

מערכות לומדות תרגיל 4

שיר שבח 322407701

12 במאי 2022

חלק תיאורטי

PAC Learnability 2.1

1. עבור אלגוריתם למידה \mathcal{A} כלשהו, התפלגות \mathcal{D} מעל \mathcal{X} ופונקציית $0-1$ loss (misclassification), הוכיחו שהבאים שווים:

$$\forall \varepsilon, \delta > 0 \quad \exists m(\varepsilon, \delta) \text{ s.t. } \forall m \geq m(\varepsilon, \delta) \quad \mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S)) \leq \varepsilon] \geq 1 - \delta(a)$$

$$\lim_{m \rightarrow \infty} \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S))] = 0 \quad (b)$$

יהי $\varepsilon > 0$. נרצה להוכיח כי $\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S))] < \varepsilon$.
עפ"י הגדרת תוחלת עבור מ"מ רציפים: $(a) \Rightarrow (b)$

$$\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S))] = \int_0^1 t \cdot \mathbb{P}_{S \sim \mathcal{D}^m} (L_{\mathcal{D}}(\mathcal{A}(S)) = t) dt$$

נבחר כי מתקיים משפט (a) עבור $\frac{\varepsilon}{2}$ ו- $\delta = \frac{\varepsilon}{2}$. לכן עבור $m_0 > m(\frac{\varepsilon}{2}, \frac{\varepsilon}{2})$ מתקיים:

$$\begin{aligned}
\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S))] &= \int_0^1 t \cdot \mathbb{P}_{S \sim \mathcal{D}^m} (L_{\mathcal{D}}(A(S)) = t) dt \\
&= \int_0^{\frac{\varepsilon}{2}} t \cdot \mathbb{P}_{S \sim \mathcal{D}^m} (L_{\mathcal{D}}(A(S)) = t) dt + \int_{\frac{\varepsilon}{2}}^1 t \cdot \mathbb{P}_{S \sim \mathcal{D}^m} (L_{\mathcal{D}}(A(S)) = t) dt \\
&\leq \int_0^{\frac{\varepsilon}{2}} \frac{\varepsilon}{2} \cdot \mathbb{P}_{S \sim \mathcal{D}^m} (L_{\mathcal{D}}(A(S)) = t) dt + \int_{\frac{\varepsilon}{2}}^1 1 \cdot \mathbb{P}_{S \sim \mathcal{D}^m} (L_{\mathcal{D}}(A(S)) = t) dt \\
&= \frac{\varepsilon}{2} \cdot \int_0^{\frac{\varepsilon}{2}} \mathbb{P}_{S \sim \mathcal{D}^m} (L_{\mathcal{D}}(A(S)) = t) dt + \int_{\frac{\varepsilon}{2}}^1 \mathbb{P}_{S \sim \mathcal{D}^m} (L_{\mathcal{D}}(A(S)) = t) dt \\
&= \frac{\varepsilon}{2} \cdot \mathbb{P}_{S \sim \mathcal{D}^m} \left(L_{\mathcal{D}}(A(S)) \leq \frac{\varepsilon}{2} \right) + \mathbb{P}_{S \sim \mathcal{D}^m} \left(L_{\mathcal{D}}(A(S)) > \frac{\varepsilon}{2} \right) \\
&= \frac{\varepsilon}{2} \cdot \left(1 - \mathbb{P}_{S \sim \mathcal{D}^m} \left(L_{\mathcal{D}}(A(S)) > \frac{\varepsilon}{2} \right) \right) + \mathbb{P}_{S \sim \mathcal{D}^m} \left(L_{\mathcal{D}}(A(S)) > \frac{\varepsilon}{2} \right) \\
&= \frac{\varepsilon}{2} \cdot 1 - \frac{\varepsilon}{2} \cdot \mathbb{P}_{S \sim \mathcal{D}^m} \left(L_{\mathcal{D}}(A(S)) > \frac{\varepsilon}{2} \right) + \mathbb{P}_{S \sim \mathcal{D}^m} \left(L_{\mathcal{D}}(A(S)) > \frac{\varepsilon}{2} \right) \\
&= \frac{\varepsilon}{2} \cdot 1 + \left(1 - \frac{\varepsilon}{2} \right) \cdot \mathbb{P}_{S \sim \mathcal{D}^m} \left(L_{\mathcal{D}}(A(S)) > \frac{\varepsilon}{2} \right) \\
&< \frac{\varepsilon}{2} \cdot 1 + \left(1 - \frac{\varepsilon}{2} \right) \cdot \frac{\varepsilon}{2} \\
&= \varepsilon - \frac{\varepsilon^2}{4} < \varepsilon
\end{aligned}$$

