

(20942) מבוא ללמידה חישובית | סיכום הרצאה 7

מנחה: שי מימון
סמסטר: 2022א'
נכתב על ידי: מתן כהן

1 Soft-Margin SVM

גם במקרה של SVM כמו עם הפרספטרון, נוכל להתקל במצבור דוגמאות שלא ניתן להפרדה מושלמת בצורה לינארית. במקרה וננסה לפתור את הבעיה עם Hard-Margin SVM נבחין כי זה בלתי אפשרי - כלומר לנסות לפתור בעיה בה קבוצת האימון שלנו לא ספרבילית בצורה מושלמת על ידי מפריד לינארי על פי המתכון של Hard-Margin SVM - **לא ניתן!**

1.1 הגדרת ה Soft-Margin SVM

ניזכר כי הבעיה הראשונית שלנו הייתה כזו:

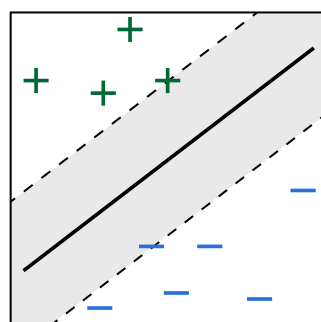
$$\begin{aligned} \underset{\mathbf{w}, b}{\text{minimize}} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{s.t} \quad & y_n \cdot (\mathbf{w}^T \mathbf{x}_n + b) \geq 1, \quad n = 1, \dots, N \end{aligned}$$

ואמרנו שהאילוצים שלנו לא יכולים להתקיים כאשר המידע שלנו לא ספרבילי. על מנת לגשר על הפער - נשנה את האילוצים כך שנרשה לחלק מהנקודות להיות בתוך השוליים שלנו ונגדיר מחדש את בעיית האופטימיזציה עם נקודות שיקראו "slack variables":

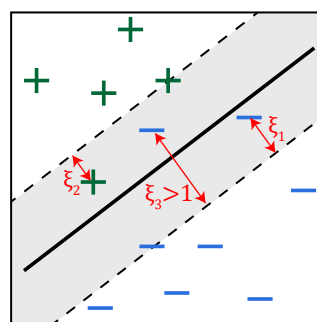
$$\begin{aligned} \underset{\mathbf{w}, b, \xi}{\text{minimize}} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \cdot \sum_{n=1}^N \xi_n \\ \text{s.t} \quad & y_n \cdot (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n, \quad n = 1, \dots, N \\ & \xi_n \geq 0, \quad n = 1, \dots, N \end{aligned}$$

כאשר:

- "נשלם" על כל נקודה שעברה את השוליים וננסה להביא את המחיר למינימום
 - ישנו trade-off בין הרחבת השוליים לבין צמצום הנקודות שנכנסו לשוליים
- כל נקודה \mathbf{x}_n עברה $\xi_n \in (0, 1)$ תניב לנו סיווג נכון אך בתוך השוליים
- כל נקודה \mathbf{x}_n עברה $\xi > 1$ הן נקודות שלא תסווגנה נכון



Hard-Margin



Soft-Margin

1.2 הגדרת הבעיה בצורה שונה

כעת, כפי שעשינו ב-Hard-Margin SVM נוכל גם כאן להשתמש בלגראנז'יאן. אסור לשכוח להביא את האילוצים לצורה $f_i \leq 0$ ולכן האילוח החדש שהגדרנו יהפוך ל:

$$y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n$$

$$\iff 1 - \xi_n - y_n(\mathbf{w}^T \mathbf{x}_n + b) \leq 0$$

$$\mathcal{L}(\mathbf{w}, b, \xi, \alpha, \beta) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \cdot \sum_{n=1}^N \xi_n + \sum_{n=1}^N \alpha_n (1 - \xi_n - y_n(\mathbf{w}^T \mathbf{x}_n + b)) - \sum_{n=1}^N \beta_n \xi_n$$

