

### מבוא למערכות לומדות – תרגיל 3

1. If we knew  $\mathcal{D}$ , our best predictor would have been assigning the class with the higher probability:

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

where the probability is over  $\mathcal{D}$ . This classifier is known as the **Bayes Optimal** classifier.

Show that

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

1. ראשית נזכר בנוסחת בייס ונוסחת ההסתברות השלמה:

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

$$P(y | x) = P(y = 1 | x) + P(y = -1 | x)$$

נשתמש בנוסחאות ונפתח את ההגדרה עבור  $h_{\mathcal{D}}$ :

$$\begin{aligned} \operatorname{argmax}_{y \in \{\pm 1\}} P(x | y) \cdot P(y) &\stackrel{\text{bayes rule}}{=} \operatorname{argmax} \frac{P(y | x) \cdot P(x)}{P(y)} \cdot P(y) = \operatorname{argmax} P(y | x) \cdot P(x) \stackrel{\text{הסתברות שלמה}}{=} \\ &= \operatorname{argmax} (P(y = 1 | x) + P(y = -1 | x)) \cdot P(x) \end{aligned}$$

כעת נבחין בין שני מקרים אפשריים:

$$\text{אם } P(y = 1 | x) \geq \frac{1}{2} \text{ אזי } P(y = -1 | x) < \frac{1}{2}$$

לפי הגדרת מרחב הסתברות כאשר יש רק שתי אפשרויות ל- $y$ , נקבל חלוקה שלמה שלו ונקבל כי:

$$\operatorname{argmax} (P(y = 1 | x) + P(y = -1 | x)) \cdot P(x) = 1$$

$$\text{במקרה השני } P(y = 1 | x) < \frac{1}{2} \text{ אזי } P(y = -1 | x) \geq \frac{1}{2} \text{ ונקבל:}$$

$$\operatorname{argmax} (P(y = 1 | x) + P(y = -1 | x)) \cdot P(x) = -1$$

קיבלנו שההגדרה שיולה לכל  $h_{\mathcal{D}}$  כנדרש.

2. Assume that  $\mathcal{X} = \mathbb{R}^d$  and that  $\mathbf{x}|y \sim \mathcal{N}(\mu_y, \Sigma)$  for some mean vector  $\mu_y \in \mathbb{R}^d$  and covariance matrix  $\Sigma \in \mathbb{R}^{d \times d}$  (that is, the covariance matrix  $\Sigma$  is the same for both  $y \in \{\pm 1\}$ , but the expectation  $\mu_y$  is different for each  $y \in \{\pm 1\}$ ). In other words,

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

where  $f$  is the density function for the multivariate normal distribution. Show that in this case, if we knew  $\mu_{+1}, \mu_{-1}$  and  $\Sigma$  then the Bayes Optimal classifier is

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

where  $\delta_{+1}$  and  $\delta_{-1}$  are functions  $\mathbb{R}^d \rightarrow \mathbb{R}$  given by

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

.2

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

וכעת, מנוסחאת בייס :

$$P(y) = f_Y(y) = \frac{f_X(x) f_{Y|X=x}(y)}{f_{X|Y=y}(x)}$$

נציב ונקבל :

$$= \operatorname{argmax} x^T \Sigma^{-1} \mu_y - \frac{1}{2} \mu_y^T \Sigma^{-1} \mu_y + \ln \left( \frac{f_X(x) f_{Y|X=x}(y)}{f_{X|Y=y}(x)} \right)$$

נפעיל חוקי  $\ln$  :

$$= \operatorname{argmax} x^T \Sigma^{-1} \mu_y - \frac{1}{2} \mu_y^T \Sigma^{-1} \mu_y + \ln(f_X(x)) + \ln(f_{Y|X=x}(y)) - \ln(f_{X|Y=y}(x))$$

כעת שנשים לב כי הביטוי  $\ln(f_X(x))$  הינו קבוע ונוכל להשמיטו :

$$= \operatorname{argmax} x^T \Sigma^{-1} \mu_y - \frac{1}{2} \mu_y^T \Sigma^{-1} \mu_y + \ln(f_{Y|X=x}(y)) - \ln(f_{X|Y=y}(x))$$