$\frac{\cdot (b) \Rightarrow (a)}{\text{יהי } \varepsilon, \delta > 0}$
 עפ"י אי שיויון מרקוב:

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

אם נסמן $a = \varepsilon$ ו- $\delta = \varepsilon \cdot \delta$, אזי יתקיים:

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

וקיבלנו

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

שזה שקול ל-

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

2. יהי $\mathcal{X} := \mathbb{R}^2, \mathcal{Y} := \{0, 1\}$ ויהי \mathcal{H} מחלקה של מעגלים קונצנטריים במישור:

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

הוכיחו כי \mathcal{H} היא למידה-PAC וה-sample complexity חסום ע"

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

בעת הוכחה, אל תשתמשו בטיעון VC-Dimension. במקום זאת הוכיחו את הטענה ישירות מהגדרת למידות PAC על ידי הצגת אלגוריתם ספציפי וניתוח מורכבות המדגם שלו.

האלגוריתם A שניקח הוא אלגוריתם שחזה את r באופן הבא:

$$\hat{r} = \frac{1}{2} \cdot \max_{x_i, x_j \in \mathcal{X}} \{d(x_i, x_j)\} \text{ s.t. } y_i = y_j = 1$$

נסמן את איזור החיזוי שלנו $(h_{\hat{r}})$ ב- \hat{I} , ואת האיזור האמיתי (h_r) ב- I . איזור הטעות הוא $I \setminus \hat{I}$. לכן

$$L_D(h_s) = \mathbb{P}(x \in I \setminus \hat{I})$$

נרצה לחשב את $\mathbb{P}\{L_D(h_s) \leq \varepsilon\}$. ראשית נחשב את $\mathbb{P}\{L_D(h_s) > \varepsilon\}$.

$$\mathbb{P}\{L_D(h_s) > \varepsilon\} = \mathbb{P}\left(\bigwedge_{i=1}^m x_i \notin I \setminus \hat{I}\right) \stackrel{iid}{=} \prod_{i=1}^m \mathbb{P}(x_i \notin I \setminus \hat{I}) \leq (1 - \varepsilon)^m \leq e^{-\varepsilon m}$$

לכן:

$$\mathbb{P}\{L_D(h_s) \leq \varepsilon\} = 1 - e^{-\varepsilon m}$$

אנו נרצה ש-

$$\begin{aligned}
1 - e^{-\varepsilon m} &\geq 1 - \delta \\
e^{\varepsilon m} &\geq \frac{1}{\delta} \\
\varepsilon m &\geq \log\left(\frac{1}{\delta}\right) \\
m &\geq \frac{\log\left(\frac{1}{\delta}\right)}{\varepsilon}
\end{aligned}$$

VC-Dimension 2.2

3. יהי $\mathcal{X} = \{0, 1\}^n$ ו- $\mathcal{Y} = \{0, 1\}$ לכל $I \subseteq [n]$ הגדירו את הפונקציה הזוגית:

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

מה מימד VC של המחלקה $\mathcal{H}_{parity} = \{h_I \mid I \subseteq [n]\}$ הוכיחו את תשובתכם.

נשים לב כי לכל מימד שקטן מ- n , ניתן לנתן את h . מספיק להראות עבור דוגמה אחת, ניקח $\{0\}^1, \dots, \{0\}^n$, אזי לכל $I \subseteq [n]$:

$$h_I(\{0\}^i) = 0$$

בנוסף, משום שגודל ההיפוטזה הוא $|\mathcal{H}| = 2^n$, זאת קבוצה סופית. ולכן $\text{VCdim}(\mathcal{H}) \leq \log_2(|\mathcal{H}|)$ ולכן $\text{VCdim}(\mathcal{H}) \leq n$. $\log_2(2^n) = n$ מה שגורר כי $\text{VCdim}(\mathcal{H}) = n$.