כעת נגזור את הלגראנז'יאן על פי \mathbf{w} , b ו- ξ ונשווה ל-0 על מנת למצוא מינימום

• על פי ξ :

$$\frac{\partial \mathcal{L}}{\partial \xi_i} = C - \alpha_i - \beta_i = 0 \Rightarrow C = \alpha_i + \beta_i$$

• על פי b :

$$\frac{\partial \mathcal{L}}{\partial b} = - \sum_{n=1}^N \alpha_n y_n = 0$$

• על פי \mathbf{w} :

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{w} - \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n \Rightarrow \mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$$

נכניס את הערכים למשוואה המקורית:

$$\begin{aligned} \mathcal{L}(\mathbf{w}, b, \xi, \alpha, \beta) &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \cdot \sum_{n=1}^N \xi_n + \sum_{n=1}^N \alpha_n (1 - \xi_n - y_n(\mathbf{w}^T \mathbf{x}_n + b)) - \sum_{n=1}^N \beta_n \xi_n \\ \backslash \beta_i = C - \alpha_i &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \cdot \sum_{n=1}^N \xi_n + \sum_{n=1}^N \alpha_n (1 - \xi_n - y_n(\mathbf{w}^T \mathbf{x}_n + b)) - \sum_{n=1}^N (C - \alpha_n) \cdot \xi_n \\ &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + \cancel{C \cdot \sum_{n=1}^N \xi_n} - \cancel{\sum_{n=1}^N \alpha_n \xi_n} + \sum_{n=1}^N \alpha_n (1 - y_n(\mathbf{w}^T \mathbf{x}_n + b)) - \cancel{C \cdot \sum_{n=1}^N \xi_n} + \cancel{\sum_{n=1}^N \alpha_n \xi_n} \\ &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{n=1}^N \alpha_n (1 - y_n(\mathbf{w}^T \mathbf{x}_n + b)) \\ &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{n=1}^N \alpha_n - \sum_{n=1}^N \alpha_n y_n b - \sum_{n=1}^N \alpha_n y_n \mathbf{w}^T \mathbf{x}_n \\ \backslash - \sum_{n=1}^N \alpha_n y_n b = 0 &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{n=1}^N \alpha_n - \mathbf{w}^T \cdot \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n \\ \backslash \mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n &= -\frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{n=1}^N \alpha_n - \mathbf{w}^T \mathbf{w} \\ &= -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y_n y_m \mathbf{x}_n^T \mathbf{x}_m + \sum_{n=1}^N \alpha_n \end{aligned}$$

קיבלנו בעיה זהה לבעית ה - Hard-Margin SVM ולכן נוכל לרשום גם בעית אופטימיזציה זהה עם האילוצים שקיבלנו בעת הגזירה והשוואה ל-0:

$$\begin{aligned} \min_{\alpha, \beta} \quad & \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y_n y_m \mathbf{x}_n^T \mathbf{x}_m - \sum_{n=1}^N \alpha_n \\ \text{s.t.} \quad & \alpha_n \geq 0, \quad n = 1, \dots, N \\ & \beta_n \geq 0, \quad n = 1, \dots, N \\ & \alpha_n + \beta_n = C, \quad n = 1, \dots, N \\ & \sum_{n=1}^N \alpha_n y_n = 0 \end{aligned}$$

נבחין כי: $\alpha_n + \beta_n = C \Rightarrow \alpha_n = C - \beta_n$ כיוון ש $\beta_n \geq 0$ הרי ש:

$$\alpha_n \leq C$$

וקיבלנו בעיה זהה לבעית ה Hard-Margin SVM למעט האילוץ על C:

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y_n y_m \mathbf{x}_n^T \mathbf{x}_m - \sum_{n=1}^N \alpha_n \\ \text{s.t.} \quad & \alpha_n \geq 0, \quad n = 1, \dots, N \\ & \alpha_n \leq C, \quad n = 1, \dots, N \\ & \sum_{n=1}^N \alpha_n y_n = 0 \end{aligned}$$

