

## עבודת הגשה מספר 4

קורס: מבוא למערכות לומדות (67577)

מגיש: דור מסיקה, ת.ז: 318391877

11 במאי 2022

### החלק התיאורטי

#### PAC Learnability

1.

עבור אלגוריתם למידה  $A$ , התפלגות  $D$  מעל  $\mathcal{X}$  ופונקציית  $loss_{0-1}$  (misclassification), נוכיח שהבאים הינם שקולים:  
(א)

$$\forall \varepsilon, \delta > 0, \exists m(\varepsilon, \delta) \text{ s.t. } \forall m \geq m(\varepsilon, \delta) \quad P_{S \sim D^m} [L_D(A(S)) \leq \varepsilon] \geq 1 - \delta$$

(ב)

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

**פתרון:** נוכיח את הטענה באמצעות הכלה דו כיוונית.

$(a \Leftarrow b)$  מתקיים כי  $\lim_{m \rightarrow \infty} \mathbb{E}_{S \sim D^m} [L_D(A(S))] = 0$ . נשתמש באי שיויון מרקוב, ונקבל כי לכל  $\varepsilon > 0$

$$P_{S \sim D^m} (L_D(A(S)) \leq \varepsilon) = 1 - P_{S \sim D^m} (L_D(A(S)) \geq \varepsilon) \geq 1 - \frac{\mathbb{E}_{S \sim D^m} [L_D(A(S))]}{\varepsilon}$$

ומכיוון שמתקיים כי  $\lim_{m \rightarrow \infty} \frac{\mathbb{E}_{S \sim D^m} [L_D(A(S))]}{\varepsilon} = 0$ , מהנתון ומאריתמטיקה של גבולות, נקבל כי

$$\lim_{m \rightarrow \infty} P_{S \sim D^m} (L_D(A(S)) \leq \varepsilon) \geq 1$$

כלומר, קיים  $m'$  גדול מספיק כך ש- $P_{S \sim D^m} (L_D(A(S)) \leq \varepsilon) \geq 1 - \delta$  לכל  $m \geq m'$ , ולכן לכל  $\delta > 0$  ו- $\varepsilon > 0$  נוכל למצוא  $m'$  גדול מספיק התלוי בהם, נסמנו  $m(\varepsilon, \delta)$  כך שלכל  $m$  הגדול ממנו יתקיים כי

$$P_{S \sim D^m} [L_D(A(S)) \leq \varepsilon] \geq 1 - \delta$$

כפי שרצינו להוכיח.

$(b \Leftarrow a)$  מתקיים כי  $\forall \varepsilon, \delta > 0, \exists m(\varepsilon, \delta) \text{ s.t. } \forall m \geq m(\varepsilon, \delta) \quad P_{S \sim D^m} [L_D(A(S)) \leq \varepsilon] \geq 1 - \delta$  ונראה שלכל  $\varepsilon, \delta > 0$  ולכל  $m \geq m(\varepsilon, \delta)$  מתקיים כי  $\mathbb{E}_{S \sim D^m} [L_D(A(S))] \leq \varepsilon + \delta$ . נשים לב שמנוסחת התוחלת נקבל יהיו  $\varepsilon, \delta > 0$  ו- $m(\varepsilon, \delta)$  משתנה התלוי בהם. נשים לב שמנוסחת התוחלת נקבל

$$\mathbb{E}_{S \sim D^m} [L_D(A(S))] \leq P_{S \sim D^m} [L_D(A(S)) \leq \varepsilon] \cdot \varepsilon + P_{S \sim D^m} [L_D(A(S)) > \varepsilon] \cdot 1 \leq$$

מההנחה מתקיים כי

$$P_{S \sim D^m} [L_D(A(S)) > \varepsilon] < \delta$$

ולכן בסך הכל נקבל עבור ביטוי זה

$$\leq \varepsilon \cdot 1 + P_{S \sim D^m} [L_D(A(S)) > \varepsilon] \leq \varepsilon + \delta$$

ולכן כפי שנראה בסעיפים הבאים נובע כי

$$\lim_{m \rightarrow \infty} \mathbb{E}_{S \sim D^m} [L_D(A(S))] \leq 0 + 0$$

אבל מאי שליליות התוחלת נקבל סך הכל

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

כפי שרצינו להוכיח.

**2.**

יהיו  $\mathcal{Y} := \{0, 1\}$  ו- $\mathcal{H}$  מחלקה של עיגולים במרחב, כלומר  $\mathcal{H} := \{h_r : r \in \mathbb{R}_+, h_r(x) = 1_{[\|x\|_2 \leq r]}\}$ . נוכיח כי  $\mathcal{H}$  הינה למידה PAC, ושה-*sample complexity* שלה חסום על ידי  $m_{\mathcal{H}(\varepsilon, \delta)} \leq \frac{\log(\frac{1}{\delta})}{\varepsilon}$ .

פתרון:

נוכיח זאת ישירות מהגדרת ה-*PAC learnability* על ידי כך שנראה אלגוריתם ספציפי וניתוח ה-*sample complexity* בדומה לאופן שראינו בתרגול עם האינטרוול. בנוסף, נזכור שלכל  $\varepsilon > 0$  מתקיים כי  $1 - \varepsilon \leq e^{-\varepsilon}$ .