4. בהינתן מספר שלם k , יהי $[a_i, b_i]_{i=1}^k$ להיות קבוצה כלשהי של קטעים על \mathbb{R} ומוגדר האיחוד $A = \bigcup_{i=1}^k [a_i, b_i]$. מחלקת ההיפוטזות $\mathcal{H}_{k-intervals}$ כוללת את הפונקציות: $h_A(x) := \mathbb{1}_{[x \in A]}$, לכל בחירות של k מקטעים. מצאו את מימד VC של $\mathcal{H}_{k-intervals}$ והוכיחו את תשובתכם. הראו כי אם אנו נותנים ל- A להיות כל איחוד סופי של מקטעים (כלומר k לא מוגבל), אז המחלקה המתקבלת $\mathcal{H}_{intervals}$ יש מימד $\text{VC} = \infty$.

נרצה להוכיח כי מימד VC עבור מחלקת היפוטזות $\mathcal{H}_{k-intervals}$ הוא $2k$.
 לכן נרצה להוכיח כי לכל $d < 2k$, קיים צמצום C , $|C| = d$, כך ש $\mathcal{H}_{k-intervals}$ מנתצת את C .
 עבור דגימה $\{x_1, x_2, \dots, x_{2k}\}$, נעטוף במקטע אחד אינדקסים שמתוייגים ברצף כ-1. לדוגמה, עבור $k = 3$, $|X| = 6$, $X = \{0, 1, 1, 0, 0, 1\}$ אז $A = [1, 2] \cup [5]$. נשים לב כי גודל הקבוצות באיחוד הוא לכל היותר k .
 משום שבמחלקת ההיפוטזות שלנו $\mathcal{H}_{k-intervals}$ זהו איחוד של k מקטעים, קיבלנו כי $C = h_A(x) \in \mathcal{H}_{k-intervals}$, כלומר $\mathcal{H}_{k-intervals}$ מנתצת את C .

נרצה להוכיח כי לכל $|C| = 2k + 1$, $\mathcal{H}_{k-intervals}$ לא מנתצת את C .
 עבור דגימה $\{x_1, \dots, x_{2k}, x_{2k+1}\}$, אם מספר הנקודות המתוייגות כ-1 גדול מ- k , אזי משובך היונים קיים מקטע שבו יש 2 נקודות. אם ביניהם נשים נקודה המתוייגת ל-0, הם לא יוכלו להיות באותו מקטע, ובכך אין פונקציה ב- $\mathcal{H}_{k-intervals}$ שיכולה לתייג נכון דגימה זו. ולכן $\mathcal{H}_{k-intervals}$ לא מנתצת את C .
 עבור המחלקת $\mathcal{H}_{intervals}$, נשים לב כי לכל k , קבוצה בגודל $2k$ מנותצת ע"י פונקציה ב- $\mathcal{H}_{k-intervals}$. קל וחומר עבור מחלקת היפוטזות $\mathcal{H}_{intervals}$. ולכן $\mathcal{H}_{intervals}$ מנתצת קבוצה מכל גודל סופי. לכן מימד VC שלה הוא ∞ .

Monotonicity 2.3

5. יהי \mathcal{H} מחלקת היפוטזות עבור סיווג בינארי. נניח כי \mathcal{H} למידה PAC וסיבוכיות המודל ניתן ע"י $m_{\mathcal{H}}(\cdot, \cdot)$. הראו כי $m_{\mathcal{H}}$ מונוטוני יורד בכל אחד מהפרמטרים שלו. זאת אומרת:

- הראו כי בהינתן $\delta \in (0, 1)$, ובהינתן $0 < \varepsilon_1 \leq \varepsilon_2 < 1$, חייב כי $m_{\mathcal{H}}(\varepsilon_1, \delta) \geq m_{\mathcal{H}}(\varepsilon_2, \delta)$.
 - בדומה, הראו כי בהינתן $\varepsilon \in (0, 1)$ ובהינתן $0 < \delta_1 \leq \delta_2 < 1$, חייב כי $m_{\mathcal{H}}(\varepsilon, \delta_1) \geq m_{\mathcal{H}}(\varepsilon, \delta_2)$.
- עבור \mathcal{H} למידה PAC אנו יודעים כי מתקיים:

$$m_{\mathcal{H}}(\varepsilon, \delta) \sim \frac{VCdim(\mathcal{H}) + \log\left(\frac{1}{\delta}\right)}{\varepsilon}$$