כעת נציב את הביטוי הנתון עבור  $\ln(f_{X|Y=y}(x))$  ונקבל :

$$= \operatorname{argmax} x^T \Sigma^{-1} \mu_y - \frac{1}{2} \mu_y^T \Sigma^{-1} \mu_y + \ln(f_{Y|X=x}(y)) - \ln \left( \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \exp \left\{ -\frac{1}{2} (x - \mu_y)^\top \Sigma^{-1} (x - \mu_y) \right\} \right)$$

נשים לב כי הביטוי  $\ln \left( \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \right)$  קבוע ונוכל להשמיטו :

$$\begin{aligned}
 &= \operatorname{argmax} x^T \Sigma^{-1} \mu_y - \frac{1}{2} \mu_y^T \Sigma^{-1} \mu_y + \ln(f_{Y|X=x}(y)) + \frac{1}{2} (x - \mu_y)^T \Sigma^{-1} (x - \mu_y) = \\
 &= \operatorname{argmax} \ln(f_{Y|X=x}(y)) + x^T \Sigma^{-1} \mu_y - \frac{1}{2} \mu_y^T \Sigma^{-1} \mu_y + \left( \frac{1}{2} x^T \Sigma^{-1} - \frac{1}{2} \mu_y^T \Sigma^{-1} \right) (x - \mu_y) = \\
 &= \operatorname{argmax} \ln(f_{Y|X=x}(y)) + x^T \Sigma^{-1} \mu_y - \frac{1}{2} \mu_y^T \Sigma^{-1} \mu_y + \frac{1}{2} x^T \Sigma^{-1} x - \frac{1}{2} \mu_y^T \Sigma^{-1} x - \frac{1}{2} x^T \Sigma^{-1} \mu_y + \frac{1}{2} \mu_y^T \Sigma^{-1} \mu_y = \\
 &= \operatorname{argmax} \ln(f_{Y|X=x}(y)) + \frac{1}{2} x^T \Sigma^{-1} \mu_y + \frac{1}{2} x^T \Sigma^{-1} x - \frac{1}{2} \mu_y^T \Sigma^{-1} x \\
 &\quad \text{ומכיוון שהביטוי } \frac{1}{2} x^T \Sigma^{-1} x \text{ הוא קבוע:}
 \end{aligned}$$

$$\begin{aligned}
 &= \operatorname{argmax} \ln(f_{Y|X=x}(y)) + \frac{1}{2} x^T \Sigma^{-1} \mu_y - \frac{1}{2} \mu_y^T \Sigma^{-1} x \\
 &\quad \text{נשים לב כי הביטוי } \Sigma^{-1} \text{ הינו מטריצה סימטרית ונקבל:}
 \end{aligned}$$

$$\begin{aligned}
 &= \operatorname{argmax} \ln(f_{Y|X=x}(y)) + \frac{1}{2} x^T \Sigma^{-1} \mu_y - \frac{1}{2} x^T \Sigma^{-1} \mu_y = \\
 &= \operatorname{argmax} \ln(f_{Y|X=x}(y))
 \end{aligned}$$

וממונוטוניות פונקציית ה- $\ln$  נקבל:

$$= \operatorname{argmax} P(y | x) \stackrel{\text{bayes rule}}{\cong} P(x | y) \cdot P(y) \stackrel{\text{לפי שאלה 1}}{\cong} h_D$$

כנדרש.

3. In practice, we don't know  $\mu_{+1}, \mu_{-1}, \Sigma$  and  $\Pr(y)$ . In order to turn the above into a classifier, given a training set  $S = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$ , we need to estimate them. Write your formula for estimating  $\mu_{+1}, \mu_{-1}, \Sigma$  and  $\Pr(y)$  based on  $S$ .