1. אלגוריתם הלמידה: בהינתן  $S = \{(x_i, y_i)\}_{i=1}^m$  אלגוריתם הלמידה יחזיר את העיגול  $\hat{C}$  שנותן את ההתאמה ההדוקה ביותר לדגימות החיוביות, כלומר העיגול בעל השטח הקטן ביותר שיכלול בתוכו את כל הדגימות עם לייבל 1 ולא יכלול בתוכו את כל הדגימות עם הלייבל 0. כמובן שאם לא קיימים כאלה נחזיר  $\emptyset$ . באופן פורמלי, יחזיר  $\hat{h}_{\hat{r}}$  כך ש- $\hat{r}$  מגדיר את רדיוס המעגל האופטימלי, שאותו נרצה להביא למינימום. באופן פורמלי, האלגוריתם יחזיר  $\hat{h}_{\hat{r}}$  כך ש-

$$\hat{r} = \max_{i: y_i=1} \|x_i\|_2$$

2. הוכחה למידות PAC: נראה שלכל עיגול נכון  $C$ , לכל הסתברות  $D$  מעל  $\mathcal{X}$  ולכל  $\varepsilon, \delta \in (0, 1)$ , אם נדגום באופן אחיד ובלתי תלוי דגימות מ- $D$  עם הסתברות של לפחות  $1 - \delta$ , העיגול  $\hat{C}$  שיוחזר על ידי האלגוריתם שהגדרנו לעיל יהיה בעל טעות של לכל היותר  $\varepsilon$ . משלב זה, כאשר נתאר מרחב בין שני עיגולים  $C$  ו- $C'$ , זו תהיה הטבעת מרחב שחיה בין שפת העיגול הפנימי לשפת העיגול החיצוני. נתחיל עם ההבחנה הבאה: העיגול המתאים ההדוק ביותר  $\hat{C}$  תמיד מוכל בעיגול הנכון  $C$ . לכן, הטעות של האלגוריתם יכולה לבוא אך ורק מדגימות בעלות לייבל חיובי שנפלו בשטח של  $C \setminus \hat{C}$ , כלומר מתקיים כי  $L_D(h_S) = \mathbb{P}(x \in C \setminus \hat{C})$ . באופן אינטואיטיבי, נרצה לטעון שבהינתן דאטה סט גדול מספיק, לא סביר שההסתברות תחת  $D$  של דגימות חיוביות ב- $C \setminus \hat{C}$  הוא גבוה, ולכן שגיאת ההכלה לא אמורה להיות גבוהה מידי. לכן, נרצה למצוא ערך של  $m$  שיבטיח שבהסתברות  $1 - \delta$  שגיאת ההכלה תהיה קטנה מ- $\varepsilon$ , לכל  $\varepsilon, \delta \in (0, 1)$ . נרצה לחשב את ההסתברות לטעות, שזו כפי שראינו מתקיימת עבור דגימה  $(\hat{r}, r)$ . נקבל כי

$$P_{S \sim D^m} [L_D(S) \leq \varepsilon] = 1 - P_{S \sim D^m} [L_D(h_S) > \varepsilon]$$

בשלב זה נשים לב כי אם ההסתברות של דגימה מ- $(\hat{r}, r)$  היא  $\varepsilon$ , אזי  $S$  לא יכלול דגימה מ- $(\hat{r}, r)$  בהסתברות של  $(1 - \varepsilon)^m$ .

כל ש- $S$  יותר גדול, כך נהיה פחות סביר שחלק גדול מ- $D$  לא יהיה ניתן לביטוי ב- $S$ , ובאופן פורמלי

$$P_{S \sim D^m} [L_D^-(h_S) > \varepsilon] = P_{S \sim D^m} \left[ \bigcap_{i=1}^m x_i \notin [\hat{r}, r] \right] = \prod_{i=1}^m P_{x_i \sim D} (x_i \notin [\hat{r}, r]) \leq (1 - \varepsilon)^m \leq e^{(-\varepsilon m)}$$

ולכן נקבל בסך הכל כי

$$P_{S \sim D^m} [L_D(S) \leq \varepsilon] \geq 1 - e^{(-\varepsilon m)} \geq 1 - \delta$$

ומכך נקבל כי כמות הדגימות הדרישה היא

$$m \geq \frac{\log\left(\frac{1}{\delta}\right)}{\varepsilon}$$

ולכן בסך הכל נסיק כי מחלקת ההיפוטזות של העיגול הינה למידה  $PAC$ , עם  $sample\ complexity$  מהצורה

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

כפי שרצינו להוכיח.

## $VC - Dimension$

.3

יהי  $\mathcal{X} := \{0, 1\}^n$  ו- $\mathcal{Y} := \{0, 1\}$  ולכל  $I \subseteq [n]$  נגדיר את פונקציית הזוגיות  $h_I(x) = \left(\sum_{i \in I} x_i\right) \bmod 2$ . נראה מהו ה- $VC - dimension$  של מחלקת ההיפוטזות  $\mathcal{H}_{parity} = \{h_I : I \subseteq [n]\}$ .