בהינתן $\delta \in (0, 1)$ ו $0 < \varepsilon_1 \leq \varepsilon_2 < 1$, מתקיים:

$$m_{\mathcal{H}}(\varepsilon_1, \delta) \sim \frac{VCdim(\mathcal{H}) + \log\left(\frac{1}{\delta}\right)}{\varepsilon_1} \geq \frac{VCdim(\mathcal{H}) + \log\left(\frac{1}{\delta}\right)}{\varepsilon_2} \sim m_{\mathcal{H}}(\varepsilon_2, \delta)$$

ובהינתן $\varepsilon \in (0, 1)$ ו $0 < \delta_1 \leq \delta_2 < 1$

$$m_{\mathcal{H}}(\varepsilon, \delta_1) \sim \frac{VCdim(\mathcal{H}) + \log\left(\frac{1}{\delta_1}\right)}{\varepsilon} \geq \frac{VCdim(\mathcal{H}) + \log\left(\frac{1}{\delta_2}\right)}{\varepsilon} \sim m_{\mathcal{H}}(\varepsilon, \delta_2)$$

כנדרש.

6. יהי \mathcal{H}_1 ו- \mathcal{H}_2 שתי מחלקות המסווגות בינארי, כך ש- $\mathcal{H}_1 \subseteq \mathcal{H}_2$. הראו כי $VCdim(\mathcal{H}_1) \leq VCdim(\mathcal{H}_2)$.

נניח בשלילה כי $\mathcal{H}_1 \subseteq \mathcal{H}_2$ אולם $VCdim(\mathcal{H}_1) > VCdim(\mathcal{H}_2)$. אזי קיים צמצום C כך ש \mathcal{H}_1 מנתצת את C אולם \mathcal{H}_2 לא. סתירה לכך שלכל $h \in \mathcal{H}_1$ חייב כי $h \in \mathcal{H}_2$.

Agnostic-PAC 2.4

7. הוכיחו כי אם \mathcal{H} יש תכונת התכנסות אחידה עם פונקציה $m_{\mathcal{H}}^{UC} : (0, 1)^2 \rightarrow \mathbb{N}$, אז \mathcal{H} היא למידה Agnostic-PAC עם סיבוכיות מודל $m_{\mathcal{H}}(\varepsilon, \delta) \leq m_{\mathcal{H}}^{UC}\left(\frac{\varepsilon}{2}, \delta\right)$.

משום של- \mathcal{H} יש תכונת UC עם פונקציית $m_{\mathcal{H}}^{UC}$, יהי S סט דגימות $\frac{\varepsilon}{2}$ -מייצג (מגודל $m > m_{\mathcal{H}}^{UC}\left(\frac{\varepsilon}{2}, \delta\right)$) עבור ההיפוטזה \mathcal{H} , התפלגות \mathcal{D} , ופונקציית ℓ loss נסמן $h_S = \text{ERM}_{\mathcal{H}}(S)$. לכן, בהסתברות של לפחות $1 - \delta$ מתקיים:

$$L_D(h_S) \leq L_D(h) + \frac{\varepsilon}{2}$$

לכן, לכל $h \in \mathcal{H}$

$$L_D(h_S) \stackrel{\clubsuit}{\leq} L_S(h_S) + \frac{\varepsilon}{2} \stackrel{\spadesuit}{\leq} \min_{h \in \mathcal{H}} L_S(h) + \frac{\varepsilon}{2} \stackrel{\clubsuit}{\leq} \min_{h \in \mathcal{H}} L_D(h) + \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \min_{h \in \mathcal{H}} L_D(h) + \varepsilon$$

\clubsuit - S הוא ε -מייצג לכן לפי ההגדרה:

$$|L_S(h) - L_D(h)| \leq \frac{\varepsilon}{2} \Leftrightarrow -\frac{\varepsilon}{2} \leq L_S(h) - L_D(h) \leq \frac{\varepsilon}{2}$$

\spadesuit - $h_S = \text{ERM}_{\mathcal{H}}(S)$
לכן קיבלנו כי \mathcal{H} היא Agnostic-PAC עם $m^{\text{UC}}(\frac{\varepsilon}{2}, \delta)$

