

## 67577) מבוא למערכות לומדות | תרגיל 4

שם: נמרוד בר גיורא | ת"ז: 207090622

### חלק תאורטי

### למידות PAC

### שאלה 1

1. For  $\mathcal{A}$  some learning algorithm,  $\mathcal{D}$  a probability distribution over  $\mathcal{X}$  and the 0-1 loss function (i.e misclassification), prove the following are equivalent:

- (a)  $\forall \epsilon, \delta > 0 \exists m(\epsilon, \delta) \text{ s.t. } \forall m \geq m(\epsilon, \delta) \mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S)) \leq \epsilon] \geq 1 - \delta$   
(b)  $\lim_{m \rightarrow \infty} \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S))] = 0$

הוכחה:  $(b \Rightarrow a)$  נניח ש- $\lim_{m \rightarrow \infty} \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S))] = 0$ . יהיו  $\epsilon, \delta > 0$ . מאי-שיויון מרקוב נובע שלכל  $m \in \mathbb{N}$  מתקיים:

$$\mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S)) > \epsilon] \leq \frac{\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S))]}{\epsilon}$$

מהגדרת הגבול נובע שקיים  $m(\epsilon, \delta) \in \mathbb{N}$  כך שלכל  $m \geq m(\epsilon, \delta)$  מתקיים ש-

$$\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S))] < \epsilon \cdot \delta$$

ולכן לכל  $m \geq m(\epsilon, \delta)$  מתקיים:

$$\mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S)) > \epsilon] < \frac{\epsilon \cdot \delta}{\epsilon} = \delta \xRightarrow{\text{המאורע המשלים}} \mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S)) \leq \epsilon] \geq 1 - \delta$$

$(a \Rightarrow b)$  נניח שלכל  $\epsilon, \delta > 0$  קיים  $m(\epsilon, \delta) \in \mathbb{N}$  כך שלכל  $m \geq m(\epsilon, \delta)$  מתקיים ש- $\mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S)) \leq \epsilon] \geq 1 - \delta$ . יהיו  $\epsilon, \delta > 0$ . נשתמש בנוסחת התוחלת השלמה על  $\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S))]$  ונקבל לכל  $m \geq m(\epsilon, \delta)$  ש-

$$\begin{aligned} \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S))] &= \mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S)) \leq \epsilon] \cdot \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S)) \mid L_{\mathcal{D}}(\mathcal{A}(S)) \leq \epsilon] + \\ &+ \mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S)) > \epsilon] \cdot \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S)) \mid L_{\mathcal{D}}(\mathcal{A}(S)) > \epsilon] \end{aligned}$$

נחסום כל אחד מהביטויים לעיל בנפרד:

מכך ש- $\mathbb{P}_{S \sim \mathcal{D}^m}$  היא פונקציית הסתברות נובע ש-

$$\mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S)) \leq \epsilon] \leq 1$$

מההנחה של  $a$  (תוך שימוש במאורע המשלים) נובע ש-

$$\mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S)) > \epsilon] = 1 - \mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S)) \leq \epsilon] \stackrel{a}{\leq} 1 - (1 - \delta) = \delta$$

ממונוטוניות התוחלת, ומכך שזו תוחלת מותנית נובע ש-

$$\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S)) \mid L_{\mathcal{D}}(\mathcal{A}(S)) \leq \varepsilon] \leq \mathbb{E}_{S \sim \mathcal{D}^m} [\varepsilon] = \varepsilon$$

לפי ההגדרה -  $L_{\mathcal{D}} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [h(\mathbf{x}) \neq f(\mathbf{x})]$ . לכן מכיוון ש-  $\mathbb{1}_{[h(\mathbf{x}) \neq f(\mathbf{x})]} \leq 1$  אז ממונוטוניות התוחלת (מופעלת על  $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}$  ואז על  $\mathbb{E}_{S \sim \mathcal{D}^m}$ ):

$$\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S)) \mid L_{\mathcal{D}}(\mathcal{A}(S)) > \varepsilon] \leq \mathbb{E}_{S \sim \mathcal{D}^m} [1] = 1$$

ולכן קיבלנו שלכל  $m \geq m(\varepsilon, \delta)$  מתקיים -

$$\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S))] \leq 1 \cdot \varepsilon + \delta \cdot 1 = \varepsilon + \delta$$

ומכאן לפי הגדרת הגבול -  $\lim_{m \rightarrow \infty} \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S))] = 0$ .

## שאלה 2

2. Let  $\mathcal{X} := \mathbb{R}^2$ ,  $\mathcal{Y} := \{0, 1\}$  and let  $\mathcal{H}$  be the class of concentric circles in the plane, i.e.,

$$\mathcal{H} := \{h_r : r \in \mathbb{R}_+\} \quad \text{where} \quad h_r(\mathbf{x}) = \mathbb{1}_{\|\mathbf{x}\|_2 \leq r}$$

Prove that  $\mathcal{H}$  is PAC-learnable and its sample complexity is bounded by

$$m_{\mathcal{H}}(\varepsilon, \delta) \leq \frac{\log(1/\delta)}{\varepsilon}$$

אלגוריתם למידה: בהינתן סט דגימות  $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$  האלגוריתם  $\mathcal{A}$  יחזיר את ההיפותרזה  $\mathcal{A}(S) = h_{\hat{r}} \in \mathcal{H}$  כאשר:

$$\hat{r} = \max_{i \in [m]: y_i = 1} \|\mathbf{x}_i\|_2$$

כאשר  $\max \emptyset = 0$ .

כלומר האלגוריתם יחזיר את המעגל סביב הראשית עם הרדיוס הגדול ביותר שמכיל את הדגימות שהלייבל שלהן הוא 1, מתוך סט הדגימות הנתון. במילים אחרות - הרדיוס של המעגל יקבע לפי הדגימה  $\mathbf{x}_i$  כך ש-  $y_i = 1$  ו-  $\|\mathbf{x}_i\|_2$  מקסימלי. הוכחת נכונות:

**הוכחה:** נראה שלכל מעגל סביב הראשית ברדיוס  $r$ , לכל הסתברות  $\mathcal{D}$  מעל  $\mathcal{X}$  ולכל  $\varepsilon, \delta \in (0, 1)$  אם נדגום  $m$  דגימות i.i.d לפי ההתפלגות  $\mathcal{D}$  אז עם הסתברות של לפחות  $1 - \delta$  המעגל סביב הראשית ברדיוס  $\hat{r}$  שיוחזר ע"י האלגוריתם לעיל הוא בעל שגיאה של לכל היותר  $\varepsilon$ . אז יהיו  $\mathcal{D}, \varepsilon, \delta$  כנ"ל. נשים לב תחילה שהרדיוס של המעגל שמוחזר ע"י האלגוריתם תמיד יהיה קטן/שווה לרדיוס האמיתי -  $\hat{r} \leq r$ , ולכן:

$$L_{\mathcal{D}}(h_{\hat{r}}) = \mathbb{P}(\|\mathbf{x}_i\|_2 \in (\hat{r}, r])$$

נרצה למצוא מספר דגימות גדול מספיק  $m$  שיבטיח ש-  $\mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_{\hat{r}}) \leq \varepsilon] \geq 1 - \delta$ .

