

תרגול 10 : למידה בצורה PAC (Probably Approximately Correct)

תקציר התאוריה

סימונים :

- X - מרחב הכניסה.
- Y - מרחב היציאה ($Y = \mathbb{R}$) לבעיית הרגרסיה, $Y = \{-1, +1\}$ לבעיית הסיווג.
- $f_0 : X \rightarrow Y$ - פונקציית המטרה אותה רוצים ללמוד (עבור בעיית הסיווג זה המסווג הנלמד).
- $p_X(x)$ - פילוג הסתברות על מרחב הכניסה X . פילוג זה קובע את פילוג ה"כניסה הטיפוסית" $x \in X$.
- $\{x_k, y_k\}_{k=1}^n$ - דוגמאות (סדרת הלימוד).
- F - מרחב ההיפותזות, מתוכו עלינו לבחור את הפונקציה \hat{f} אשר משערכת את f_0 .

קריטריון ביצועים :

- מגדירים פונקציית מחיר $\ell(\hat{y}, y)$. עבור בעיית הרגרסיה נהוג לבחור ב-
 $\ell(\hat{y}, y) = (\hat{y} - y)^2$, ואילו עבור בעיית הסיווג ב- $\ell(\hat{y}, y) = \mathbb{I}\{\hat{y} \neq y\}$.
קריטריון ביצועים עבור פונקציה (היפותזה) כלשהי $\hat{f} \in F$ הוא המחיר הממוצע,
הנתון ע"י :

$$L(\hat{f}) \triangleq \mathbb{E} \{ \ell(\hat{f}(x), y) \}$$

כאשר התוחלת היא לפי $p_X(x)$ ו- $y = f_0(x)$.

- מטרתנו למצוא פונקציה (היפותזה) $\hat{f} \in F$ אשר מביאה את קריטריון הביצועים למינימום :

$$\hat{f}^* \in \arg \min_{f \in F} L(\hat{f})$$

נשים לב כי קריטריון זה לא ניתן לחישוב במקרה ולא יודעים את $p_X(x)$ ו- $f_0(x)$. במקרה זה משתמשים בסדרת הלימוד ומגדירים פונקציית המחיר האמפירי:

$$\hat{L}_n(\hat{f}) \triangleq \frac{1}{n} \sum_{k=1}^n \ell(\hat{f}(x_k), y_k)$$

ומחפשים:

$$\hat{f}_n = \arg \min_{f \in F} \hat{L}_n(\hat{f})$$

לסיכום, סימונים קריטריון ביצועים:

פונקציית מחיר	$\ell(\hat{y}, y)$
המחיר ממוצע	$L(\hat{f}) \triangleq \mathbb{E} \{ \ell(\hat{f}(x), y) \}$
המחיר הממוצע האמפירי	$\hat{L}_n(\hat{f}) \triangleq \frac{1}{n} \sum_{k=1}^n \ell(\hat{f}(x_k), y_k)$
הפונקציה האופטימאלית ב- $\hat{f} \in F$, והמחיר הממוצע עבור פונקציה זו.	$\hat{f}^* \in \arg \min_{f \in F} L(\hat{f})$ $L^* \triangleq L(\hat{f}^*) = \min_{f \in F} L(\hat{f})$
הפונקציה האופטימאלית האמפירית ב- $\hat{f} \in F$	$\hat{f}_n = \arg \min_{f \in F} \hat{L}_n(\hat{f})$

חסם PAC עבור קבוצת היפוטזות סופית ($|F| < \infty$):

זהו חסם הסתברותי המציין עד כמה הקירוב שלנו, \hat{f}_n , קרוב לפונקציה האופטימלית, \hat{f}^* , במובן קריטריון הביצועים $L(\cdot)$.
נסמן $L^* \triangleq \min_{f \in F} L(\hat{f})$. אזי:

$$(1) \quad \forall \varepsilon > 0: \mathbb{P} \{ L(\hat{f}_n) - L^* > \varepsilon \} < 2|F|e^{-\varepsilon^2 n/2}$$