את בעיה זו ניתן לנסח בתור convex quadratic programming ולקבל $\alpha^* \leftarrow QP(Q_0, \mathbf{p}_0, A_0, \mathbf{c}_d)$ כמו כן:

$$\mathbf{w}^* = \sum_{\alpha_s^* > 0} \alpha_s^* \cdot y_s \cdot \mathbf{x}_s$$

כמו-כן ה-Complementary Slackness של ה-KKT:

1. $\alpha_n \cdot (1 - \xi_n - y_n(\mathbf{w}^T \mathbf{x}_n + b)) = 0, \quad n = 1, \dots, N$
2. $\beta_n \cdot \xi_n = (C - \alpha_n) = 0, \quad n = 1, \dots, N$

מכך נסיק:

$$\begin{aligned} \alpha_s^* > 0 : y_s(\mathbf{w}^{*T} \mathbf{x}_s + b^*) &= 1 - \xi_s \leq 1 \quad (1) \\ \beta > 0 \iff (C - \alpha_s^*) > 0 \iff \alpha_s^* < C : \xi_s = 0 \Rightarrow y_s(\mathbf{w}^{*T} \mathbf{x}_s + b^*) &\geq 1 \quad (2) \\ \text{מ- (1) ו-(2) נסיק שאם: } 0 < \alpha_s^* < C : y_s(\mathbf{w}^{*T} \mathbf{x}_s + b^*) &= 1 \quad (\alpha) \\ \text{לכן הנקודה } \mathbf{x}_s, y_s \text{ נמצאת על השוליים (margin support vector), נוכל לחלץ את } b^* : & \\ y_s(\mathbf{w}^{*T} \mathbf{x}_s + b^*) = 1 \iff b^* = y_s - \mathbf{w}^{*T} \mathbf{x}_s \quad (\alpha) & \\ \text{(5) לכל הנקודות עבורן } \alpha_s^* = C \text{ נקרא bounded/non-margin support vectors} & \end{aligned}$$

2 מידע ניתן להפרדה לא לינארית

2.1 שימוש בטרנספורמציות

ניזכר כי כאשר רצינו לבצע סיווג בעזרת פרספטרון למידע שניתן להפרדה אך לא בצורה לינארית, עשינו טרנספורמציה לפיצ'רים לסט פיצ'רים אחר או ממימד גבוה יותר ושם פתרנו את הבעיה בעזרת פרספטרון.

גם כאן נוכל להשתמש באותו "טריק".

נגדיר טרנספורמציה (לא בהכרח לינארית):

$$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{\tilde{d}}$$

כך שעבור סט פיצ'רים $\mathbf{x}_n \in \mathbb{R}^d$ נקבל:

$$\mathbf{z}_n = \phi(\mathbf{x}_n) \in \mathbb{R}^{\tilde{d}}$$

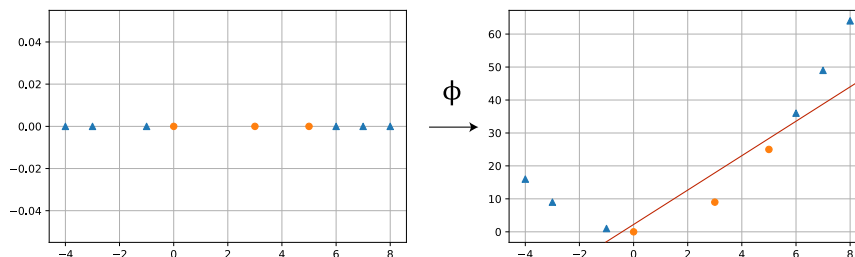
ונפתור את בעיית ה-Hard-Margin SVM:

$$\begin{aligned} \min_{\tilde{\mathbf{b}} \in \mathbb{R}, \tilde{\mathbf{w}} \in \mathbb{R}^{\tilde{d}}} \quad & \frac{1}{2} \tilde{\mathbf{w}}^T \tilde{\mathbf{w}} \\ \text{s.t.} \quad & y_n (\tilde{\mathbf{w}}^T \mathbf{z}_n + \tilde{b}) \geq 1, \quad n = 1, \dots, N \end{aligned}$$