נתבונן לצורך כך בהסתברות של המאורע המשלים -  $\mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_{\hat{r}}) > \varepsilon]$ . נשים לב שאם ההסתברות לדגום נקודה  $\mathbf{x}_i \in \mathbb{R}^2$  שנמצאת בהפרש בין המעגל האמיתי (ברדיוס  $r$ ) למעגל שהחזיר  $\mathcal{A}$  (ברדיוס  $\hat{r}$ ) היא גדולה מ- $\varepsilon$ , כלומר אם  $\mathbb{P}(\|\mathbf{x}_i\|_2 \in (\hat{r}, r]) > \varepsilon$ , אז ההסתברות שבסט הדגימות  $S$  לא תהיה אף נקודה באיזור הזה היא לפחות  $(1 - \varepsilon)^m$ . פורמלית:

$$\mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_{\hat{r}}) > \varepsilon] = \mathbb{P}_{S \sim \mathcal{D}^m} \left[ \bigwedge_{i=1}^m \|\mathbf{x}_i\|_2 \notin (\hat{r}, r] \right] \stackrel{\text{i.i.d}}{=} \prod_{i=1}^m \mathbb{P}_{\mathbf{x}_i \sim \mathcal{D}} (\|\mathbf{x}_i\|_2 \notin (\hat{r}, r]) \leq (1 - \varepsilon)^m \leq e^{-\varepsilon m}$$

מכאן נובע ש-

$$\mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_{\hat{r}}) \leq \varepsilon] = 1 - \mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_{\hat{r}}) > \varepsilon] \geq 1 - e^{-\varepsilon m}$$

לכן מספר הדגימות צריך לקיים:

$$1 - e^{-\varepsilon m} \geq 1 - \delta \implies e^{-\varepsilon m} \leq \delta \implies -\varepsilon m \leq \ln(\delta) \implies m \geq \frac{-\ln(\delta)}{\varepsilon} = \frac{\ln(\delta^{-1})}{\varepsilon} = \frac{\ln(\frac{1}{\delta})}{\varepsilon}$$

כלומר, לכל  $m \geq \frac{\ln(\frac{1}{\delta})}{\varepsilon}$  נקבל ש-

$$\mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_{\hat{f}}) \leq \varepsilon] \geq 1 - \delta$$

ומכיון ש- $\mathcal{D}$ ,  $\varepsilon$ ,  $\delta$  היו שרירותיים -  $\mathcal{H}$  היא למידה-PAC וה-sample complexity שלה מקיים:

$$m_{\mathcal{H}}(\varepsilon, \delta) \leq \frac{\ln(\frac{1}{\delta})}{\varepsilon}$$

■

## VC-Dimension

### שאלה 3

3. Let  $\mathcal{X} = \{0, 1\}^n$  and  $\mathcal{Y} = \{0, 1\}$ , for each  $I \subseteq [n]$  define the parity function:

$$h_I(\mathbf{x}) = \left( \sum_{i \in I} x_i \right) \bmod 2.$$

What is the VC-dimension of the class  $\mathcal{H}_{\text{parity}} = \{h_I \mid I \subseteq [n]\}$ ? Prove your answer.

ה-VC-dimension של המחלקה  $\mathcal{H}_{\text{parity}} = \{h_I : I \subseteq [n]\}$  הוא  $n$ .

**הוכחה:** נשים לב תחילה שלכל  $i \subseteq [n]$  מתקיים לכל  $\mathbf{x} \in \{0, 1\}^n$  ש-

$$h_I(\mathbf{x}) = \left( \sum_{i \in I} x_i \right) \bmod 2 = \langle \mathbf{x} \mid v_I \rangle \bmod 2$$

כאשר  $v_I \in \{0, 1\}^n$  הוא וקטור המוגדר ע"י:

$$[v_I]_i = \begin{cases} 1 & i \in I \\ 0 & i \notin I \end{cases}$$

כלומר הוא וקטור שהקואורדינטה ה- $i$  שלו מציינת האם האינדקס  $i$  שייך לקבוצת האינדקסים  $I$ .

לכן לכל  $I \subseteq [n]$  מתקיים ש- $h_I \in \mathcal{H}$  אם ורק אם קיים לה וקטור  $v_I$  כנ"ל.

נראה שהקבוצה  $C = \{0, 1\}^n \supseteq \{e_1, \dots, e_n\}$  (כאשר לכל  $i \in [n]$  הוקטור  $e_i$  הוא וקטור היחידה ה- $i$ ) מנותצת ע"י  $\mathcal{H}$ :

לכל  $n$  לייבלים  $y_1, \dots, y_n \in \{0, 1\}$  עבור הוקטורים ב- $C$  נגדיר

$$v_I = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

ונקבל לכל  $i \in [n]$  ש-

$$h_I(e_i) = \langle e_i \mid v_I \rangle \bmod 2 = [v_I]_i \bmod 2 = y_i \bmod 2 \stackrel{y_i \in \{0, 1\}}{=} y_i$$

כלומר  $\text{VC-DIM}(\mathcal{H}) \geq n$  ולכן  $C$  הנ"ל, ולכן  $\text{VC-DIM}(\mathcal{H}) \geq n$ .  
 תהי כעת קבוצה  $\{0, 1\}^n \supseteq C = (\mathbf{x}_1, \dots, \mathbf{x}_{n+1})$  בגודל  $n+1$ . נשים לב ש- $\{0, 1\}^n \subseteq \mathbb{R}^n$ , ולכן (כאשר חושבים על  $\mathbb{R}^n$  כעל מרחב אוקלידי)  $C$  היא בהכרח קבוצה תלויה ליניארית. לכן קיים צירוף ליניארי לא-טריוויאלי מתאפס של איברי  $C$ :

$$\sum_{i=1}^{n+1} a_i \mathbf{x}_i = \bar{0}$$

נניח בה"כ ש- $a_{n+1} \neq 0$  ונקבל ש-

$$\mathbf{x}_{n+1} = \sum_{i=1}^n b_i \mathbf{x}_i \quad \text{☀}$$

כאשר  $b_i := -\frac{a_i}{a_{n+1}}$  לכל  $i \in [n]$ . נתבונן כעת בלייבלים הבאים עבור  $C$ :

$$\forall i \in [n+1]: \quad y_i = \begin{cases} b_i \bmod 2 & i \in [n] \\ (1 + \sum_{j=1}^n (b_j \bmod 2)) \bmod 2 & i = n+1 \end{cases}$$

נניח בשלילה שקיים  $v_I \in \{0, 1\}^n$  כך שלכל  $i \in [n+1]$  מתקיים:

$$y_i = \langle \mathbf{x}_i \mid v_I \rangle \bmod 2 = \begin{cases} b_i \bmod 2 & i \in [n] \\ (1 + \sum_{j=1}^n (b_j \bmod 2)) \bmod 2 & i = n+1 \end{cases}$$