**פתרון:** ראשית, נשים לב שגודל מחלקת ההיפוטזות הוא  $2^n$ , שכן קיימת פונקציה  $h_I$  עבור כל תת קבוצה  $I$  ב- $n$ , ואנו יודעים שקיימות  $2^n$  תתי קבוצות כאלו. נרצה להראות כי ה- $VC - dimension$  של  $\mathcal{H}_{parity}$  הוא  $n$ . לשם כך, נצטרך להראות ש- $\mathcal{H}$  מנטצת קבוצה בגודל  $n$ , ושהיא לא מנטצת כל קבוצה בגודל  $n + 1$ .  
1. נקח את קבוצת וקטורי הבסיס הסטנדרטי  $C = \{e_1, \dots, e_n\}$ . עבור כל קבוצת לייבלים  $(y_1, \dots, y_n) \in \{0, 1\}^n$  נשים לב שנוכל לקבל אותה באופן הבא: עבור כל לייבל  $j$  שערך  $1$ , ההיפוטזה  $h_{\{j\}}$  תקבל אותו, שכן היא סוכמת את כל הכניסות ה- $j$  של וקטורי הבסיס הסטנדרטי, וברור ששכום זה שווה ל- $1$ . באופן כללי, לקבל קבוצת לייבלים עם ערכי  $1$  במקומות שונים, ההיפוטזה תהיה זו שעבורה  $I$  הינה קבוצת הערכים שהינם בדיוק אלו שערךם אחד בכניסות המתאימות. לכן הראינו כי  $C$  מנוטצת על ידי מחלקת ההיפוטזות  $\mathcal{H}$ , כלומר מתקיים כי  $VC - Dim(\mathcal{H}) \geq n$ .  
2. מצד שני, ראינו בתרגול הוכחה לטענה כי עבור מחלקת היפוטזות  $\mathcal{H}$  בגודל סופי, מתקיים כי  $VC - Dim(\mathcal{H}) \leq \log_2 |\mathcal{H}|$ . מכיוון שמתקיים במקרה זה כי  $|\mathcal{H}| = 2^n$ , הרי שנקבל מטענה זו כי  $VC - Dim(\mathcal{H}) \leq n$ , ולכן קיבלנו בסך הכל כי

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

כפי שרצינו להוכיח.

.4

בהינתן מספר  $k$ , יהי  $([a_i, b_i])_{i=1}^k$  קבוצה של  $k$  אינטרוולים על  $\mathbb{R}$ , ונגדיר את האיחוד שלהם  $A = \bigcup_{i=1}^k [a_i, b_i]$ . מחלקת ההיפוטזות  $H_{k-intervals}$  כוללת את הפונקציות  $h_A(x) = 1_{[x \in A]}$  לכל בחירות של  $k$  אינטרוולים. נמצא את ה- $VC - dimension$  של  $H_{k-intervals}$  ונוכיח את התשובה. בנוסף, נראה כי אם ניתן ל- $A$  להיות איחוד סופי כלשהו של אינטרוולים (כלומר  $k$  אינו מוגבל), אזי הקלאס המתקבל  $H_{intervals}$  הנו בעל  $VC - dimension = \infty$ .

**פתרון:** נרצה להראות כי ה- $VC - Dimension$  של  $H_{k-intervals}$  הוא  $2k$ . לשם כך, נצטרך להראות ש- $\mathcal{H}$  מנטצת קבוצה בגודל  $2k$ , ושהיא לא מנטצת כל קבוצה בגודל  $2k + 1$ .  
1. נקח  $C$  קבוצה כלשהי של דגימות בגודל  $2k$  המסודרת בסדר עולה, ונראה שכל קבוצת לייבלים  $(y_1, \dots, y_{2k}) \in \{0, 1\}^{2k}$  נוכל לקבל באופן הבא: ראשית, נשים לב שבמקרה המורכב ביותר הלייבלים יהיו מסודרים בזיגזג, כלומר  $0, 1, 0, 1, \dots, 0, 1$ . עבור  $k$  ערכי ה- $1$  נצטרך אינטרוול שיתפוס כל דגימה בנפרד, וביניהם עבור דגימות שיקבלו לייבל  $0$  יהיו קטעים ללא אינטרוול קיים. שנית, קל לשים לב שבכל מקרה אחר, בו יש שני לייבלים סמוכים שערךם זהה, אם מדובר בערך  $1$  נוכל להגדיר אינטרוול אחד שיתפוס את שניהם, או אם מדובר בערך  $0$  פשוט לא יהיה שום אינטרוול בקטע זה, ולכן נוכל בצורה

קלה יותר לקבל את קבוצת הלייבלים, ובמקרים מסוימים אפילו לא יהיה צורך בכל  $k$  האינטרוולים כדי לקבל את קבוצת הלייבלים. לכן נקבל כי  $VC - Dim(\mathcal{H}_{k-intervals}) \geq 2k$ .

2. נראה כי ה- $VC - dimension$  של  $\mathcal{H}_{k-interval}$  הוא לכל היותר  $2k$ . באופן אינטואיטיבי, גודל הלייבלים שלא נוכל כבר לקבל צריך לכלול  $k + 1$  ערכי 1 עם רווחים של לייבלים בעלי ערך 0 ביניהם, שכן במקרה זה נצטרך  $k + 1$  אינטרוולים זרים.

לכן, נניח בשלילה שקיימת קבוצת לייבלים בגודל  $2k + 1$  המנוטצת על ידי  $\mathcal{H}_{k-intervals}$ , כלומר, נוכל לקבל את קבוצת הלייבלים  $(1, 0, 1, \dots, 0, 1)$ . זאת כמובן בסתירה שאין ברשותנו  $k + 1$  אינטרוולים זרים, שכן במקרה זה קיימות לנו  $k + 1$  דגימות שנצטרך לקבל, כשבין כל אחת לשנייה דגימה שנצטרך לא לקבל. לכן קיבלנו בסך הכל כי  $VC - Dim(\mathcal{H}_{k-interval}) = 2k$ .

בנוסף לכל, נוכל כעת לשים לב שאם  $k$  כלל אינו מוגבל עבורנו, נוכל עבור כל קבוצת לייבלים  $y$  להתאים אינטרוול עבור כל לייבל שערכו 1 ולדאוג שלא יותאמו אינטרוולים לדגימות שערכן 0, זאת בצורה הכי ישירה שיש, ולכן במקרה זה  $\mathcal{H}_{k-interval}$  תוכל לנטץ כל גודל של קבוצת דגימות  $C$ . כלומר במצב זה יתקיים כי  $VC - dimension(\mathcal{H}_{k-interval}) = \infty$ .

