left(\frac{1}{2} \lambda \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m l^{\text{hinge}}(y_i \langle \mathbf{w}, \mathbf{x}_i \rangle) \right)$$

כאשר:

$$l^{\text{hinge}}(y_i \langle \mathbf{w}, \mathbf{x}_i \rangle) = \max \{0, 1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle\}$$

הוכחה:

מתקיים:

$$\sum_{i=1}^m l^{\text{hinge}}(y_i \langle \mathbf{w}, \mathbf{x}_i \rangle) = \sum_{i, y_i \langle \mathbf{w}, \mathbf{x}_i \rangle < 1} (1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle)$$

כלומר הבעיה הנ"ל שקולה ל-

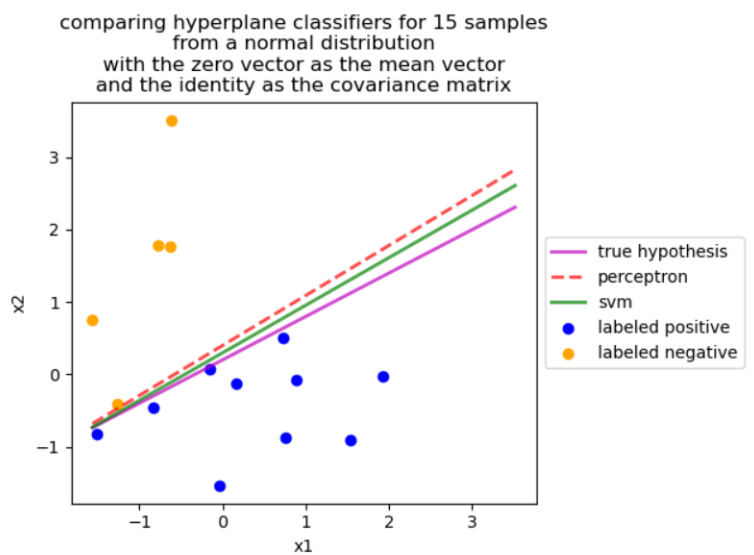
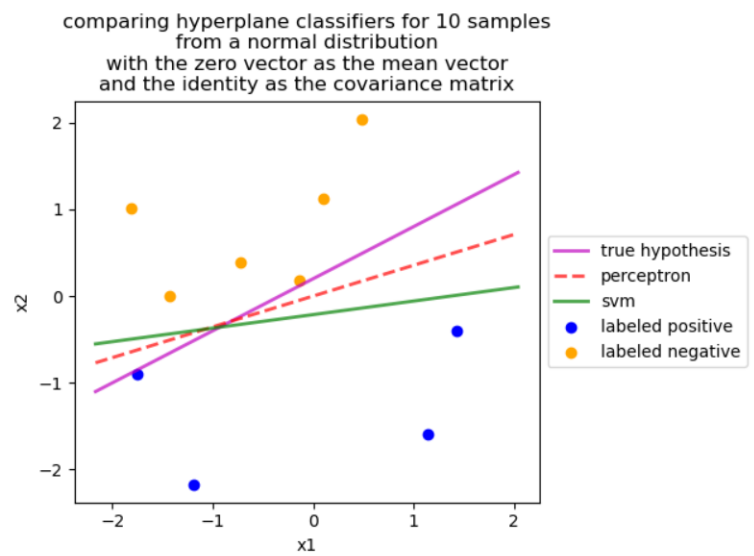
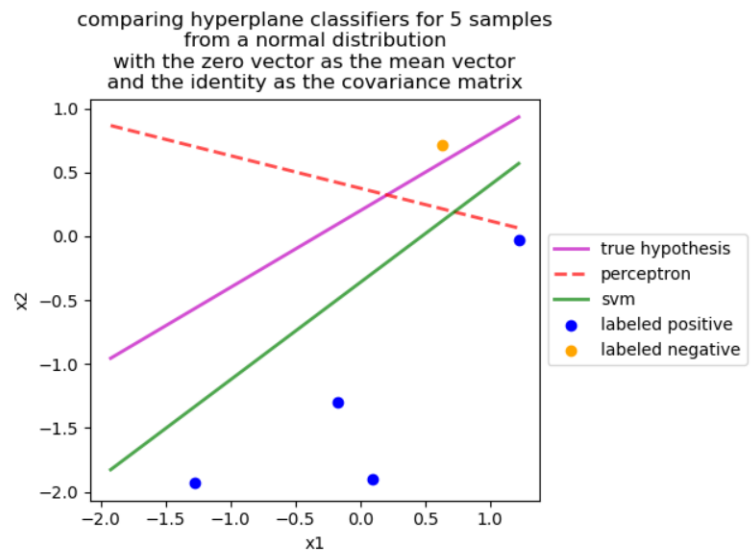
$$\operatorname{argmin}_{\mathbf{w}} \left(\frac{1}{2} \lambda \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i, y_i \langle \mathbf{w}, \mathbf{x}_i \rangle < 1} (1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle) \right)$$

