

Machine Learning from Data – IDC

HW6 – Theory

Nitai Aharoni - 203626742

1 VC-Dimension

Compute the VC-dimension of the following hypothesis classes:

1. (10 pt.) Assume the instance space X satisfies $|X| = \infty$. The space of binary hypotheses, which given a training set, returns the target y of \mathbf{x} if the pair (\mathbf{x}, y) was observed in the training set, and $+1$ otherwise. Formally, compute $VC(H)$ of $H = \{h : X \rightarrow \{-1, +1\} : h \text{ equals } -1 \text{ on a finite subset of } X \text{ and } +1 \text{ elsewhere}\}$.

- נתון מרחב X כך ש- $|X| = \infty$
- נגדיר את d להיות ה- training set כך ש- $d \subseteq X$ (תת קבוצה סופית)
- מהגדרת H , h תנתן באופן מושלם את קבוצת האימון d
- היות ו- $d \subseteq X$ אז $|d| \leq |X| = \infty$ ולכן עבור כל תת קבוצה d בכל גודל קיימת h שתנתן אותה באופן מושלם
- לכן $VC(H) = \infty$

2. (15 pt.) n -Interval classifiers of length ≥ 2 . Let $X = \mathbb{R}$,
 $H = \{x \rightarrow +1 \iff x \in [a_1, b_1] \cup [a_2, b_2] \cup \dots \cup [a_n, b_n] : a_1 + 2 \leq b_1, \dots, a_n + 2 \leq b_n\}$.

- $VC(H) = 2n$ נסביר:
 - נראה $VC(H) \geq 2n$
 - נתון מרחב X כך ש- $X = \mathbb{R}$
 - נגדיר את d להיות ה- training set כך ש- $d \subseteq X$ (קבוצה סופית)
 - מוגדר לנו classifier המורכב מ- n אינטרוולים על הישר באורך גדול שווה ל-2
 - נחלק את הנקודות ב- d לזוגות - כך שעבור כל $x \in d$ הזוגות (x_i, x_{i+1}) ממוקמים על הישר במרחק 1 ביניהם ובמרחק 10 מהזוג הבא.
 - עבור כל זוג ניתן לקבוע אינטרוול $[a_i, b_i]$ - כך ש- $a_i + 2 \leq b_i$, האינטרוול מנפץ את הזוג, וגם אינטרוול זה לא משפיע על ניפוצ הזוגות השכנים.
 - לכן בעזרת n אינטרוולים כאלה ניתן לנפץ $2n$ נקודות ולכן $VC(H) \geq 2n$
 - נראה $VC(H) < 2n + 1$
 - נוסף נקודה אחת ונסה לנפץ
 - כעת, היות ויש לנו $2n + 1$ נקודות (כלומר מספר אי זוגי של נקודות) ננסה לנפץ שלישית נקודות באמצעות אינטרוול אחד
 - נשים לב שאם נרצה לסווג את הנקודות באופן הבא: $+, -, +, -$ - לא ניתן לעשות זאת באמצעות אינטרוול אחד
 - כמובן שבאותו אופן גם רביעיית נקודות וכך הלאה לא ניתן לנפץ באמצעות אינטרוול אחד.
 - לכן $VC(H) = 2n$