## Monotonicity

5.

תהי  $\mathcal{H}$  מחלקת היפוטוזות עבור קלסיפיקציה בינארית. נניח כי  $\mathcal{H}$  הינה למידה  $PAC$  וכי ה- $sample complexity$  שלה נתון על ידי  $m_{\mathcal{H}}(\cdot, \cdot)$ .

נראה כי  $m_{\mathcal{H}}$  הינה מונוטונית יורדת בכל אחד מהפרמטרים שלה. במילים אחרות, נראה כי:

(א) בהינתן  $\delta \in (0, 1)$  ו- $0 < \varepsilon_1 \leq \varepsilon_2 < 1$ , מתקיים כי  $m_{\mathcal{H}}(\varepsilon_1, \delta) \geq m_{\mathcal{H}}(\varepsilon_2, \delta)$ .

(ב) בהינתן  $\varepsilon \in (0, 1)$  ו- $0 < \delta_1 \leq \delta_2 < 1$ , מתקיים כי  $m_{\mathcal{H}}(\varepsilon, \delta_1) \geq m_{\mathcal{H}}(\varepsilon, \delta_2)$ .

פתרון: א. נניח בשלילה כי  $\varepsilon_1 \leq \varepsilon_2$  אבל מתקיים כי  $m_{\mathcal{H}}(\varepsilon_1, \delta) < m_{\mathcal{H}}(\varepsilon_2, \delta)$ .

על פי ההגדרה,  $m_{\mathcal{H}}(\varepsilon, \delta)$  מייצג את מספר הדגימות המינימלי הדרוש על מנת לקבל כי  $P_{S \sim D^m}(L_D(A(S)) \leq \varepsilon) \geq 1 - \delta$ , ובמילים אחרות להבטיח בהסתברות  $1 - \delta$  שגיאת ההכללה תהיה קטנה מ- $\varepsilon$ . לכן לפי ההנחה עבור  $\delta \in (0, 1)$  נקבל כי

$$P_{S \sim D^m}(L_D(A(S)) \leq \varepsilon_2) \geq P_{S \sim D^m}(L_D(A(S)) \leq \varepsilon_1) \geq 1 - \delta$$

כלומר השתמשנו בפחות דגימות על מנת לקבל חסם טוב יותר על על הסתברות שגיאה ההכללה, וזו סתירה.

ב. נניח בשלילה כי  $\delta_1 \leq \delta_2$  אבל מתקיים כי  $m_{\mathcal{H}}(\varepsilon, \delta_1) < m_{\mathcal{H}}(\varepsilon, \delta_2)$ .

על פי ההגדרה,  $m_{\mathcal{H}}(\varepsilon, \delta)$  מייצג את מספר הדגימות המינימלי הדרוש על מנת לקבל כי  $P_{S \sim D^m}(L_D(A(S)) \leq \varepsilon) \geq 1 - \delta$ , ובמילים אחרות להבטיח בהסתברות  $1 - \delta$  שגיאת ההכללה תהיה קטנה מ- $\varepsilon$ . לכן לפי ההנחה עבור  $\varepsilon \in (0, 1)$  נקבל כי

$$P_{S \sim D^m}(L_D(A(S)) \leq \varepsilon) \geq 1 - \delta_1 \geq 1 - \delta_2$$

כלומר השתמשנו בפחות דגימות על מנת לקבל חסם טוב יותר על על הסתברות שגיאה ההכללה, וזו סתירה.

6.

תהינה  $\mathcal{H}_1$  ו- $\mathcal{H}_2$  שתי מחלקות עבור קלסיפיקציה בינארית, כך ש- $\mathcal{H}_1 \subseteq \mathcal{H}_2$ . נראה כי  $VC - dim(\mathcal{H}_1) \leq VC - dim(\mathcal{H}_2)$ .

פתרון: נסמן  $VC - dim(\mathcal{H}_1) = k_1$  ו- $VC - dim(\mathcal{H}_2) = k_2$ , ונניח בשלילה כי  $k_1 > k_2$ . מההגדרה של  $VC - Dimension$ , קיים קבוצת דגימות  $C = \{c_1, \dots, c_{k_1}\}$  כך ש- $\mathcal{H}_1$  מנוטצת את  $C$  אבל  $\mathcal{H}_2$  לא, כלומר קיימת היפוטזה  $h \in \mathcal{H}_1$  כך שבהינתן קבוצת הדגימות  $C$  יכולה לקבל כל וקטור לייבלים בינארי אפשרי בגודל  $k_1$ . בנוסף, נתון כי  $\mathcal{H}_1 \subseteq \mathcal{H}_2$ , כלומר מתקיים כי  $h \in \mathcal{H}_2$ , כלומר גם  $\mathcal{H}_2$  מנוטצת את  $C$ , וזו סתירה.

.7

נוכיח כי אם  $\mathcal{H}$  בעלת התכנסות במידה שווה עם הפונקציה  $m_{\mathcal{H}}^{UC} : (0, 1)^2 \rightarrow \mathbb{N}$ , אזי  $\mathcal{H}$  הינה Agnostic – PAC למידה עם  $m_{\mathcal{H}}(\varepsilon, \delta) \leq m^{UC}(\frac{\varepsilon}{2}, \delta)$  sample complexity.