עתה:

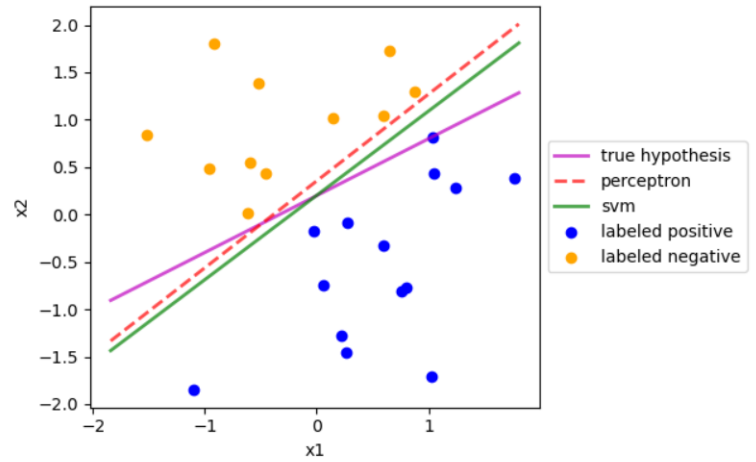
$$y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1 - \xi_i \implies \xi_i \geq 1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle$$

ולכן כשאנחנו ממזערים $\sum_{i=1}^m \xi_i$ אנחנו למעשה גם ממזערים את $1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle$ ולכן המינימומים של שתי הבעיות מתקבלים באותן נקודות.

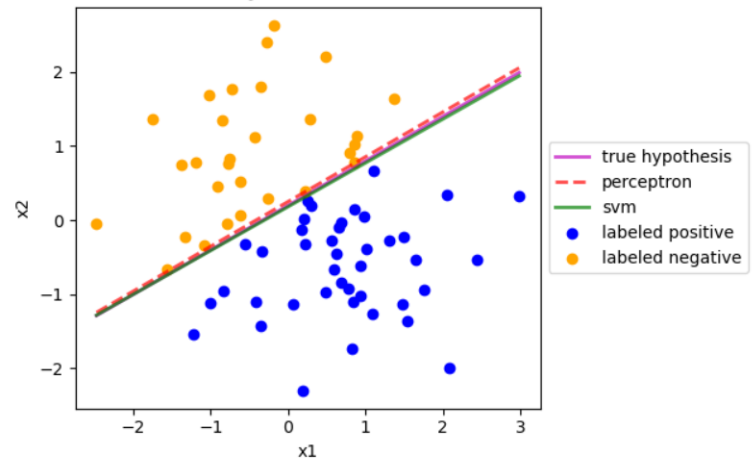
מש"ל



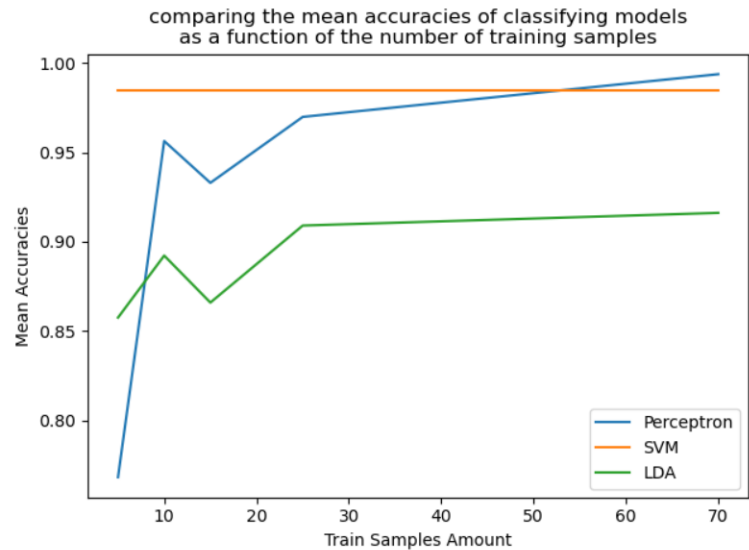
comparing hyperplane classifiers for 25 samples
from a normal distribution
with the zero vector as the mean vector
and the identity as the covariance matrix



comparing hyperplane classifiers for 70 samples
from a normal distribution
with the zero vector as the mean vector
and the identity as the covariance matrix