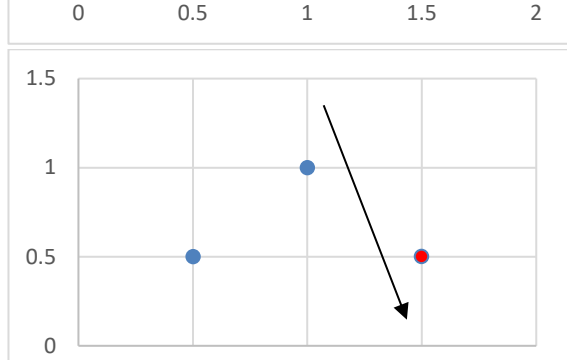
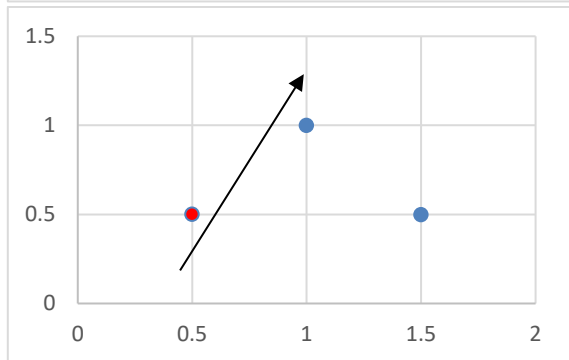
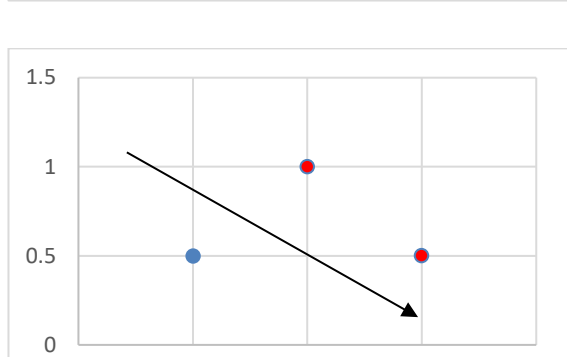
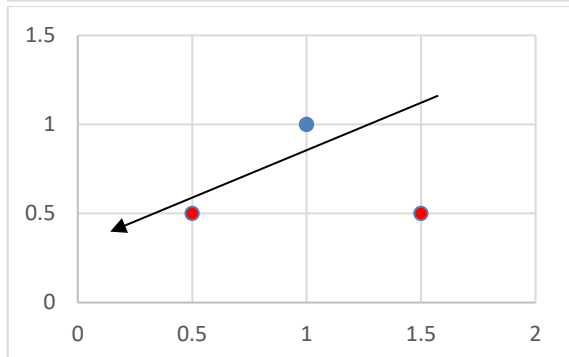
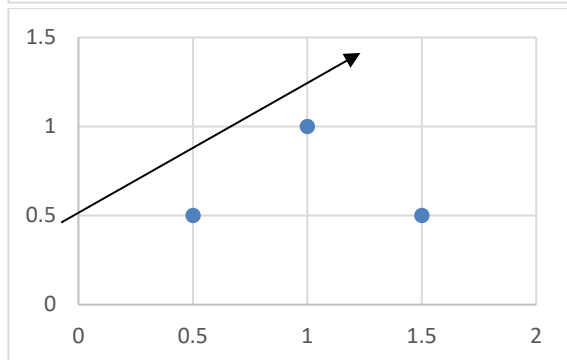
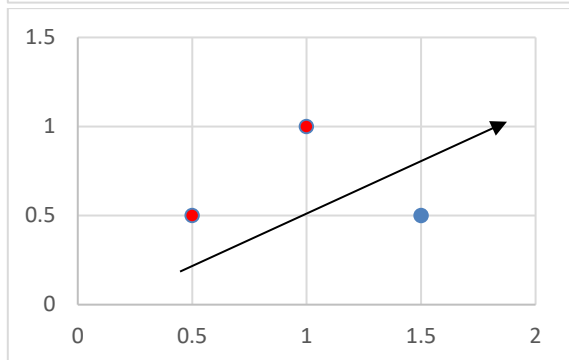
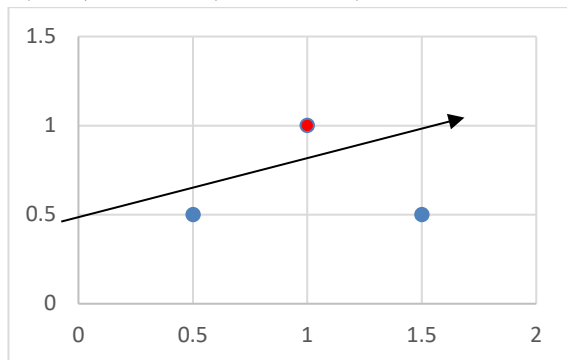
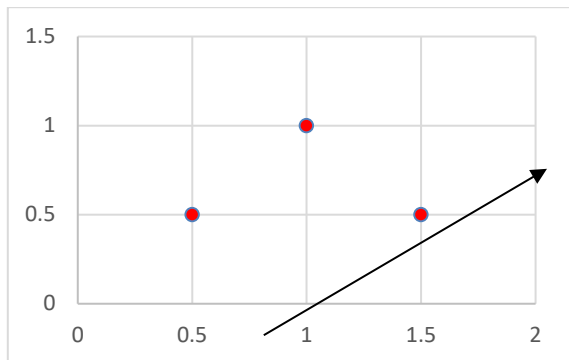
3. (20 pt.) Linear classifiers in the plain. Let $X = \mathbb{R}^2$,
 $H = \{(x_1, x_2) \rightarrow +1 \mid w_1 x_1 + w_2 x_2 + b > 0 \text{ or } -1 \mid w_1 x_1 + w_2 x_2 + b \leq 0\}$

$: w_1, w_2, b \in \mathbb{R}$

Show that $\text{VCdim}(\mathcal{H}) = 3$:

(a) Find a set of size 3 that \mathcal{H} shatters.

נראה פיזור של 3 נקודות כאלה כך ש-H מפץ אותן



(b) Show that no set of size 4, $A = (z_1, z_2, z_3, z_4)$, $z_i \in \mathbb{R}^2$ can be shattered by H .