**פתרון:** מתכונת ההתכנסות במידה שווה של  $H$  עם  $m_{\mathcal{H}}^{UC}$  מתקיים כי לכל  $\varepsilon, \delta \in (0, 1)$  ולכל פונקציית התפלגות  $D$  מעל  $\mathcal{X}$  כי אם  $S = \{(x_i, y_i)\}_{i=1}^m$  סט דגימות בגודל  $m \geq m_{\mathcal{H}}^{UC}$  הנדגמות באופן אחיד ובלתי תלוי מ- $D$  מתקיים כי  $1 - \delta \leq \Pr\{S \text{ is } \varepsilon \text{ representative}\}$ , כלומר מהגדרת  $\varepsilon$  representative לכל  $h \in \mathcal{H}$   $|L_{S_m}(h) - L_D(h)| \leq \varepsilon$ . מכיוון שהטענה נכונה לכל  $\varepsilon \in (0, 1)$ , נקח  $\varepsilon' = \frac{\varepsilon}{2} \leq \varepsilon$  ונרצה להראות כי על פי הגדרת למידות Agnostic PAC כי

$$D^m \left( S \in Z^m : \exists h \in \mathcal{H}, L_{S_m}(h) \leq \min_{h \in \mathcal{H}} L_D(h) + \varepsilon \right) \geq 1 - \delta$$

ראשית, נקבל עבור  $\varepsilon'$  representative  $L_D(h_S) \leq \min_{h \in \mathcal{H}} L_D(h) + \varepsilon'$  עבור  $h_S = \operatorname{argmin}_{h \in \mathcal{H}} L_S(h)$ , ולכן מההתכנסות במידה שווה מתקיים

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

ועל כל קיבלנו כי עבור  $h$  כנ"ל  $\mathcal{H}$  הינה למידה Agnostic – PAC עם sample complexity של  $m_{\mathcal{H}}(\varepsilon, \delta) \leq m^{UC}(\varepsilon' = \frac{\varepsilon}{2}, \delta)$ . כנדרש.

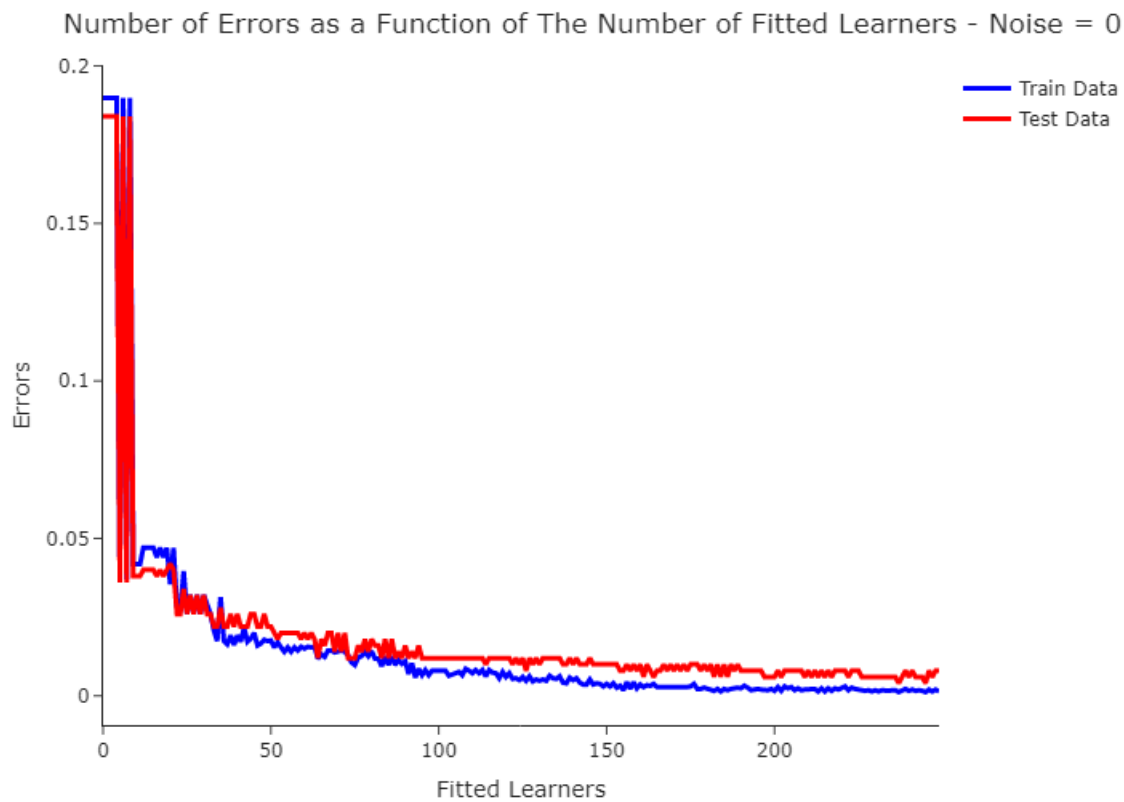
.8

תהי  $\mathcal{H}$  מחלקה היפוטזות מעל  $\mathcal{Z} = \mathcal{X} \times \left\{ \begin{smallmatrix} +1 \\ -1 \end{smallmatrix} \right\}$ , עם פונקציית  $loss_{0-1}$ , ונניח שקיימת פונקציה  $m_{\mathcal{H}}$  כך שלכל התפלגות  $D$  מעל  $\mathcal{Z}$  קיים אלגוריתם  $A$  עם התכונה הבאה: כאשר מריצים את  $A$  על  $m \geq m_{\mathcal{H}}$  דגימות המתקבלות מ- $D$  באופן אחיד ובלתי תלוי, מובטח שתוחזר בהסתברות לפחות  $1 - \delta$  היפוטזה  $h_s : \mathcal{X} \rightarrow \left\{ \begin{smallmatrix} +1 \\ -1 \end{smallmatrix} \right\}$  עם  $L_D(h_s) \leq \min_{h \in \mathcal{H}} L_D(h) + \varepsilon$ . נקבע האם  $\mathcal{H}$  הינה למידה Agnostic – PAC באמצעות הוכחה או דוגמה נגדית.