10. הגרף להלן:

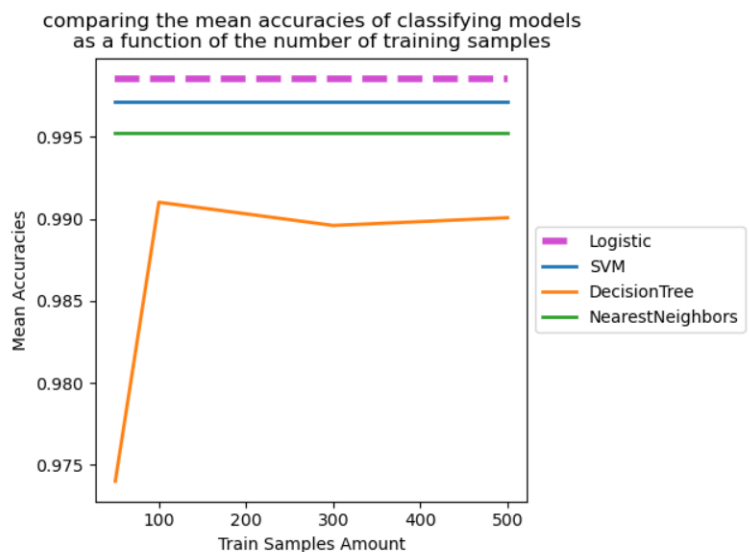


11. בממוצע, ל- SVM יש את הביצוע הטוב ביותר, אף על פי שבאזור ה- 55 דגימות, הפרספטרון עוקף אותו.

הסיבה לכך היא ש- SVM מוצא את הקו שממקסם את את רוחב השוליים, מה שנותן לו מעין "טווח ביטחון" עבור הדגימות הבאות. אכן, ככל שהשוליים צרים יותר כך גדלה ההסתברות לשגיאה, ולכן שוליים מקסימליים ממזערים את ההסתברות לכך (קרי שטח גדול יותר שמאפשר תנודות של הדגימות).

14. הגרף להלן:

עבור המימוש של Logistic, SVM, DecisionTree לקחתי את אותו המימוש מהקובץ models.py כאשר עבור עומק מקסימלי של DecisionTree לקחתי עומק של 5. כמו כן, עבור NearestNeighbors לקחתי את המימוש של NearestCentroid שעבור כל דגימה לוקח את הלייבל של קבוצת הדגימות שהצנטרואיד הממוצע שלהם הוא הקרוב ביותר לדגימה שמסתכלים עליה.



זמני הריצה להלן:

```
Logistic, 50 training samples: 0.17 seconds
Logistic, 100 training samples: 0.19 seconds
Logistic, 300 training samples: 0.19 seconds
Logistic, 500 training samples: 0.17 seconds
SVM, 50 training samples: 3.58 seconds
SVM, 100 training samples: 3.56 seconds
SVM, 300 training samples: 3.53 seconds
SVM, 500 training samples: 3.53 seconds
DecisionTree, 50 training samples: 0.37 seconds
DecisionTree, 100 training samples: 0.39 seconds
DecisionTree, 300 training samples: 0.37 seconds
DecisionTree, 500 training samples: 0.36 seconds
NearestNeighbors, 50 training samples: 0.47 seconds
NearestNeighbors, 100 training samples: 0.48 seconds
NearestNeighbors, 300 training samples: 0.47 seconds
NearestNeighbors, 500 training samples: 0.47 seconds
```

ניתן לראות שבאחוזים אלגוריתם ה- SVM לוקח משמעותית יותר זמן משאר האלגוריתמים.

מכאן ניתן לשער שייתכן כי הדבר טמון בכך שמקסום השוליים יותר קשה לחישוב משאר החישובים שהאלגוריתמים האחרים עושים.