8. יהי \mathcal{H} מחלקת היפוטזות מעל תחום $\mathcal{Z} = \mathcal{X} \times \{\pm 1\}$, ועבור פונקציית $0-1$ loss. הניחו כי קיימת פונקציה $m_{\mathcal{H}}$, שעבורה היא מחזיקה שלכל התפלגות \mathcal{D} מעל \mathcal{Z} קיים אלגוריתם \mathcal{A} עם התכונות הבאות: כאשר מריצים את \mathcal{A} על $m > m_{\mathcal{H}}$ דוגמאות iid המתפלגות מ- \mathcal{D} , מובטח שיחזור, עם הסתברות של לפחות $1-\delta$, היפוטזה $h_S : \mathcal{X} \rightarrow \{\pm 1\}$ עם $L_{\mathcal{D}}(h_S) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \varepsilon$. האם \mathcal{H} למידת Agnostic-PAC? הוכיחו או הראו דוגמה מספרית.

נוכיח כי H אינה למידת agnostic PAC.

ההגדרה של למידות agnostic PAC:

יהיו $\varepsilon, \delta \in (0, 1)$, נאמר כי אלגוריתם למידה A הוא למיד agnostic PAC עם בטיחות δ ודיוק ε ביחס לפונקציית ℓ loss, מחלקת היפוטזות \mathcal{H} והתפלגות D :

$$D^m \left(\left\{ S \mid L_D(h_S) \leq \min_{h' \in \mathcal{H}} L_D(h') + \varepsilon \right\} \right) \geq 1 - \delta$$

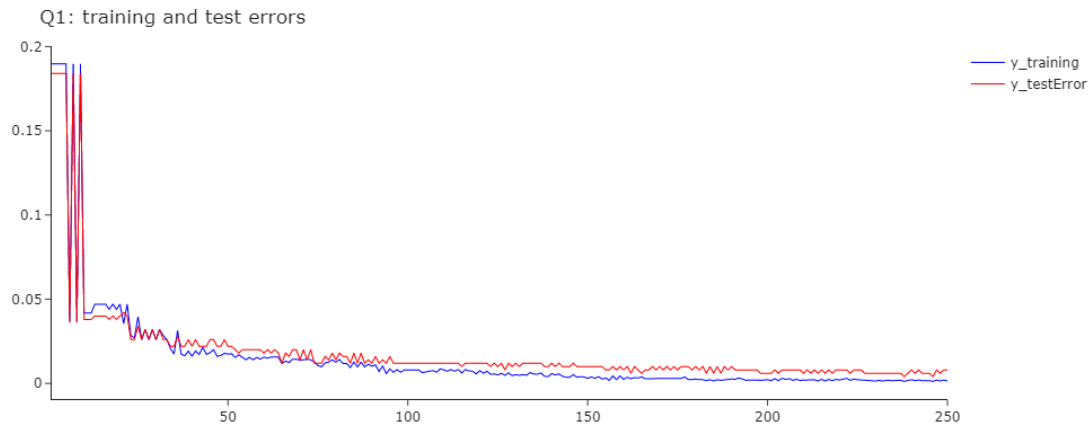
נשים לב כי לפי הגדרה אלגוריתם A אינו תלוי ב- D . לעומת זאת, A הנתון בשאלה תלוי ב- D .

דוגמה מספרית: ניקח H מחלקת היפוטזות שאינה למידה PAC , לדוגמה H היא כל הפונקציות הבינאריות מעל הממשיים. אם H למידה agnostic PAC אזי קיים אלגוריתם A עם סיבוכיות מקום $m_H(\varepsilon, \delta)$, כך שלכל התפלגות D $\mathbb{P}(A(S) \leq \min_{h \in H} L_D(h) + \varepsilon) \geq 1 - \delta$.