ונקבל (תוך שימוש באריתמטיקה של mod) ש-

$$\begin{aligned} \langle \mathbf{x}_{n+1} \mid v_I \rangle \bmod 2 &\stackrel{\text{☀}}{=} \left\langle \sum_{i=1}^n b_i \mathbf{x}_i \mid v_I \right\rangle \bmod 2 = \left( \sum_{i=1}^n b_i \langle \mathbf{x}_i \mid v_I \rangle \right) \bmod 2 = \sum_{i=1}^n ((b_i \bmod 2) \langle \mathbf{x}_i \mid v_I \rangle \bmod 2) = \\ &= \sum_{i=1}^n ((b_i \bmod 2) (b_i \bmod 2)) = \sum_{i=1}^n (b_i \bmod 2)^2 = \sum_{i=1}^n (b_i \bmod 2) = \\ &\quad \text{סדרה נכונה בשלילה} \quad \quad \quad \text{לפי הבחירה של הלייבלים} \\ &= \left( \sum_{i=1}^n (b_i \bmod 2) \right) \bmod 2 \neq \left( 1 + \sum_{i=1}^n (b_i \bmod 2) \right) \bmod 2 = y_{n+1} \end{aligned}$$

וזה סתירה. לכן לא ניתן לנתץ תת-קבוצות בגודל  $n+1$  של  $\{0, 1\}^n$ , וזה אומר ש- $\text{VC-DIM}(\mathcal{H}) < n+1$ , ובסך הכל קיבלנו ש-

$$\text{VC-DIM}(\mathcal{H}) = n$$

■

4. Given an integer  $k$ , let  $([a_i, b_i])_{i=1}^k$  be any set of  $k$  intervals on  $\mathbb{R}$  and define their union  $A = \cup_{i=1}^k [a_i, b_i]$ . The hypothesis class  $\mathcal{H}_{k\text{-intervals}}$  includes the functions:  $h_A(x) := \mathbb{1}_{[x \in A]}$ , for all choices of  $k$  intervals. Find the VC-dimension of  $\mathcal{H}_{k\text{-intervals}}$  and prove your answer. Show that if we let  $A$  be any finite union of intervals (i.e.  $k$  is unlimited), then the resulting class  $\mathcal{H}_{\text{intervals}}$  has VC-dimension  $\infty$ .

ה-VC-dimension של המחלקה  $\mathcal{H}_{k\text{-intervals}}$  הוא  $2k$ .

**הוכחה:** נראה שהקבוצה  $C = \{x_1, \dots, x_{2k}\} \subseteq \mathbb{R}$  כאשר  $x_i = i$  לכל  $i \in [2k]$  (קבוצת כל המספרים הטבעיים עד  $2k$ , כולל) מנותצת ע"י  $\mathcal{H}_{k\text{-intervals}}$ .

לכל  $2k$  לייבלים  $y_1, \dots, y_{2k}$  עבור איברי  $C$  נוכל לבנות היפותזה  $h_A \in \mathcal{H}_{k\text{-intervals}}$  כך ש-  $h_A(x_i) = y_i$  לכל  $i \in [2k]$ , באופן הבא:  
אם  $y_i = 0$  לכל  $i \in [2k]$  אז נגדיר את  $A$  להיות איחוד של  $k$  קטעים שמוכלים ב-  $(-\infty, 0]$ , ונקבל ש-  $x_i \notin A$  לכל  $i \in [2k]$ , כלומר ש-  $h_A(x_i) = 0 = y_i$  לכל  $i$ .

אחרת, הקבוצה  $A$  שעליה מבוססת  $h_A$  תהיה איחוד של קטעים המכסים רצפים של  $x_i$ ים מתוך  $C$  שהלייבל שלהם הוא 1. ניתן לבנות אותה באופן הבא:

בכל שלב  $\ell \in [k]$ , אם קיים  $i \in [2k] \setminus [b_{\ell-1}]$  כך ש-  $y_i = 1$  נגדיר את הקטע  $[a_\ell, b_\ell]$  ע"י-

$$a_\ell = \min \{i \in [2k] \setminus [b_{\ell-1}] : y_i = 1\}$$

$$b_\ell = \max \{i \in [2k] \setminus [a_\ell] : y_j = 1 \ \forall j \in [i] \setminus [a_\ell], y_{i+1} = 0\}$$

כאשר מגדירים  $b_0 = 0$  (ולכן  $[b_0] = \emptyset$ ) ו-  $y_{2k+1} = 0$ .

אם לא קיים  $i$  כנ"ל אז הבניה תעצר (ופורמלית נגדיר  $[a_\ell, b_\ell] = \emptyset$ ). מוגדרת ע"י:

$$A = \bigcup_{\ell=1}^k [a_\ell, b_\ell]$$

מכיוון שבמקרה הזה אנחנו מניחים שקיים  $i \in [2k]$  כך ש-  $y_i = 1$ , אז מתקבל בבניה הזאת לפחות קטע  $[a_\ell, b_\ell]$  אחד.  
נשים לב שמספר הקטעים הוא לכל היותר  $k$  (זו הסיבה לכך ש-  $\ell \in [k]$ ) כי בין כל שני קטעים  $[a_\ell, b_\ell], [a_{\ell+1}, b_{\ell+1}]$  קיים לפי הבניה  $i \in [a_{\ell+1} - 1] \setminus [b_\ell]$  (לפחות אחד) כך ש-  $y_i = 0$ . לכן אם נניח בשלילה שניתן לבנות כך  $k+1$  קטעים, נקבל שיש ב-  $C$  לפחות  $k$  נקודות  $x_i$  כך ש-  $y_i = 0$  ובנוסף בכל אחד מהקטעים יש לפחות נקודה אחת  $x_i$  כך ש-  $y_i = 1$ , ומכאן נקבל שיש ב-  $C$  לפחות  $2k+1$  נקודות, בסתירה להגדרתה.  
עכשיו מהגדרת  $A$  נובע שלכל  $x_i \in C$  מתקיים ש-  $x_i \in A$  אם ורק אם  $y_i = 1$ , או באופן שקול  $x_i \notin A$  אם ורק אם  $y_i = 0$ , ולכן

$$h_A(x_i) = y_i$$

כלומר  $\mathcal{H}_{k\text{-intervals}}$  מנתצת את  $C$ , ולכן  $\text{VC-DIM}(\mathcal{H}_{k\text{-intervals}}) \geq 2k$ .

תהי כעת קבוצה  $C = \{x_1, \dots, x_{2k+1}\} \subseteq \mathbb{R}$  כלשהי. נתבונן בלייבלים הבאים עבור הנקודות ב-  $C$ :

$$y_i = \begin{cases} 1 & i \equiv 1 \pmod{2} \\ 0 & i \equiv 0 \pmod{2} \end{cases}$$

$$\frac{1}{y_1}, \frac{0}{y_2}, \frac{1}{y_3}, \dots, \frac{0}{y_{2k}}, \frac{1}{y_{2k+1}}$$

כלומר הלייבלים הם 0 לכל  $x_i$  עם אינדקס זוגי ו-1 לכל  $x_i$  עם אינדקס אי-זוגי. לכן יש ב-  $C$  בדיוק  $k$  נקודות עם לייבל 0 ו-  $k+1$  נקודות עם לייבל 1.