Guidance: First prove the following lemma:

Lemma 1. Suppose a linear classifier h obtains prediction $y \in \{-1, +1\}$ on a set of points $z, z' \in \mathbb{R}^2$ ($h(z) = h(z') = y$). Then it also obtains the same prediction on any intermediate point. Namely,

$$\forall \alpha \in [0, 1] \quad h((1 - \alpha)z + \alpha z') = y.$$

And use it in each of the following 3 possible cases:

- The convex hull of A forms a line.
- The convex hull of A forms a triangle.
- The convex hull of A forms a quadrilateral.

• יהי $z = (x, y)$, $z' = (x', y')$

• נתון $h(z) = h(z') = y$

• נראה כי לכל $\alpha \in [0, 1]$ מתקיים $h((1 - \alpha)z + \alpha z') = y$

$$\begin{aligned} h((1 - \alpha)z + \alpha z') &= h((1 - \alpha)(x, y) + \alpha(x', y')) = h((1 - \alpha)x + (1 - \alpha)y + \\ &+ \alpha x' + \alpha y') = w_1((1 - \alpha)x + \alpha x') + w_2((1 - \alpha)y + \alpha y') + b = w_1((1 - \alpha)x + \\ &+ \alpha x') + w_2((1 - \alpha)y + \alpha y') + (1 - \alpha)b + \alpha b = (1 - \alpha)(w_1x + w_2y + b) + \\ &+ \alpha(w_1x' + w_2y' + b) = (1 - \alpha)h(z) + \alpha h(z') = (1 - \alpha)y + \alpha y = y \end{aligned}$$

• יהיו 4 נקודות: $z_1 = (x_1, y_1)$, $z_2 = (x_2, y_2)$, $z_3 = (x_3, y_3)$, $z_4 = (x_4, y_4)$

• נשתמש בכך עבור שלושת המקרים:

○ קו ישר:

▪ מהיות קו ישר אזי לכל i , $y_i = 0$

▪ ונניח $x_1 < x_2 < x_3 < x_4$

▪ ו- $h(z_1) = 1, h(z_2) = -1, h(z_3) = 1, h(z_4) = 1$

▪ לפי הלמה שהוכחנו אם $h(z_1) = 1$ וגם $h(z_3) = 1$ אזי בהכרח $h(z_2) = 1$ בסתירה

▪ לכן הקבוצה A לא יכולה להיות מנופצת ע"י H

○ משולש:

▪ כאשר z_1, z_2, z_3 יוצרות משולש ו- z_4 מוכלת במשולש

▪ ו- $h(z_1) = 1, h(z_2) = 1, h(z_3) = 1, h(z_4) = -1$

▪ אם z_4 נמצא על אחת מצלעות המשולש - בה"כ על הצלע שבין z_1 ל- z_3 :

• לפי הלמה שהוכחנו אם $h(z_1) = 1$ וגם $h(z_3) = 1$ אזי בהכרח $h(z_4) = 1$ בסתירה

• לכן הקבוצה A לא יכולה להיות מנופצת ע"י H

▪ אם z_4 נמצא בתוך המשולש:

• אזי קיימת נקודה z' כך ש- z' שנמצאת על אחת מצלעות המשולש, בה"כ בין z_1 ל- z_3 , כך ש z_4

נמצאת על הישר שבין z' לבין z_2 .

• לפי הלמה $h(z_1) = 1$ וגם $h(z_3) = 1$ אזי בהכרח $h(z') = 1$

• ולכן לפי הלמה שהוכחנו אם $h(z') = 1$ וגם $h(z_2) = 1$ אזי בהכרח $h(z_4) = 1$ בסתירה

• לכן הקבוצה A לא יכולה להיות מנופצת ע"י H

○ מרובע:

- כאשר $z1, z2, z3, z4$ יוצרות מרובע
- מדובר בבעיית XOR אותה לא ניתן לנפץ בעזרת מפריד לינארי כפי שהראינו בכיתה

2 Learning Conjunctions of Literals

(30 pt.) Let $X = \{0, 1\}^n$ (all Boolean strings of length n), let $C = H$ = the set of all conjunctions on X (e.g. $x_1 \wedge \neg x_3 \wedge x_n$ is in C and H). Define an algorithm L so that C is PAC-learnable by L using H . Prove all your steps.