3. בקורס הסתברות למדנו כי אומד בלתי מוטה לתוחלת צריך להיות הממוצע, ובהתאם נגדיר את האומדים הנ"ל:

$$\hat{\mu}_{+1} = \frac{1}{|S_{+1}|} \cdot \sum_{i \in S_{+1}} x_i, \quad \hat{\mu}_{-1} = \frac{1}{|S_{-1}|} \cdot \sum_{i \in S_{-1}} x_i$$

ולכן האומד הבלתי תלוי לשונות יהיה:

$$\hat{\Sigma} = \frac{1}{m-1} \sum_{i=1}^m (x_i - \bar{x})(x_i - \bar{x})^T$$

לבסוף עבור הסתברות  $P(y)$  נוכל לבחור בהסתברות ברנולי עם פרמטר  $p$  כך ש-

$$p = \frac{|S_{+1}|}{m}$$

## Spam

[Relevant material - Lecture 3]

4. You are building a spam filter - a classifier that receives an email and decides whether it's a spam message or not. What are the two kinds of errors that your classifier could make? Which of them is the error we really don't want to make? Which of the labels {spam, not-spam} should be the **negative** label and which should be the **positive** label, if we want the false-positive error (Type-I error) to be the error we really don't want to make?

4. כשאנו רוצים לסנן מיילים לפי הפרמטר האם הם ספאם או לא נצטרך להבחין בין שתי עלויות: הראשונה, כשאנו קובעים על מייל מסוים שהוא ספאם כאשר בפועל הוא אינו כזה, לכן אנו "מפספסים" משהו שהיה יכול להיות חשוב לנו. השניים, שאנו מדווחים על מייל מסוים שהוא תקין אבל בפועל הוא ספאם, ובזבזנו עליו את הזמן שלנו. אם כך, הטעות שנרצה להימנע ממנה יותר היא הטעות הראשונה, תיוג מייל חשוב כספאם, לכן נסמן את +1 להיות התגית "ספאם" ו-1- לתגית "אינו ספאם".

## SVM- Formulation

[Relevant material - Recitation 4]

5. The canonical form of a Quadratic Program (QP) is:

$$\begin{aligned} \underset{\mathbf{v} \in \mathbb{R}^n}{\operatorname{argmin}} \quad & \frac{1}{2} \mathbf{v}^\top Q \mathbf{v} + \mathbf{a}^\top \mathbf{v} \\ \text{s.t.} \quad & A \mathbf{v} \leq \mathbf{d}, \end{aligned}$$

where  $Q \in \mathbb{R}^{n \times n}$ ,  $A \in \mathbb{R}^{m \times n}$ ,  $\mathbf{a} \in \mathbb{R}^n$ ,  $\mathbf{d} \in \mathbb{R}^m$  are fixed vectors and matrices.

Write the Hard-SVM problem as a QP problem in canonical form. Specifically, using the Hard-SVM problem formulation

$$\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \|\mathbf{w}\|^2 \text{ s.t. } \forall i, y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1.$$

what are the values of  $Q, A, \mathbf{a}, \mathbf{d}$  that express this problem as a QP in canonical form?

**Why is it interesting?** QP software solvers take QP in canonical form. To use a QP solver, you'll need to express the SVM problem as QP in canonical form as above.

5. ראשית, ניזכר כי מתקיים  $\|w\|^2 = w^T w = w^T I_n w$ , מעבר לכך, נשים לב שכל קורדינאטה  $y_i \langle w, x_i \rangle$  מיוצגת ע"י  $(x_i y_i)^T$  ואנו מנסים למזער את השגיאה כמה שיותר ולכן נבחר  $a = 0$  ונקבל:

$$\begin{aligned} \underset{(w, b)}{\operatorname{argmin}} \quad & \|w\|^2 \text{ s.t. } \forall i, y_i (\langle w, x_i \rangle + b) \geq 1 \\ = \underset{(w, b)}{\operatorname{argmin}} \quad & (w \ b) I_n \begin{pmatrix} w \\ b \end{pmatrix} \text{ s.t. } \begin{bmatrix} x_1 y_1 & 1 \\ \vdots & \vdots \\ x_m y_m & 1 \end{bmatrix} \begin{pmatrix} w \\ b \end{pmatrix} \geq \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \end{aligned}$$