נניח בשלילה שקיימת  $h_A \in \mathcal{H}_{k\text{-intervals}}$  כך ש- $h_A(x_i) = y_i$  לכל  $i \in [2k+1]$ . מכיוון שיש  $k+1$  נקודות ב- $C$  עם לייבל 1, ב- $A$  יש  $k$  קטעים, וכל נקודה עם לייבל 1 חייבת להיות מוכלת בקטע כלשהו ב- $A$  אז (מעיקרון שובך היונים) קיים קטע  $[a_\ell, b_\ell] \subseteq A$  וקיימות  $x_i, x_j \in C$  **שווים** (בה"כ  $x_i < x_j$ ) כך ש-

$$x_i, x_j \in [a_\ell, b_\ell] \quad \wedge \quad y_i = y_j = 1$$

מבחרת הלייבלים נובע שקיימת נקודה  $x_d \in C$  כך ש- $x_i < x_d < x_j$  וגם  $y_d = 0$ , ולכן:

$$y_d \in [a_\ell, b_\ell] \subseteq A$$

כלומר מתקיים ש- $h_A(x_d) = 1 \neq y_d$ , וזו סתירה.

לכן כל קבוצה בגודל  $2k+1$  לא ניתנת לניתוח ע"י  $\mathcal{H}_{k\text{-intervals}}$ , כלומר  $\text{VC-DIM}(\mathcal{H}_{k\text{-intervals}}) < 2k$  ובסה"כ קיבלנו ש-

$$\text{VC-DIM}(\mathcal{H}_{k\text{-intervals}}) = 2k$$

■

מימד ה-VC של מחלקת ההיפותוזות  $\mathcal{H}_{\text{intervals}}$  (שבה אין חסם על מספר הקטעים ב- $A$ ) הוא  $\infty$  כי: לכל  $n \in \mathbb{N}$  הקבוצה  $\mathbb{R} \supseteq C = \{x_1, \dots, x_{2n}\}$  כאשר  $x_i = i$  לכל  $i \in [2n]$  (קבוצת כל המספרים הטבעיים עד  $2n$ , כולל) **ניתנת לניתוח** ע"י  $\mathcal{H}_{\text{intervals}}$ , כי לכל סדרת לייבלים  $y_1, \dots, y_{2n}$  עבור הנקודות ב- $C$  נוכל לבנות את  $A$  כאיחוד של קטעים **בדיוק** כמו שעשינו בחלק הראשון של ההוכחה לעיל. נצטרך להשתמש ב- $n$  קטעים לכל היותר בבניה הזאת, וזה אפשרי מכיוון שאין חסם על מספר הקטעים ב- $A$  במחלקת ההיפותוזות הזו.

מכיוון שלכל  $n \in \mathbb{N}$  קיימת קבוצה  $C \subseteq \mathbb{R}$  בגודל  $|C| = 2n$  הניתנת לניתוח ע"י  $\mathcal{H}_{\text{intervals}}$ , אז

$$\forall n \in \mathbb{N} : \quad \text{VC-DIM}(\mathcal{H}_{\text{intervals}}) \geq 2n \quad \implies \quad \text{VC-DIM}(\mathcal{H}_{\text{intervals}}) = \infty$$

## מונוטוניות

### שאלה 5

5. Let  $\mathcal{H}$  be a hypothesis class for a binary classification task. Suppose that  $\mathcal{H}$  is PAC learnable and its sample complexity is given by  $m_{\mathcal{H}}(\cdot, \cdot)$ . Show that  $m_{\mathcal{H}}$  is monotonically non-increasing in each of its parameters. That is:

- Show that given  $\delta \in (0, 1)$ , and given  $0 < \varepsilon_1 \leq \varepsilon_2 < 1$ , we have that  $m_{\mathcal{H}}(\varepsilon_1, \delta) \geq m_{\mathcal{H}}(\varepsilon_2, \delta)$ .
- Similarly, show that given  $\varepsilon \in (0, 1)$ , and given  $0 < \delta_1 \leq \delta_2 < 1$ , we have that  $m_{\mathcal{H}}(\varepsilon, \delta_1) \geq m_{\mathcal{H}}(\varepsilon, \delta_2)$ .

• **הוכחה:** נתונה  $\mathcal{H}$  למידה PAC. יהי  $\mathcal{A}$  אלגוריתם שלומד אותה, יהי  $\delta \in (0, 1)$  ויהיו  $0 < \varepsilon_1 \leq \varepsilon_2 < 1$ .

לפי ההגדרה של למידות PAC ושל ה- $\text{sample complexity}$  מתקיים לכל  $m \in \mathbb{N}$ :

$$\text{I)} \quad m \geq m_{\mathcal{H}}(\varepsilon_1, \delta) \quad \iff \quad \mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S)) \leq \varepsilon_1] \geq 1 - \delta$$

$$\text{II)} \quad m \geq m_{\mathcal{H}}(\varepsilon_2, \delta) \quad \iff \quad \mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S)) \leq \varepsilon_2] \geq 1 - \delta$$

ומכך ש- $\varepsilon_1 \leq \varepsilon_2$  נובע ש- $\{L_{\mathcal{D}}(\mathcal{A}(S)) \leq \varepsilon_1\} \subseteq \{L_{\mathcal{D}}(\mathcal{A}(S)) \leq \varepsilon_2\}$ .

לכן עבור  $m = m_{\mathcal{H}}(\varepsilon_1, \delta)$  מתקיים לפי I:

$$1 - \delta \leq \mathbb{P}_{S \sim \mathcal{D}^{m_{\mathcal{H}}(\varepsilon_1, \delta)}} [L_{\mathcal{D}}(\mathcal{A}(S)) \leq \varepsilon_1] \leq \mathbb{P}_{S \sim \mathcal{D}^{m_{\mathcal{H}}(\varepsilon_1, \delta)}} [L_{\mathcal{D}}(\mathcal{A}(S)) \leq \varepsilon_2]$$

ומכאן לפי II -

$$m_{\mathcal{H}}(\varepsilon_1, \delta) \geq m_{\mathcal{H}}(\varepsilon_2, \delta)$$

• **הוכחה:** נתונה  $\mathcal{H}$  למידה PAC. יהי  $\mathcal{A}$  אלגוריתם שלומד אותה, יהי  $\varepsilon \in (0, 1)$  ויהיו  $0 < \delta_1 \leq \delta_2 \leq 1$ .