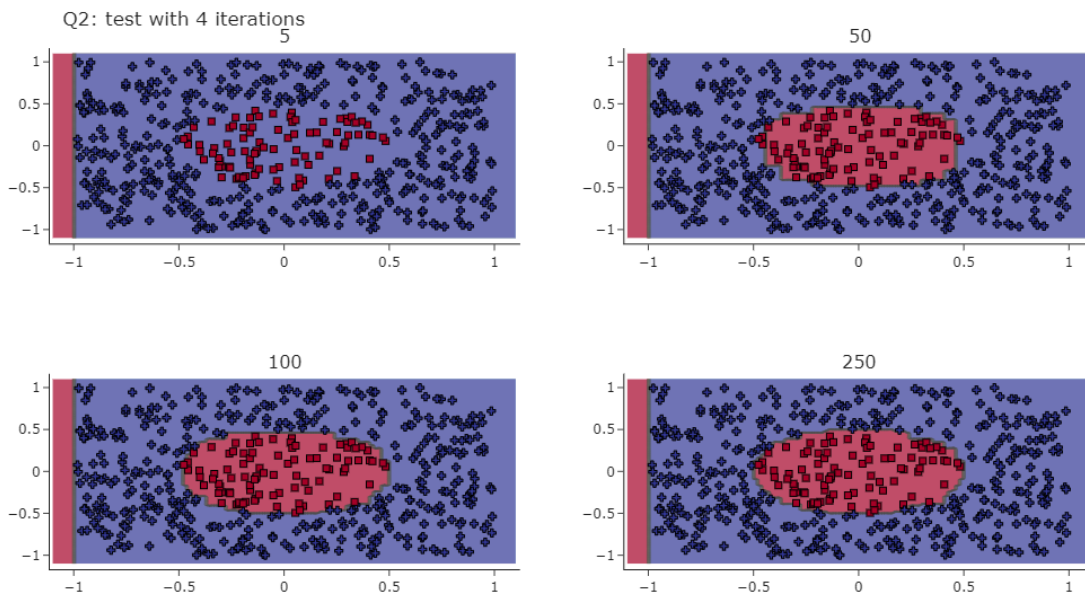
למדנו בתרגול כי H אינה למידה PAC אזי היא גם לא למידה agnostic PAC, סתירה.

מערכות לומדות חלק פרקטי

1. ניתן לראות לפי הגרף כי הטעות של train אכן דומה לtest בכל איטרציה ואכן האלגוריתם עובד נכון.

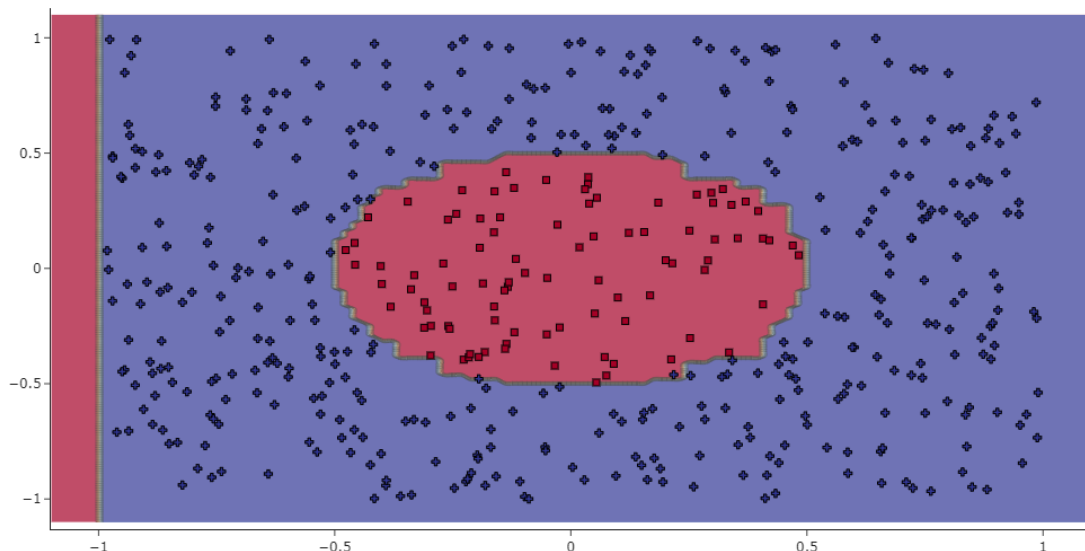


2. ניתן לראות לפי הגרפים שיצאו כי אכן ככל שמתקדמים באיטרציות החיזוי טוב יותר, אולם בין 50-100-250 השינויים מינוריים.