• אלגוריתם:

- נגדיר $h = x_1 \wedge \overline{x_1} \wedge x_2 \wedge \overline{x_2} \wedge \dots \wedge x_n \wedge \overline{x_n}$
- נעבור על כל ה- $instances$:
- נתבונן ב- $instance$ ה- i .
- אם $C(instance_i)=1$, וכל האלמנטים ב- h הנוכחי קיימים ב- $attributes$ של ה- $instance$ אז h הוא קונסיסטנטי עם ה- $instance$ הנוכחי.
- אם $C(instance_i)=1$, ואחד או יותר מהאלמנטים ב- h הנוכחי לא מופיעים בו \leftarrow אז נסיר מ- h את כל ה- $attributes$ שקיימים ב- h ולא ב- $instance$.
- אחרי מעבר על כל ה- $instances$ נקבל h שהוא קונסיסטנטי עם C
- ואם קיבלנו h כזה אז C הוא PAC-learnable באמצעות L .

• נכונות:

• לומד קונסיסטנטי:

- עבור $C(instance_i)=1$ היות ומדובר בביטוי AND, כל אחד מה- $attributes$ של אותו $instance$ אמורים לחזות את ערכו האמיתי – ולכן אם קיימים ביטויים ב- h שלא מסכימים עם ה- $attributes$ של ה- $instance$, הם יוסרו מ- h
- עבור $C(instance_i)=0$ היות ומדובר בביטוי AND, אם כל הביטויים ב- h יופיעו ב- $attributes$ של ה- $instance$ אז h יהיה קונסיסטנטי עם אותו $instance$.
- נניח בשלילה שעבור $instance$ מסוים כל הביטויים ב- h קיימים ב- $attributes$ של ה- $instance$ אזי h לא קונסיסטנטי עם C . אבל נתון ש- $H=C$ ולכן קיים $h \in H$ כך ש- h קונסיסטנטי עם C בסתירה.
- אם עברנו על כל ה- $instances$ ו- h קונסיסטנטי עם כולם אז h קונסיסטנטי עם C .

• sample complexity

- $|H| = 3^n$ - ה- $attribute$ קיים בביטוי כ- $true$, קיים כ- $false$ או לא קיים. n =attribute num
- ולכן $m \geq \frac{n}{3} \ln 3 + \frac{1}{\epsilon} \ln(\frac{1}{\delta})$ כלומר sample complexity פולינומיאלי

• בדיקת זמן ריצה:

- עבור כל $instance$ נשווה את ה- $attributes$ שלו אל מול האלמנטים ב- h – ניתן לבדוק ב- $O(n)$
- נעבור על כל ה- $instances$ כלומר m
- סה"כ זמן ריצה $O(n*m)$ – פולינומיאלי
- לכן C הוא PAC-learnable באמצעות L

3 (Almost) PAC-learnability

(25 pt.) Let C denote the class of all possible target concepts defined over a set of instances X . Suppose that H is a space of binary hypotheses containing the constant concept c_1 defined by $c_1(x) = +1$ for all $x \in X$, and having the property that $C \setminus \{c_1\}$ is PAC-learnable by an algorithm L using H with sample complexity $m(\delta, \epsilon)$. Provide a learning algorithm L' that uses L , so that C (including c_1) is PAC-learnable by L' using H with sample complexity $\max\{m(\delta, \epsilon), \left\lceil \frac{\log(\frac{1}{\delta})}{\epsilon} \right\rceil\}$. Prove all your steps.

- אלגוריתם L' :
 - נעבור על כל ה- D instances ובדוק:
 - אם קיים $x \in D$ כך ש- $c(x) = 0 \leftarrow$ אז נריץ את L הידוע בתור PAC-learnable ונחזיר את h
 - אם לכל $x \in D$ מתקיים $c(x) = 1 \leftarrow$ אז נחזיר $h=1$
- נראה L' הינו PAC-learnable:
 - ידוע ש- L הינו PAC-learnable עם $m(\delta, \epsilon)$ sample complexity ולכן אין צורך להראות במקרה בו קיים $x \in D$ כך ש- $c(x) = 0$.
 - נראה עבור המקרה שלכל $x \in D$ מתקיים $c(x) = 1$
 - לומד קונסיסטנטי:
 - ידוע שלכל $x \in D$ מתקיים $c(x) = 1$ ולכן $h=1$ קונסיסטנטי עם c
 - sample complexity:
 - $|H| = 1$
 - ולכן $m \geq \frac{1}{\epsilon} (\ln|H| + \ln(\frac{1}{\delta})) \geq \frac{1}{\epsilon} \ln(\frac{1}{\delta})$ כלומר sample complexity פולינומיאלי
 - בדיקת זמן ריצה:
 - עבור כל $x \in D$ נבדוק את $c(x)$ לכן מדובר ב- m בדיקות $O(m)$
 - לכן C הוא PAC-learnable באמצעות L'
 - ידוע שבמקרה של c_1 מתקיים $m \geq \frac{1}{\epsilon} \ln(\frac{1}{\delta})$
 - לכן ה- sample complexity הוא $\max\{m(\delta, \epsilon), \left\lceil \frac{\log(\frac{1}{\delta})}{\epsilon} \right\rceil\}$