לפי ההגדרה של למידות-PAC מתקיים לכל  $m \in \mathbb{N}$ :

$$\begin{aligned} \text{I)} \quad m \geq m_{\mathcal{H}}(\varepsilon, \delta_1) &\iff \mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S)) \leq \varepsilon] \geq 1 - \delta_1 \\ \text{II)} \quad m \geq m_{\mathcal{H}}(\varepsilon, \delta_2) &\iff \mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S)) \leq \varepsilon] \geq 1 - \delta_2 \end{aligned}$$

מכך ש- $\delta_1 \leq \delta_2$  נובע ש- $1 - \delta_1 \geq 1 - \delta_2$  ולכן עבור  $m = m_{\mathcal{H}}(\varepsilon, \delta_1)$  מתקיים לפי I:

$$\mathbb{P}_{S \sim \mathcal{D}^{m_{\mathcal{H}}(\varepsilon, \delta_1)}} [L_{\mathcal{D}}(\mathcal{A}(S)) \leq \varepsilon] \geq 1 - \delta_1 \geq 1 - \delta_2$$

ומכאן לפי II -

$$m_{\mathcal{H}}(\varepsilon, \delta_1) \geq m_{\mathcal{H}}(\varepsilon, \delta_2)$$

## שאלה 6

6. Let  $\mathcal{H}_1$  and  $\mathcal{H}_2$  be two classes for binary classification, such that  $\mathcal{H}_1 \subseteq \mathcal{H}_2$ . Show that  $VC\text{-dim}(\mathcal{H}_1) \leq VC\text{-dim}(\mathcal{H}_2)$ .

**הוכחה:** נניח בשלילה ש- $VC\text{-DIM}(\mathcal{H}_1) > VC\text{-DIM}(\mathcal{H}_2)$ .

נסמן ב- $C_1$  את הקבוצה בגודל המקסימלי המנותצת ע"י  $\mathcal{H}_1$ , וב- $C_2$  את הקבוצה בגודל המקסימלי המנותצת ע"י  $\mathcal{H}_2$ . אז:

$$|C_1| = VC\text{-DIM}(\mathcal{H}_1) > VC\text{-DIM}(\mathcal{H}_2) = |C_2|$$

נסמן  $C_1 = \{x_1, \dots, x_m\}$  אז לפי ההגדרה של צמצום של מחלקת היפתזות:

$$\begin{aligned} \mathcal{H}_1|_{C_1} &= \{(h(x_1), \dots, h(x_m)) : h \in \mathcal{H}_1\} \\ \mathcal{H}_2|_{C_1} &= \{(h(x_1), \dots, h(x_m)) : h \in \mathcal{H}_2\} \end{aligned}$$

ולכן  $\mathcal{H}_1|_{C_1} \subseteq \mathcal{H}_2|_{C_1}$ . אבל מכך ש- $\mathcal{H}_1$  מנתצת את  $C_1$  נובע ש-

$$2^{|C_1|} = |\mathcal{H}_1|_{C_1}| \leq |\mathcal{H}_2|_{C_1}| \leq 2^{|C_1|} \implies |\mathcal{H}_2|_{C_1}| = 2^{|C_1|}$$

ולכן  $\mathcal{H}_2$  מנתצת את  $C_1$ , וזו סתירה לכך ש- $C_2$  היא הקבוצה המקסימלית המנותצת על ידיה (כי ראינו ש- $|C_1| > |C_2|$ ).  
לכן  $VC\text{-DIM}(\mathcal{H}_1) \leq VC\text{-DIM}(\mathcal{H}_2)$ .

7. Prove that if  $\mathcal{H}$  has the uniform convergence property with function  $m_{\mathcal{H}}^{UC} : (0, 1)^2 \rightarrow \mathbb{N}$ , then  $\mathcal{H}$  is Agnostic-PAC learnable with sample complexity  $m_{\mathcal{H}}(\varepsilon, \delta) \leq m_{\mathcal{H}}^{UC}(\varepsilon/2, \delta)$ .

**הוכחה:** תהי  $\mathcal{H}$  מחלקת היפותזות שמקיימת את תכונת ההתכנסות בהחלט עם  $m_{\mathcal{H}}^{UC} : (0, 1)^2 \rightarrow \mathbb{N}$ , ביחס לפונקציית loss כלשהי  $\ell : \mathcal{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, \infty)$ . יהיו  $\varepsilon, \delta \in (0, 1)$ , ותהי  $\mathcal{D}$  פונקציית התפלגות על  $\mathcal{X} \times \mathcal{Y}$ . אז לפי ההגדרה לכל  $m \geq m_{\mathcal{H}}^{UC}(\frac{\varepsilon}{2}, \delta)$  מתקיים ש-

$$\mathcal{D}^m \left( \left\{ S \in (\mathcal{X} \times \mathcal{Y})^m \mid S \text{ is } \frac{\varepsilon}{2}\text{-representative} \right\} \right) \geq 1 - \delta$$

נבחר את אלגוריתם הלמידה  $\mathcal{A}_{\text{ERM}}$ , כלומר את:

$$\mathcal{A}_{\text{ERM}} : S \mapsto h \quad \text{s.t.:} \quad h \in \left\{ \underset{h \in \mathcal{H}}{\operatorname{argmin}} L_S(h) \right\}$$

קעת לפי למה 4.5.2 מספר הקורס (שראינו גם בהרצאה), אם סט דגימות  $S$  הוא  $\frac{\varepsilon}{2}$ -מייצג את  $\mathcal{D}, \mathcal{H}, \ell$  אז  $\mathcal{A}_{\text{ERM}}(S) := h_s$  מקיים ש-

$$L_{\mathcal{D}}(h_s) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \varepsilon$$

ולכן:

$$\left\{ S \in (\mathcal{X} \times \mathcal{Y})^m \mid S \text{ is } \frac{\varepsilon}{2}\text{-representative} \right\} \subseteq \left\{ S \in (\mathcal{X} \times \mathcal{Y})^m \mid L_{\mathcal{D}}(h_s) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \varepsilon \right\}$$

ולכן ממונוטוניות פונקציית ההתפלגות נובע שמתקיים לכל  $m \geq m_{\mathcal{H}}^{UC}(\frac{\varepsilon}{2}, \delta)$ :

$$1 - \delta \leq \mathcal{D}^m \left( \left\{ S \in (\mathcal{X} \times \mathcal{Y})^m \mid S \text{ is } \frac{\varepsilon}{2}\text{-representative} \right\} \right) \leq \mathcal{D}^m \left( \left\{ S \in (\mathcal{X} \times \mathcal{Y})^m \mid L_{\mathcal{D}}(h_s) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \varepsilon \right\} \right)$$

ובפרט זה מתקיים עבור  $m = m_{\mathcal{H}}^{UC}(\frac{\varepsilon}{2}, \delta)$ .

כלומר הראינו שלכל  $\varepsilon, \delta \in (0, 1)$ , ולכל  $\mathcal{D}$  פונקציית התפלגות על  $\mathcal{X} \times \mathcal{Y}$ , קיים אלגוריתם למידה  $\mathcal{A}_{\text{ERM}}$  - המקיים לכל  $m \geq m_{\mathcal{H}}^{UC}(\frac{\varepsilon}{2}, \delta)$ :

$$\mathcal{D}^m \left( \left\{ S \in (\mathcal{X} \times \mathcal{Y})^m \mid L_{\mathcal{D}}(\mathcal{A}_{\text{ERM}}(S)) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \varepsilon \right\} \right) \geq 1 - \delta$$

■

ולכן לפי ההגדרה  $\mathcal{H}$  היא Agnostic-PAC למידה עם sample complexity  $m_{\mathcal{H}}(\varepsilon, \delta) \leq m_{\mathcal{H}}^{UC}(\frac{\varepsilon}{2}, \delta)$ . כנדרש.