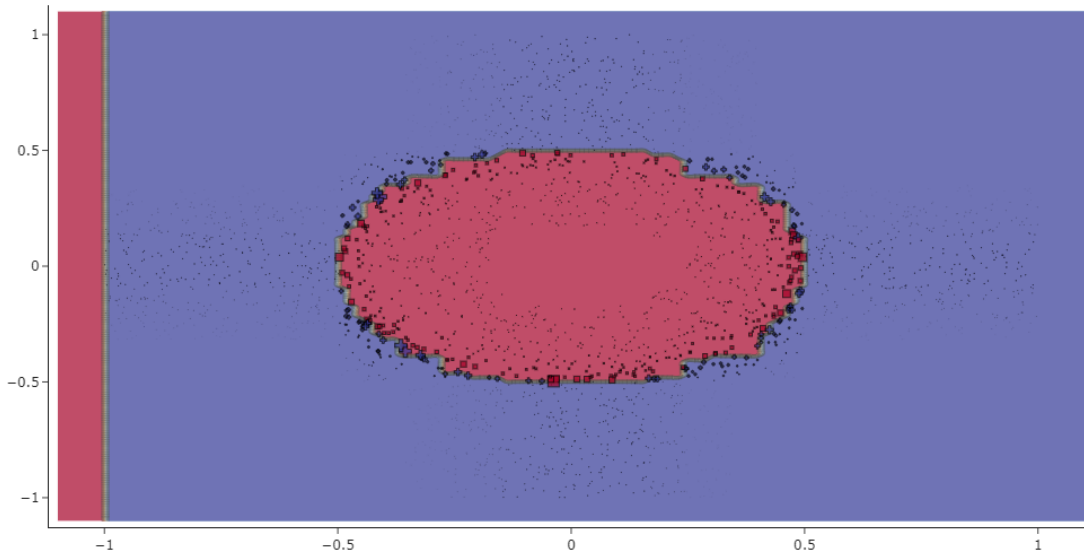
3.

Q3: Decision Surface of Ensemble with lowest Error is the size 238 with the accuracy is 0.996



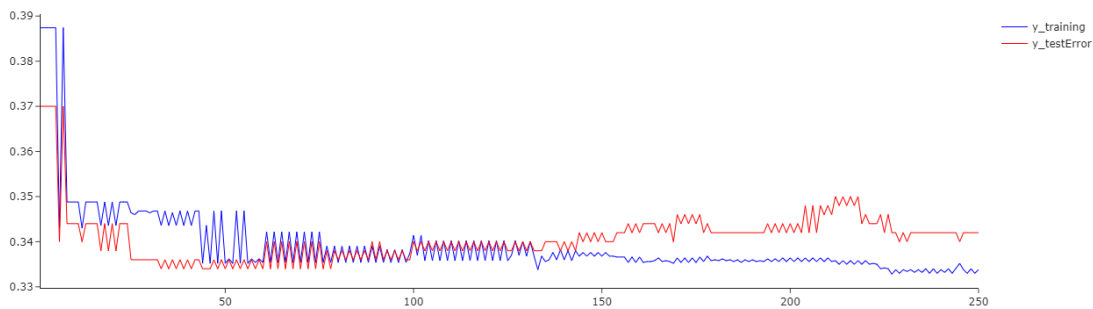
4. ניתן לראות לפי הגרף שאכן הנקודות שיותר קשה לחזות אותן, משום שהן קרובות לthreshold, אכן המשקל שלהן גבוה יותר כי במהלך האלגוריתם הן נהפכו להיות קריטיות יותר. ואכן הנקודות הרחוקות מכלל ההחלטה קטנות יותר, ז"א המשקל שלהן קטן יותר, כי הן קלות יותר לחיזוי לפי האלגוריתם.

Q4: Final Decision Boundary with proportional size training data point



5. ניתן לראות בגרף הראשון כי בהתחלה לtrain יש טעות גדולה יותר מהtest, כי אכן הוא מתחשב בנקודות שהן רעש, אולם בהמשך, כשיש הרבה איטרציות לאגוריתם adaboost, הוא כבר מתחיל לא להתייחס לנקודות רעש ויש טעות קטנה הרבה יותר.

Graph Q1 training and test errors, noise=0.4



כאן ניתן לראות כי משום שיש רעש, אכן החיזוי קשה יותר.

Q4: Final Decision Boundary with proportional size training data point, noise=0.4