אם בנוסף יודעים כי $L^* = 0$, אזי יש חסם הדוק יותר:

$$(2) \quad \forall \varepsilon > 0: \mathbb{P} \{ L(\hat{f}_n) > \varepsilon \} < |F|e^{-\varepsilon n}$$

שאלות

שאלה 1

נניח כי מעוניינים ללמוד מסווג בינארי $f_0 : \{-1, 1\} \times \{-1, 1\} \mapsto \{-1, 1\}$.
לצורך הלמידה, הוגרלו $n = 100$ דוגמאות בלתי תלויות, ונתונה משפחת
ההיפוטזות ממנה עלינו ללמוד את f_0 :

$$F = \left\{ f : \{-1, 1\} \times \{-1, 1\} \mapsto \{-1, 1\} \mid \sum_{a,b \in \{-1, 1\}} f(a, b) < 2 \right\}$$

א. אם ידוע כי המסווג הנלמד f_0 נתון ע"י:

a	b	f_0
-1	-1	1
-1	1	1
1	-1	1
1	1	-1

מצאו ε קטן ככל האפשר, כך שיתקיים:

$$\mathbb{P}\{L(\hat{f}_n) - L^* \leq \varepsilon\} > 0.98$$

ב. חשבו את L^* במקרה של סעיף א', אם ידוע כי $p_X \sim U(\{-1, 1\} \times \{-1, 1\})$.

ג. חזרו על א', אם ידוע כי המסווג הנלמד f_0 נתון ע"י:

a	b	f_0
-1	-1	-1
-1	1	-1
1	-1	-1
1	1	1

ד. אם נתון כי $p_X \sim U(\{-1, 1\} \times \{-1, 1\})$, מצאו את הסתברות לשגיאה 0.

פתרון 1

א. נחשב תחילה את $|F|$:

סה"כ יש $2^2 = 16$ פונקציות בינאריות בשני משתנים. F מכילה פונקציות המקיימות $\sum_{a,b \in \{-1,1\}} f(a,b) < 2$, ולכן F לא מכילה פונקציות בינאריות עם תוצאה אחת או פחות של -1 . מכאן, $|F| = 16 - 4 - 1 = 11$.

מכיוון שנדרש למצוא ε קטן ככל האפשר, נבדוק האם נוכל להשתמש בחסם ההדוק יותר ב- (2). כלומר נבדוק האם $f_0 \in F$. אבל, לפי הנתון, $\sum_{a,b \in \{-1,1\}} f_0(a,b) = 2$, ולכן $f_0 \notin F$. כלומר $L^* > 0$. מכאן שעלינו להשתמש בחסם (1) לעיל:

$$\mathbb{P}\{L(\hat{f}_n) - L^* \leq \varepsilon\} > 1 - 2|F|e^{-\varepsilon^2 n/2}$$

כלומר נדרוש:

$$1 - 2|F|e^{-\varepsilon^2 n/2} > 0.98$$

$$2|F|e^{-\varepsilon^2 n/2} < 0.02$$

$$\varepsilon > 0.374$$

(אפשר להשתמש ישיר בנוסחה של "מרווח הבטחון" של רשימות ההרצאה).

ב. לכל $f \in F$,

$$\begin{aligned} L(f) &= \mathbb{P}\{f(x) \neq f_0(x)\} = \sum_{x: f(x) \neq f_0(x)} p_X(x) \\ &= \frac{1}{4} |x: f(x) \neq f_0(x)|. \end{aligned}$$

$$\text{אבל: } \min_{f \in F} |x: f(x) \neq f_0(x)| = 1$$

והוא מתקבל עבור פונקציות $f \in F$ הבאות:

a	b	f	a	b	f	a	b	f
-1	-1	-1	-1	-1	1	-1	-1	1
-1	1	1	-1	1	-1	-1	1	1
1	-1	1	1	-1	1	1	-1	-1
1	1	-1	1	1	-1	1	1	-1

ולכן :

$$L^* = \min_{f \in F} L(f) = \frac{1}{4}$$

ג. במקרה זה כמובן $f_0 \in F$, כלומר $L^* = 0$, ולכן נוכל להשתמש בחסם הדוק יותר (2). כלומר :

$$\mathbb{P}\{L(\hat{f}_n) \leq \varepsilon\} > 1 - |F|e^{-\varepsilon n}$$

לכן נדרש :

$$1 - |F|e^{-\varepsilon n} > 0.98$$

$$|F|e^{-\varepsilon n} < 0.02$$

$$\varepsilon > 0.0631$$

כצפוי ה- ε המינימלי קטן יותר במקרה זה מאשר בסעיף א'.