8. Let  $\mathcal{H}$  be a hypothesis class over a domain  $\mathcal{Z} = \mathcal{X} \times \{\pm 1\}$ , and consider the 0-1 loss function. Assume that there exists a function  $m_{\mathcal{H}}$ , for which it holds that for every distribution  $\mathcal{D}$  over  $\mathcal{Z}$  there is an algorithm  $\mathcal{A}$  with the following property: when running  $\mathcal{A}$  on  $m \geq m_{\mathcal{H}}$  i.i.d. examples drawn from  $\mathcal{D}$ , it is guaranteed to return, with probability at least  $1 - \delta$ , a hypothesis  $h_S : \mathcal{X} \rightarrow \{\pm 1\}$  with  $L_{\mathcal{D}}(h_S) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \varepsilon$ . Is  $\mathcal{H}$  agnostic PAC learnable? Prove or show a counter example.

מחלקת היפותזות  $\mathcal{H}$  שמקיימת את התנאי הזה היא לא בהכרח למידה-agnostic-PAC עם פונקציית ה-loss  $\ell_{0-1}$ .



הסיבה לכך היא שהתנאי הזה מאפשר לאלגוריתם הלמידה  $\mathcal{A}$  "להכיר" את ההתפלגות  $\mathcal{D}$  על  $\mathcal{Z}$ , ובפרט מאפשר לו לחשב את  $L_{\mathcal{D}}(h)$  לכל  $h \in \mathcal{H}$ , דבר שאינו אפשרי ב-framework של agnostic PAC. דוגמא נגדית:

נגדיר  $\mathcal{X} = \mathbb{R}$ , וניקח את  $\mathcal{H}$  להיות מחלקת כל הפונקציות מ- $\mathcal{X}$  ל- $\{\pm 1\}$  (Unstructured Hypotheses Class), כלומר  $\mathcal{H} = \mathbb{R}^{\{\pm 1\}}$ . מקיימת את התנאי הנתון:

נגדיר  $m_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathbb{N}$  ע"י  $m_{\mathcal{H}}(\varepsilon, \delta) = 1$  לכל  $\varepsilon, \delta \in (0, 1)$ .

לכל ההתפלגות  $\mathcal{D}$  על  $\mathcal{Z}$  ניקח אלגוריתם  $\mathcal{A} : \mathcal{Z}^m \rightarrow \mathcal{H}$  המוגדר לכל  $S = \{(x_i, y_i)\}_{i=1}^m \subseteq \mathcal{Z}^m$  ע"י:

$$\mathcal{A}(S) = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \{L_{\mathcal{D}}(h)\}$$

נקבל שלכל  $\varepsilon, \delta \in (0, 1)$ , כאשר מפעילים את  $\mathcal{A}$  על  $S \in \mathcal{Z}^m$  עבור  $m \geq m_{\mathcal{H}}(\varepsilon, \delta)$  הוא מחזיר בהסתברות  $1 - \delta \leq 1$  היפותזה  $h_S : \mathcal{X} \rightarrow \{\pm 1\}$  כך ש-

$$L_{\mathcal{D}}(h_S) \stackrel{\text{מוגדרת } \mathcal{A}}{=} \min_{h \in \mathcal{H}} \{L_{\mathcal{D}}(h)\} \leq \min_{h \in \mathcal{H}} \{L_{\mathcal{D}}(h)\} + \varepsilon$$

מצד שני, ראינו בהרצאה ש- $\mathcal{H}$  היא לא למידה-PAC.

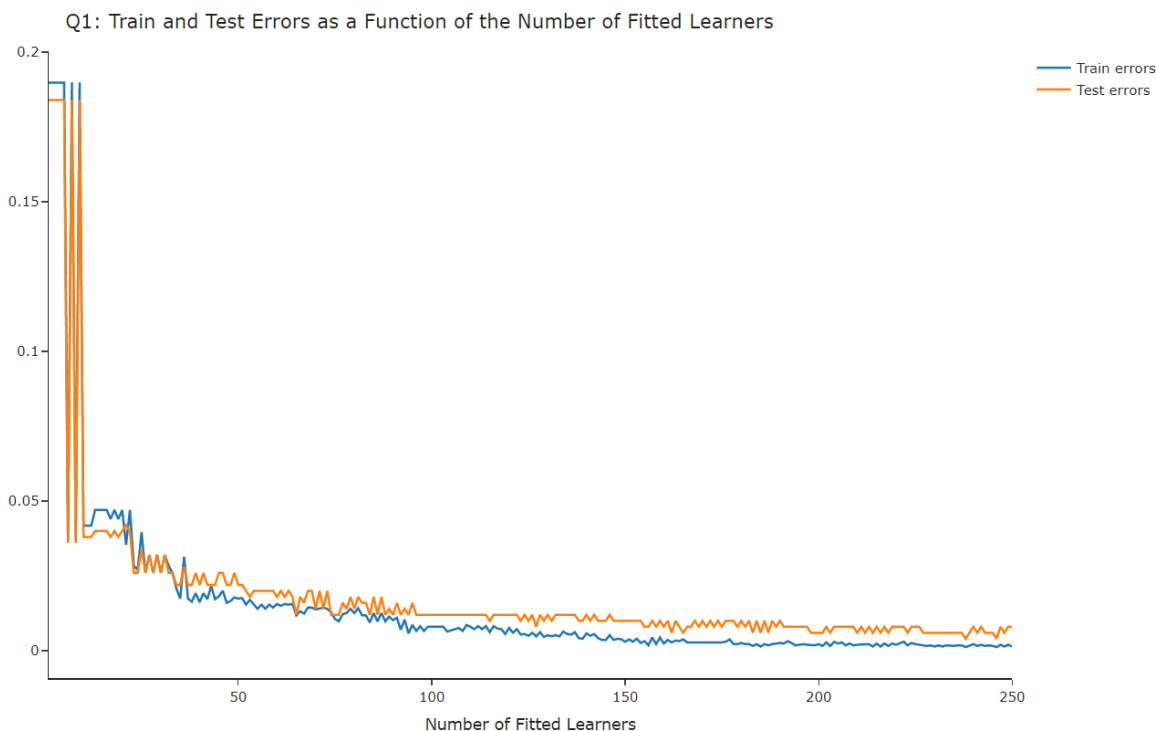
ראינו בנוסף שאם מחלקת היפותזות היא למידה-agnostic-PAC עם פונקציית ה-loss  $\ell_{0-1}$  אז היא גם למידה PAC.

לכן מכך ש- $\mathcal{H}$  לא למידה PAC נובע שהיא **לא** למידה-agnostic-PAC עם  $\ell_{0-1}$ .