**פתרון:**  $\mathcal{H}$  איננה למידה Agnostic PAC, ונציג לכך דוגמה נגדית. ראשית, נשים לב שההבדל בתיאור אל מול הגדרת למידות Agnostic PAC כפי שראינו היא שבמקרה זה קיים אלגוריתם  $A$  לכל התפלגות  $D$ , לעומת ההגדרה שבה קיים אלגוריתם לומד  $A$  לכל התפלגות  $D$  מעל  $\mathcal{Z}$ . כעת, תהי  $\mathcal{H}$  מחלקת היפוטזות של פונקציות המתאימות לכל מספר טבעי לייבל ב- $\{-1, 1\}$ . מצד אחד, נוכל לראות כי כל קבוצת דגימות  $C$  בגודל סופי שניקח איננה מנוטצת על ידי  $\mathcal{H}$ , ולכן ישירות מההגדרה שראינו בהרצאה מתקיים כי  $\mathcal{H}$  איננה למידה Agnostic – PAC. מצד שני, בהינתן פונקציה  $m_{\mathcal{H}}$  והתפלגות שרירותית  $D$  מעל  $\mathcal{Z}$ , נוכל לקחת אלגוריתם שיעבור על כל הדגימות ויחזיר את הליבלים הנכונים. במקרה זה השגיאה תהיה מינימלית לכל  $\varepsilon, \delta$ , כלומר התנאים הנ"ל מתקיימים, אף על פי שראינו ש- $\mathcal{H}$  כנ"ל איננה למידה Agnostic PAC, ולכן קיבלנו את מה שרצינו להוכיח.

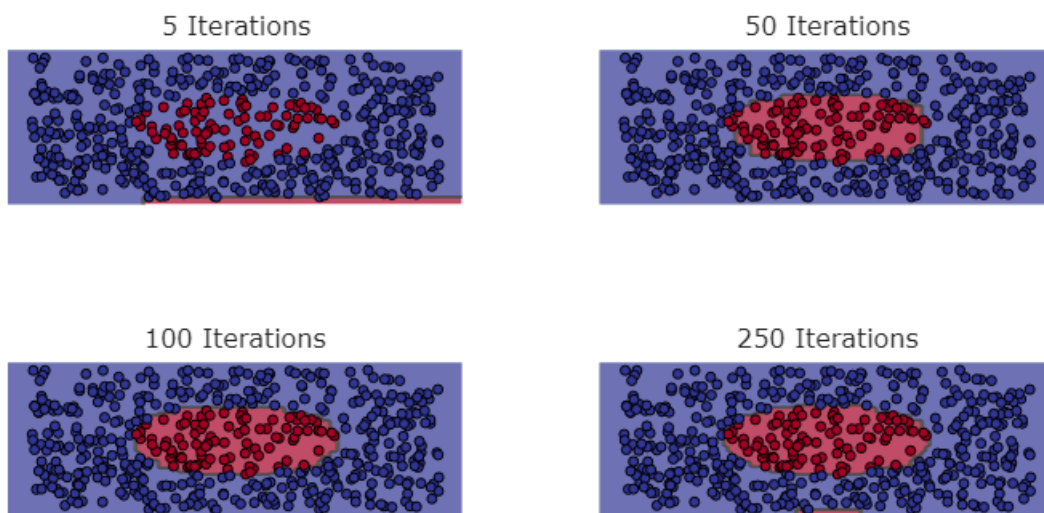
## החלק המעשי

.1



נוכל לראות שככל שנגדיל את מספר ה-*Weak learners* שאנו משתמשים בהם ב-*Adaboost*, כך נקטין את השגיאה הן על ה-*Test* והן על ה-*Train*, כאשר לא כדאי להסתמך על כמות נמוכה מאוד של לומדים.

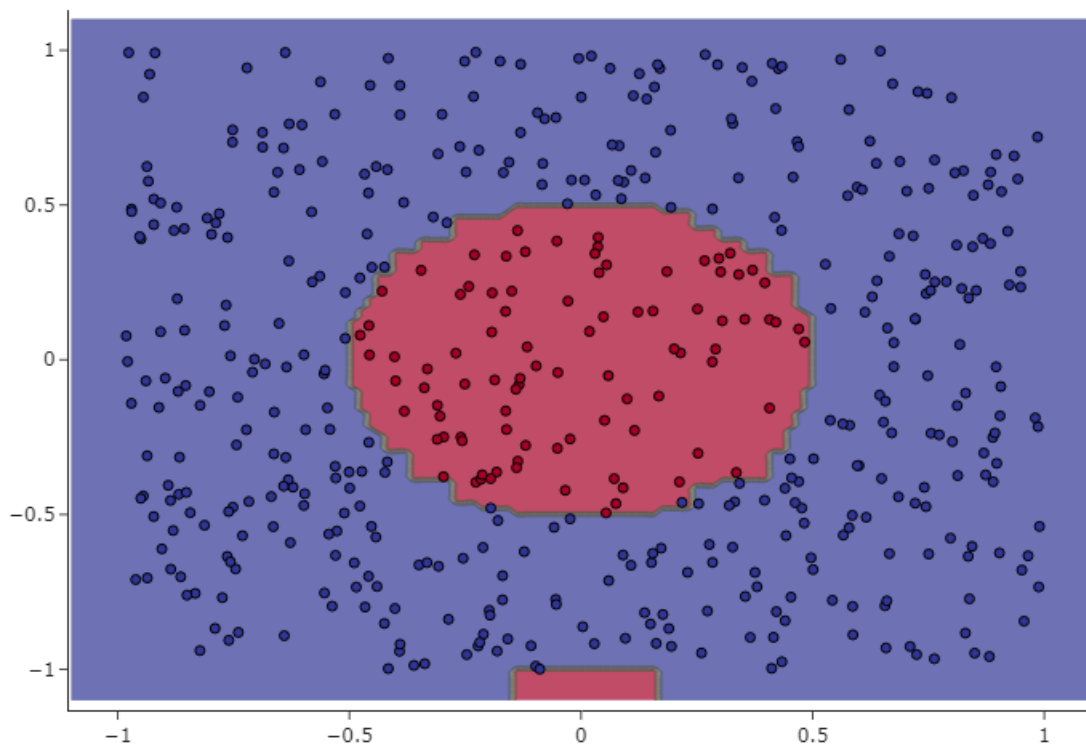
## Decision Boundary with Increasing Number of Iterations - with Noise=0