נקבל QP $\begin{bmatrix} \tilde{b}^* \\ \tilde{\mathbf{w}}^* \end{bmatrix} \leftarrow$ והחזאי האופטימלי:

$$g(\mathbf{x}) = \text{sign}(\tilde{\mathbf{w}}^{*T} \phi(\mathbf{x}) + \tilde{b}^*)$$

דוגמה: שימוש בטרנספורמציה $\phi(x) = (x, x^2)$



בפועל - בד"כ נשתמש בטרנספורמציה לא לינארית וב-Soft-Margin SVM!

הערה: יש להזהר עם אופן הגדלת המימד מ-2 סיבות עיקריות:

- ככל שנגדיל את המימד ונקבל מימד שגודלו כמספר הנקודות שיש לנו בסט הדוגמאות - כך אנחנו עלולים לגרום ל-overfit - נדבר על כך בהמשך הקורס.
- במצבים מסויימים נגדיל את המימד עד כדי כך שהבעיה הופכת להיות מורכבת מבחינה חישובית, בהכרח אם נפתור בעיות ממרחב אינסופי

נפתור את בעיית גודל המימד בתתי הסעיפים הבאים.

2.1.1 התמודדות עם מימדים גדולים

על מנת להתמודד עם מימדים גדולים נוכל להשתמש בדואליות - שם אנחנו תלויים במספר הדוגמאות ולא במימדן! ניזכר כי הביטוי לבעיה הדואלית של ה-SVM הוא:

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y_n y_m \mathbf{x}_n^T \mathbf{x}_m - \sum_{n=1}^N \alpha_n \\ \text{s.t.} \quad & \alpha_n \geq 0, \quad n = 1, \dots, N \\ & \sum_{n=1}^N \alpha_n y_n = 0 \end{aligned}$$

כאשר:

- $\alpha \in \mathbb{R}^N$
- ישנם $N + 1$ אילוצים

לאחר המעבר למרחב $\mathbb{R}^{\tilde{d}}$ נרשום את הבעיה:

$$\begin{aligned} \min_{\tilde{\alpha} \in \mathbb{R}^{\tilde{d}}} \quad & \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \tilde{\alpha}_n \tilde{\alpha}_m y_n y_m \mathbf{z}_n^T \mathbf{z}_m - \sum_{n=1}^N \tilde{\alpha}_n \\ \text{s.t.} \quad & \tilde{\alpha}_n \geq 0, \quad n = 1, \dots, N \\ & \sum_{n=1}^N \tilde{\alpha}_n y_n = 0 \end{aligned}$$

נשתמש ב QP על מנת לפתור את הבעיה ולקבל $\tilde{\mathbf{w}}^* = \sum_{n=1}^N \tilde{\alpha}_n^* \cdot y_n \cdot \mathbf{z}_n$ ו- \tilde{b}^* ולקבל חזאי

$$\begin{aligned} g(\mathbf{x}) &= \text{sign}(\tilde{\mathbf{w}}^T \phi(\mathbf{x}) + \tilde{b}^*) \\ &= \text{sign}\left(\sum_{n=1}^N \tilde{\alpha}_n^* \cdot y_n \cdot \phi(\mathbf{x}_n) \cdot \phi(\mathbf{x}) + \tilde{b}^*\right) \end{aligned}$$

עם זאת, צריך להבחין בעובדה שהמכפלה $\mathbf{z}_n^T \mathbf{z}_m$ היא עדיין מכפלה של וקטורים ממימד \tilde{d} ונתמודד עם העניין בעזרת **Kernel Trick**