## חלק פרקטי

### שאלה 1

גרף המציג את השגיאה (misclassification error) של Adaboost על דגימות האימון (בכחול) ועל דגימות המבחן (בכתום) כפונקציה של מספר המודלים החלשים שנעשה בהם שימוש באלגוריתם:

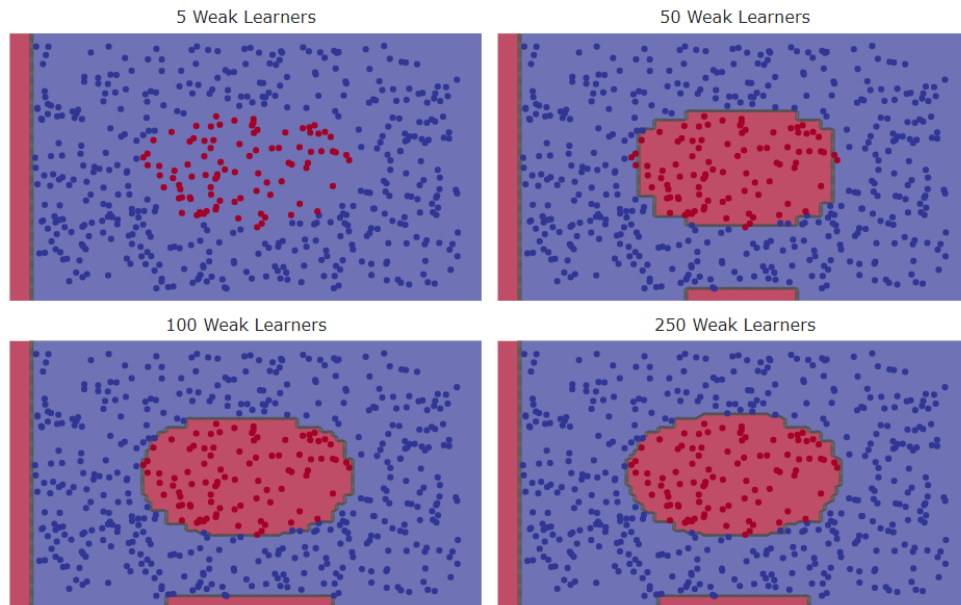


רואים מהגרף שככל שמספר המודלים החלשים שמשתתפים באלגוריתם **גדול** יותר - השגיאה נהיית **קטנה** יותר. בנוסף, רואים שהשגיאה של המודל על דגימות האימון קטנה יותר מהשגיאה שלו על דגימות המבחן. עם זאת, הפער בשגיאה לא גדול כ"כ ולכן אפשר אולי להסיק שהאלגוריתם מכליל לא רע ולא מתאים את עצמו מדי לדאטא האימון (כלומר לא סובל מ-over fitting).

## שאלה 2

גרף המציג את ה־decision boundary של המודל עבור 5, 50, 100, 250 לומדים חלשים שמשתתפים בו, כאשר הנקודות הן דגימות המבחן:

Q2: Decision Boundary Obtained by Using the Ensemble Up to Iteration [5, 50, 100, 250]

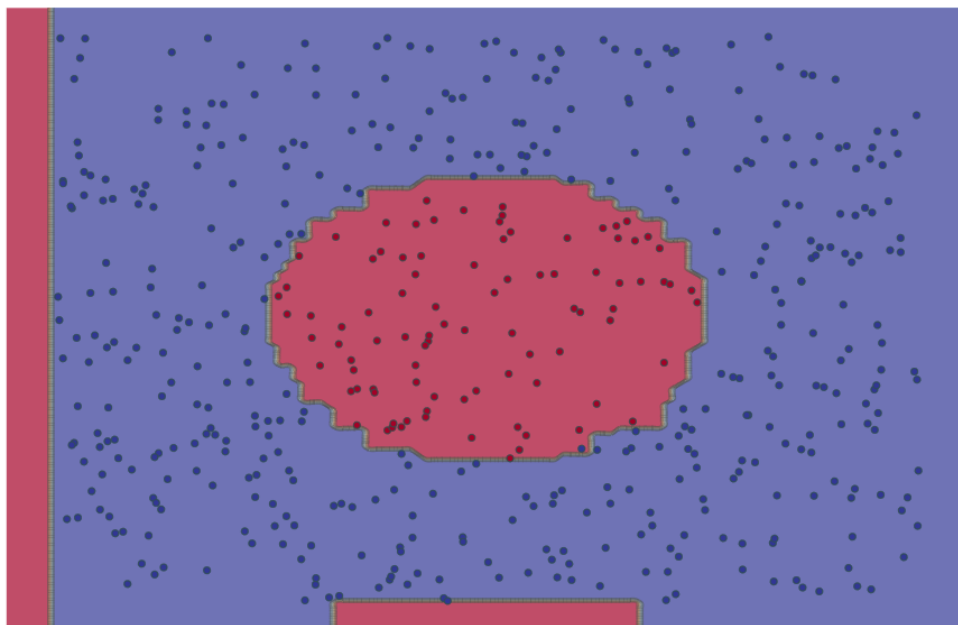


רואים מהגרף שככל שיותר לומדים חלשים משתתפים במודל, הוא נהיה אקספרסיבי יותר ומצליח לסווג טוב יותר את דגימות המבחן.

## שאלה 3

גרף המציג את ה־decision boundary של המודל עבור האנסמבל בגודל האופטימלי (מתוך 250 האפשרויות):

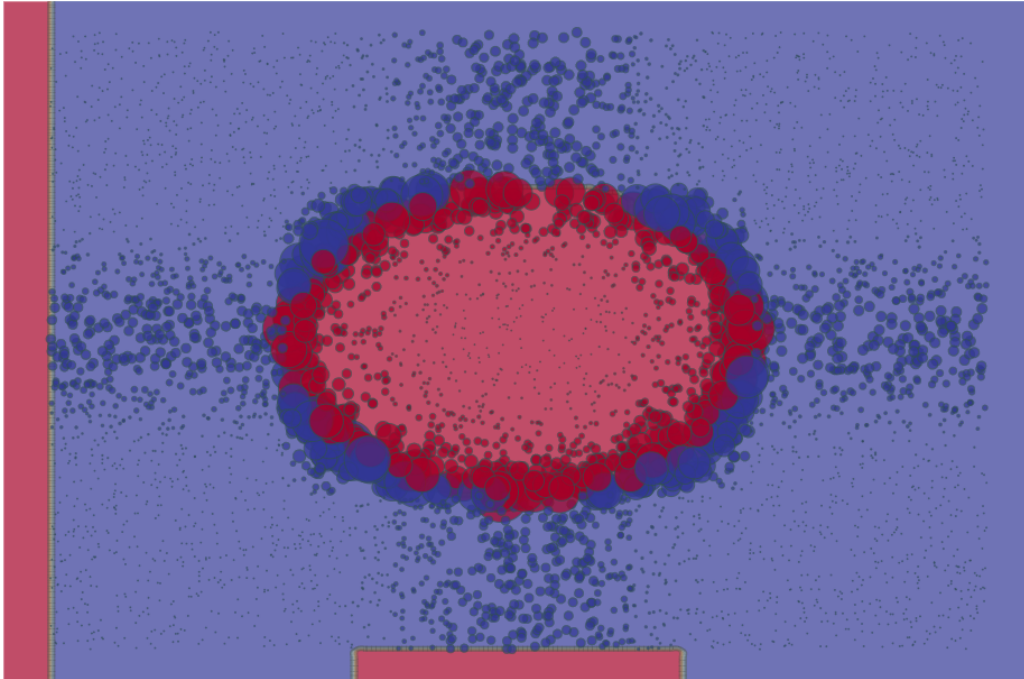
Q3: Decision Boundary of the Best Performing Ensemble, With Size = 238 and Accuracy = 0.996