גם כאן אנו רואים שהגדלת כמות האיטרציות שאנו מבצעים, המובילה לשימוש ביותר *weak learners*, תשפר את ה-*Decision Boundary*, בפרט על ה-*Test Set* ולכן נראה שלא מדובר מדובר ב-*Over fitting*. נשים את הדגש על כך שמדובר בדגימות ללא רעש, מה שלא מתרחש במציאות, ולכן יותר פשוט ל-*Adaboost* להשתמש ב-*weak learners* ולקבל תוצאה כזו טובה.

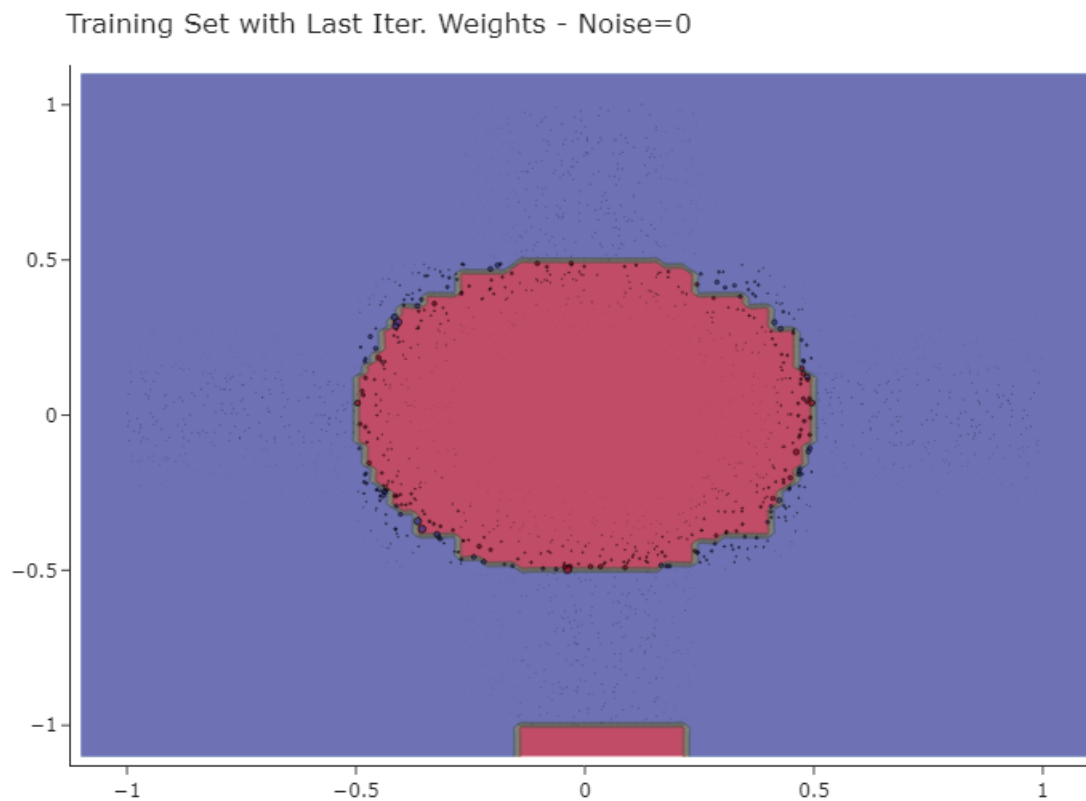
.3

Ensemble with Lowest Test Error Is On Size 238 with Accuracy 0.996 - Noise=



ניתן לראות מגרף זה שעבור שימוש ב-238 *weak learners* נקבל את הדיקט הטוב ביותר, או במילים השגיאיות השגיאה תהיה הקטנה ביותר שהיא 0.004.  
במקרה זה כפי שניתן, כאשר אין רעש, *Adaboost* מצליח לאמוד בצורה טובה מאוד את התפלגות ה-*Test data*.