ד. נתון כי כל כניסה מופיעה בהסתברות $1/4$:

הסתברות	b	a
$1/4$	-1	-1
$1/4$	1	-1
$1/4$	-1	1
$1/4$	1	1

נסמן את n הדגימות האקראיות (i.i.d.) ב $\{X_k\}_{k=1}^n$. אזי :

$$\begin{aligned} \mathbb{P}\left\{\begin{array}{l} \text{error}=0 \text{ after} \\ \text{seeing } n \text{ samples} \end{array}\right\} &= \mathbb{P}\left\{\begin{array}{l} n \text{ samples contains} \\ \text{all 4 possible inputs} \end{array}\right\} \\ &= 1 - \mathbb{P}\left\{\begin{array}{l} n \text{ samples doesn't contain} \\ \text{all 4 possible inputs} \end{array}\right\} \\ &= 1 - \mathbb{P}\left\{\begin{array}{l} \{X_1 \neq (-1, -1), \dots, X_n \neq (-1, -1)\} \cup \{X_1 \neq (-1, 1), \dots, X_n \neq (-1, 1)\} \cup \\ \cup \{X_1 \neq (1, -1), \dots, X_n \neq (1, -1)\} \cup \{X_1 \neq (1, 1), \dots, X_n \neq (1, 1)\} \end{array}\right\} \\ &\geq 1 - \mathbb{P}\{X_1 \neq (-1, -1), \dots, X_n \neq (-1, -1)\} - \mathbb{P}\{X_1 \neq (-1, 1), \dots, X_n \neq (-1, 1)\} \\ &\quad - \mathbb{P}\{X_1 \neq (1, -1), \dots, X_n \neq (1, -1)\} - \mathbb{P}\{X_1 \neq (1, 1), \dots, X_n \neq (1, 1)\} \\ &= 1 - \left(\frac{3}{4}\right)^n - \left(\frac{3}{4}\right)^n - \left(\frac{3}{4}\right)^n - \left(\frac{3}{4}\right)^n = 1 - 4\left(\frac{3}{4}\right)^n \end{aligned}$$

כאשר האי-שוויון נובע מחסם איחוד, והשוויון הלפני אחרון נובע מהעובדה כי הדגימות הן i.i.d. עבור $n = 100$ נקבל אם כן :

$$\mathbb{P}\left\{\begin{array}{l} \text{error}=0 \text{ after} \\ \text{seeing 100 samples} \end{array}\right\} \geq 1 - 4\left(\frac{3}{4}\right)^{100} = 0.9999... \approx 1$$

$$\Rightarrow \boxed{\mathbb{P}\left\{\begin{array}{l} \text{error}=0 \text{ after} \\ \text{seeing 100 samples} \end{array}\right\} \approx 1}$$

תוצאה זו מדגישה כי חסמי PAC הם כלליים (נכונים לכל פילוג של כניסה ופונקציית מטרה), וכי אם יש לנו מידע נוסף לגבי המודל, ניתן לקבל חסמים הרבה יותר טובים.

שאלה 2

כצעד ראשון לקראת הרחבה של חסמי PAC עבור קבוצות היפוטזות אינסופיות נרצה להעריך את ה"עושר"/"גיוון" של אוסף היפוטזות מסוים. נתחיל בלהראות את התכונה הבאה:

משפט (Steele Dudley) יהי G מרחב וקטורי של פונקציות ממשיות ב- R^d , נסמן את מימדו ב- $r = \dim(G)$. נגדיר את משפחת המסווגים (היפוטזות) הבאה:

$$G' = \{g'(x) = \text{sign}(g(x)) \mid g \in G\}$$

לכל קבוצה $\{x_i\}_{i=1}^m$ של $m = r + 1$ תבניות קיימת סדרת תיוגים בינאריים $\{y_i\}_{i=1}^m \in \{-1, 1\}^m$ כך שלא ניתן לסווג את הנקודות $\{x_i, y_i\}_{i=1}^m$ בצורה מושלמת ע"י מסווגים ב- G' .