$$= \operatorname{argmin}_{(w,b)} \frac{1}{2} \overbrace{\begin{pmatrix} w \\ b \end{pmatrix}}^{v^T} \overbrace{\begin{pmatrix} 1 \\ 0 \end{pmatrix}}^q \overbrace{\begin{pmatrix} w \\ b \end{pmatrix}}^v + \overbrace{0}^a \overbrace{\begin{pmatrix} w \\ b \end{pmatrix}}^v \quad s. t$$

$$\underbrace{\begin{bmatrix} x_1 y_1 & 1 \\ \vdots & \vdots \\ x_m y_m & 1 \end{bmatrix}}_A \underbrace{\begin{pmatrix} w \\ b \end{pmatrix}}_v \leq \underbrace{\begin{pmatrix} 1 \\ 1 \end{pmatrix}}_d$$

6. In the Soft-SVM we defined the problem:

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

Show that this problem is equivalent to the problem (namely that these problem have the same solutions)

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

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

6. ראשית, נשים לב כי השוני בין הביטויים נובע מהאיבר השני בחיבור ולא מהאיבר הראשון  $\left(\frac{\lambda}{2} \|\mathbf{w}\|^2\right)$ . הזהה בשניהם. על פי הביטוי הראשון מתקיים כי

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

וע"י העברת אגפים נקבל :

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

מצד שני, נתון כי  $\xi_i \geq 0$ , ולכן נוכל להגדיר :

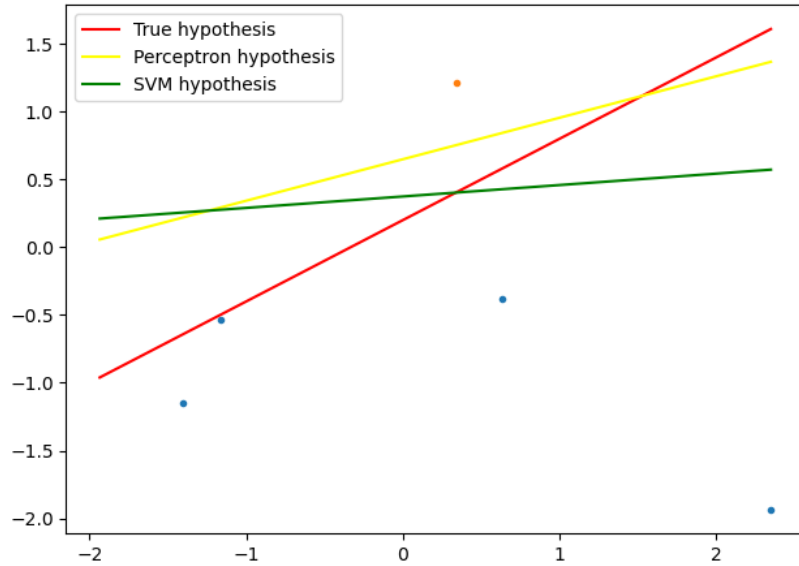
$$\xi_i \geq \max\{0, 1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle\}$$

נרצה למזער את האיבר השני בביטוי ולכן נוכל לבחור  $a = y_i \langle \mathbf{w}, \mathbf{x}_i \rangle$  ונקבל :