#### שאלה 4

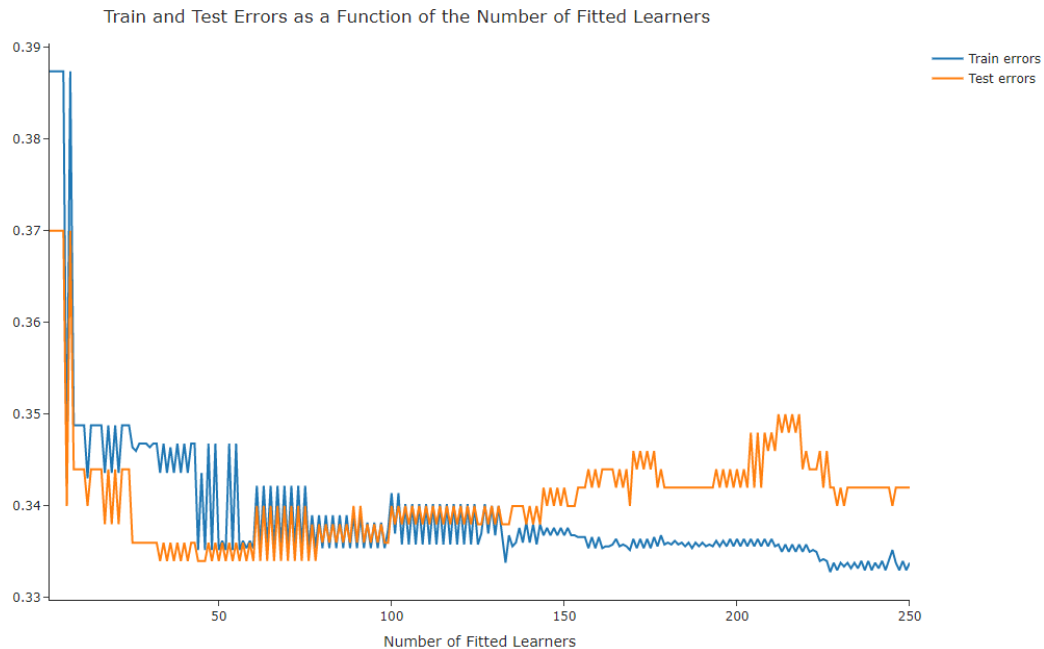
גרף המציב את ה- $\text{decision boundary}$  של המודל המלא על דאטא **האימון**, כאשר גודל הנקודות פורפוציונלי להתפלגות שלהן בסוף ריצת ה- $\text{adaboost}$ :

Q4: Decision Boundary of full Ensemble With Point Size Proportional to its Weight



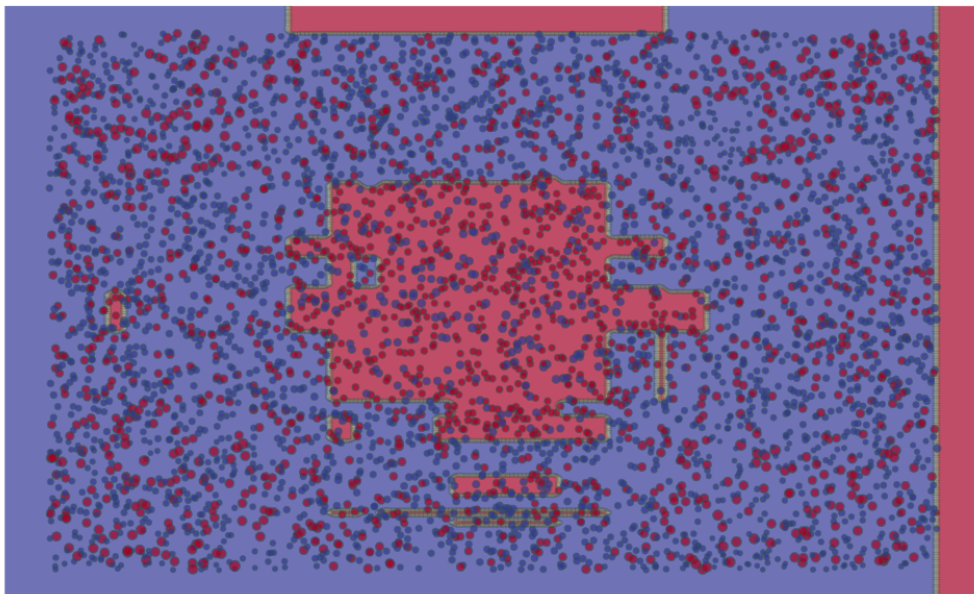
ניתן לראות שהנקודות עם המשקל הגדול ביותר בסוף ריצת האלגוריתם הן הנקודות שנמצאות על "הגבול" בין המחלקות. המשמעות של זה היא שמרבית הטעויות שהלומדים החלשים במודל עשו במהלך ה- $\text{fit}$  נעשו על הנקודות הללו. אפשר להסיק מכך (וזו מסקנה גיונית לחלוטין בלי קשר לגרף) שהנקודות שהכי "מאתגרות" את המודל הן הנקודות שנמצאות על הגבול בין המחלקות, ואילו הנקודות שהכי "קל" למודל לסווג הן הנקודות שנמצאות בסביבה הומוגנית יותר, כלומר רחוקות מנקודות עם לייבל ששונה מהלייבל שלהן.

הגרפים המתאימים לשאלות 1 ו-4 עבור דאטא שנוצר עם רעש ברמה 0.4:



באופן דומה למקרה ללא הרעש, ככל שמספר הלומדים החלשים משתתפים במודל - השגיאה על סט האימון קטנה. אבל במקרה הזה, בשלב מסויים ניתן לראות בגרף שהשגיאה על סט המבחן מתחילה לגדול. ניתן להסביר זאת בכך שהמודל מתחיל לעשות overfitting לדאטא וללמוד את הרעש. ניתן לראות שהפער בין השגיאה של המודל על סט האימון והשגיאה שלו על סט המבחן גדל ככל שהוא עושה שימוש ביותר לומדים חלשים. במונחי bias-variance אפשר לומר שככל שהמורכבות של המודל גדולה יותר, כלומר ככל שהוא עושה שימוש ב-ensemble גדול יותר, ה-variance שלו גדל בעוד שה-bias שלו קטן.

Decision Boundary of full Ensemble With Point Size Proportional to its Weight



ניתן לראות שכאשר הדאטא נדגם עם רעש (ברמה דיי גבוהה) המודל מתקשה יותר להפריד אותו.  
מכיוון שיש נקודות רבות שהלומדים החלשים טועים עליהן, אז הגודל של רוב הנקודות בגרף השני הוא יחסית דומה.