$$\text{sign}(x) = \begin{cases} 1 & x \geq 0 \\ -1 & x < 0 \end{cases} \cdot \underline{\text{תזכורת:}}$$

סקיצת ההוכחה:

1. תהי קבוצת נקודות כלשהיא ב- R^d . $D = \{x_i\}_{i=1}^m$. בהמשך ההוכחה נבנה סדרת תיוגים $\{y_i\}_{i=1}^m$ כך שלא יהיה ניתן לסווג את הנקודות $\{x_i, y_i\}_{i=1}^m$ בצורה מושלמת ע"י מסווגים ב- G' .

2. נגדיר אופרטור לינארי $L : G \rightarrow R^m$ באופן הבא: $L(g) = (g(x_1) \dots g(x_m))^T$. נשים לב כי האופרטור $L(\cdot)$ הוא לינארי (הומוגני ואדיטיבי) למרות שהפונקציות $g \in G$ אינן בהכרח לינאריות.

3. נגדיר $L(G)$ להיות הטווח של האופרטור L מ- G , כלומר,

$$L(G) = \{x \in R^m \mid \exists g \in G : L(g) = x\}$$

$L(G)$ הינו תת מרחב של R^m , כלומר, סגור לחיבור וכפל בסקלר.

4. נשים לב כי מימד המרחב $L(G)$ אינו יכול להיות יותר גבוה ממימד של G .

לדוגמא, ניתן להראות בקלות כי אם $\{g_1, g_2, \dots, g_r\}$ הוא בסיס של G אזי אוסף

הוקטורים $\{L(g_1), \dots, L(g_r)\}$ פורש את המרחב $L(G)$. כלומר,

$$\dim L(G) \leq \dim(G) = r = m - 1 \text{ או במילים אחרות } \dim L(G) < m.$$

5. מהעובדות $L(G) \subseteq R^m$ ו- $\dim L(G) < \dim(R^m)$ נובע כי קיים $P = (P_1, \dots, P_m) \neq 0$ כך

ש- $L(g) \perp P$ לכל $g \in G$. באופן שקול ניתן לרשום כי

$$\sum_{i=1}^m P_i g(x_i) = 0 \quad (*)$$

לכל $g \in G$. נניח כי קיים לפחות רכיב אחד של P שלילי ממש, אחרת ניתן

להסתכל על וקטור $(-P)$ שהוא גם מאונך ל- $L(G)$.

6. נבחר את התיוגים הבאים: $y_i = \text{sign}(P_i)$. ע"מ לסיים את ההוכחה נותר להראות

כי לא קיים מסווג $g' \in G'$ שמצליח לסווג את הסדרה $\{x_i, y_i\}_{i=1}^m$ בצורה מושלמת.

נניח בשלילה כי קיים $\hat{g} \in G$ כך ש- $\hat{g} = \text{sign}(g')$ מסווג בצורה המושלמת את

הסדרה $\{x_i, y_i\}_{i=1}^m$. מזה נובע כי $\text{sign}(P_i) = \text{sign}(\hat{g}(x_i))$ לכל i , או, באופן שקול, ש-

$P_i \hat{g}(x_i) \geq 0$ לכל i . לכן, עבור j כך ש- $P_j < 0$ (לפי ההנחה סעיף 5 קיים אחד כזה)

מתקיים כי $g(x_j) < 0$ ולכן $P_j \hat{g}(x_j) > 0$.

נתבונן בסכום $\sum_{i=1}^m P_i \hat{g}(x_i)$ - כאן כל אחד מהגורמים הוא אי-שלילי ואחד

מהגורמים הוא חיובי ממש, לכן $\sum_{i=1}^m P_i \hat{g}(x_i) > 0$. אך זה עומד בסתירה לתנאי

הניצבות (*) בסעיף 5.

מש"ל