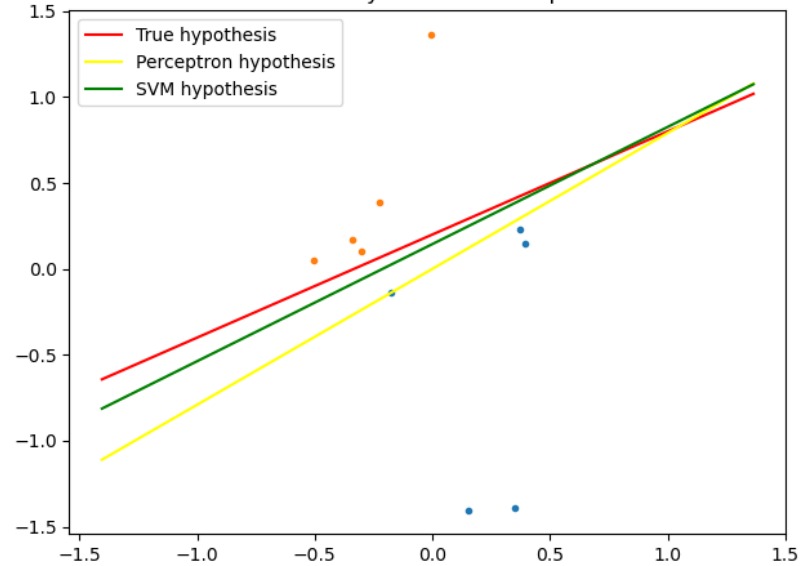
$$\operatorname{argmin}_w \frac{\lambda}{2} \|w\|^2 + \frac{1}{m} \sum_{i=1}^m \ell^{\text{hinge}}(a)$$

כנדרש.

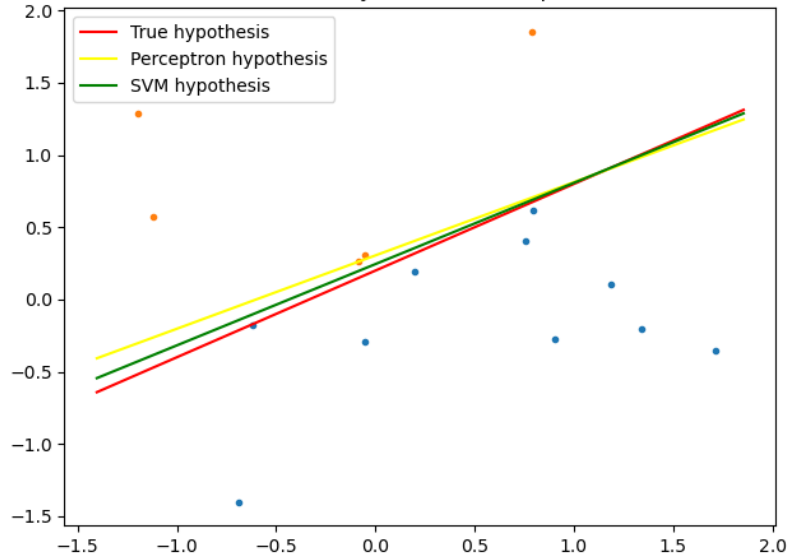
Data analysis over 5 samples



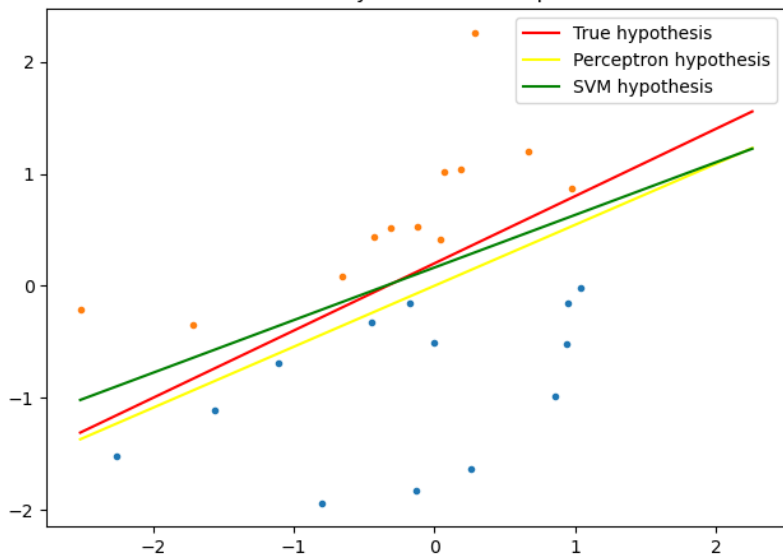
Data analysis over 10 samples .9



Data analysis over 15 samples



Data analysis over 25 samples



Data analysis over 70 samples