2.1.2 Kernel Trick

ראינו שעל מנת לפתור את בעית האופטימיזציה יש צורך לפתור פונקציות מהצורה הבאה:

$$\underbrace{K(\mathbf{x}', \mathbf{x}'')}_{\text{Kernel}} = \phi^T(\mathbf{x}') \cdot \phi(\mathbf{x}'')$$

בצורה הישירה היינו מחשבים את $\phi(\mathbf{x}')$ ואת $\phi(\mathbf{x}'')$ ומכפילים

אך ברצוננו לעשות דבר יעיל יותר, ו- Kernel Trick בא לענות על השאלה:

האם ניתן לחשב את המכפלה הפנימית מבלי לעבור דרך המרחב החדש שיצרנו?

התשובה לשמחתינו - היא כן!

דוגמה: נתבונן ב: $K(\mathbf{x}', \mathbf{x}'') = (1 + \mathbf{x}'^T \cdot \mathbf{x}'')^2$, $\mathbf{x}' \in \mathbb{R}^2, \mathbf{x}'' \in \mathbb{R}^2$
ונפתור:

$$\begin{aligned} K(\mathbf{x}', \mathbf{x}'') &= (1 + \mathbf{x}'^T \cdot \mathbf{x}'')^2, \quad \mathbf{x}' \in \mathbb{R}^2, \mathbf{x}'' \in \mathbb{R}^2 \\ &= \left(1 + \begin{bmatrix} x'_1 & x'_2 \end{bmatrix} \cdot \begin{bmatrix} x''_1 \\ x''_2 \end{bmatrix} \right)^2 \\ &= (1 + x'_1 x''_1 + x'_2 x''_2)^2 \\ &= 1 + (x'_1)^2 (x''_1)^2 + (x'_2)^2 (x''_2)^2 + 2x'_1 x''_1 + 2x'_2 x''_2 + 2 \cdot x'_1 x'_2 \cdot x''_1 x''_2 \\ &= \underbrace{\begin{bmatrix} 1 \\ \sqrt{2}x'_1 \\ \sqrt{2}x'_2 \\ (x'_1)^2 \\ (x'_2)^2 \\ \sqrt{2}x'_1 x'_2 \end{bmatrix}^T}_{\phi^T(\mathbf{x}')} \cdot \underbrace{\begin{bmatrix} 1 \\ \sqrt{2}x''_1 \\ \sqrt{2}x''_2 \\ (x''_1)^2 \\ (x''_2)^2 \\ \sqrt{2}x''_1 x''_2 \end{bmatrix}}_{\phi(\mathbf{x}'')} \end{aligned}$$

מסקנה: הצלחנו לחשב מכפלה פנימית שעולה לנו $O(d)$ ($d = 2$) ולהגיע למכפלה פנימית בממד הגבוה יותר \tilde{d}

ניתן להרחיב את מה שהראינו ל **Kernel פולינומיאלי**:

$$\begin{aligned} K(\mathbf{x}', \mathbf{x}'') &= \phi(\mathbf{x}') \cdot \phi(\mathbf{x}'') \\ &= (\xi + \zeta \cdot \mathbf{x}'^T \cdot \mathbf{x}'')^Q - \text{Valid Kernel} \end{aligned}$$

ניתן גם להרחיב עבור **Kernel גאוזי (RBF - Radial Basis Function)** שמדמה מעבר למרחב אינסופי:

$$K(\mathbf{x}', \mathbf{x}'') = \exp \left\{ -\gamma \cdot \|\mathbf{x}' - \mathbf{x}''\|^2 \right\}$$

הערה: על מנת ש-Kernel יהיה ולידי (Valid) יש צורך שיקיים:

- סימטריות
- לכל N ולכל x_1, \dots, x_N

$$\begin{bmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & K(\mathbf{x}_1, \mathbf{x}_2) & \dots & K(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \ddots & \vdots \\ K(\mathbf{x}_N, \mathbf{x}_1) & K(\mathbf{x}_N, \mathbf{x}_2) & \dots & K(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix} \succeq 0$$