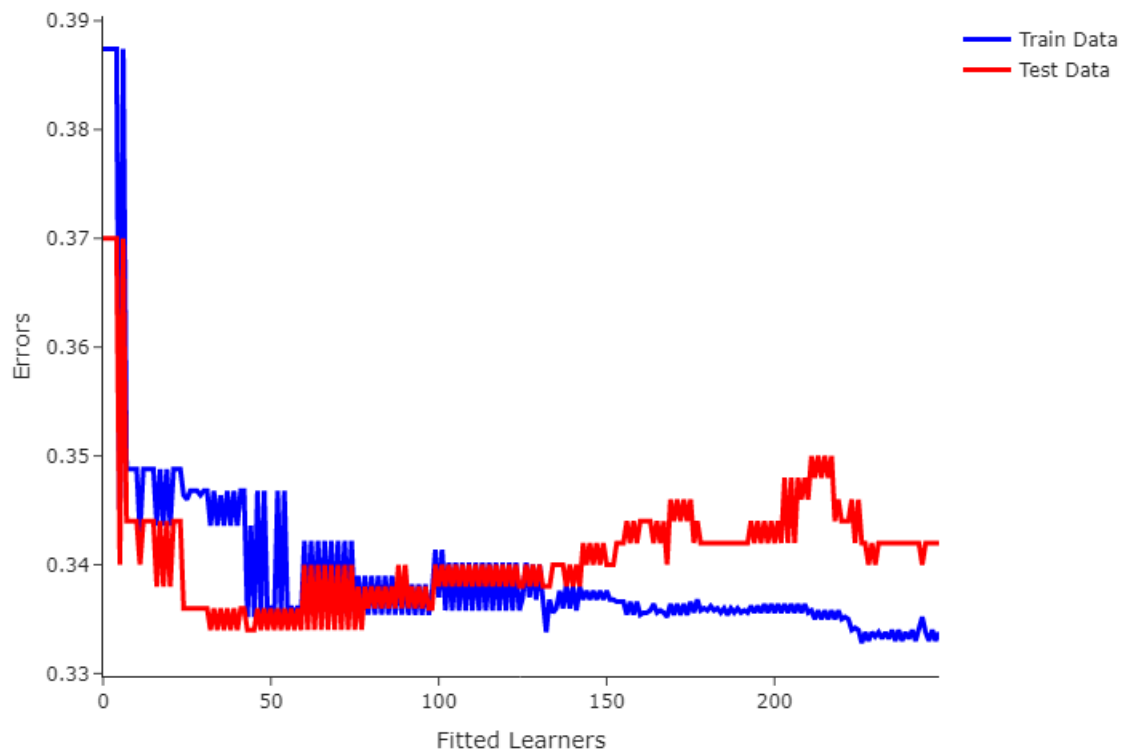




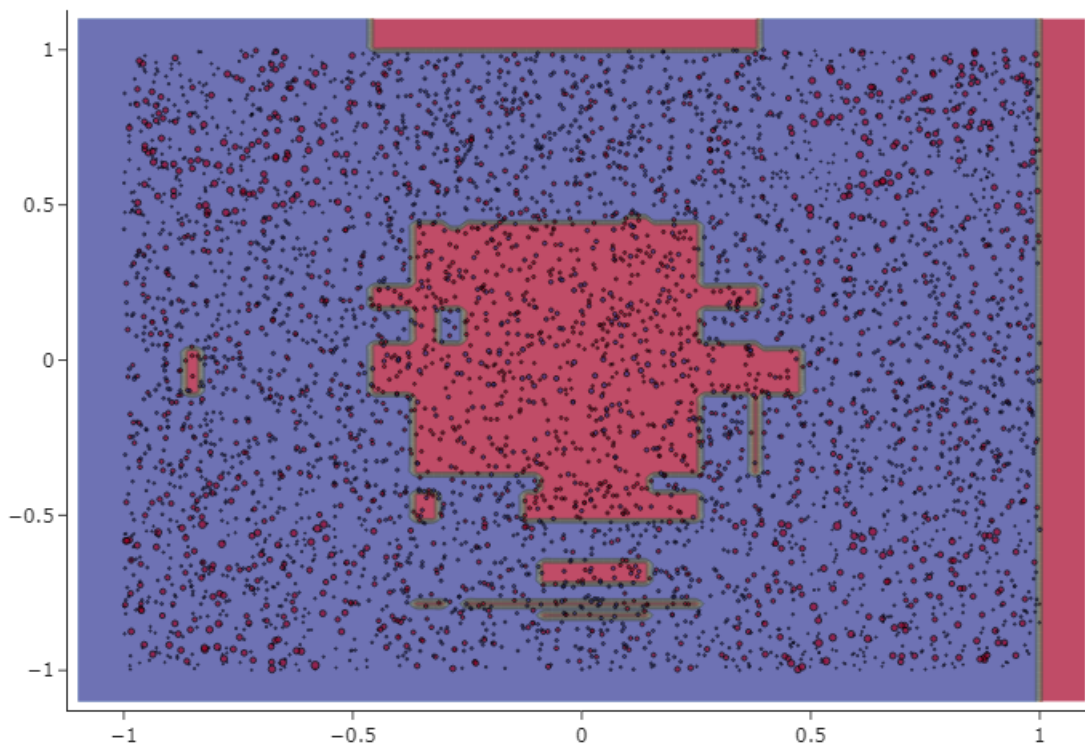
כפי שראינו גם בתרגול ובהרצאה, האלגוריתם של *Adaboost* פועל כך שהמשקל יגדל עבור דגימות האומד לא מצליח לסווג נכון, ולכן בגרף זה נוכל לראות דגימות אלה כגדולות ביותר. על כן הדגימות שקלות עבור המודל הן אלו הקטנות יתר, שרובן בצורת + בשטח הכחול הרחוק, ואלו שמאתגרות את האומד הן אלו שנמצאות בעיקר על שפת השטח האדום, ומוצאות כגדולות.

.5

Number of Errors as a Function of The Number of Fitted Learners - Noise = 0.



Training Set with Last Iter. Weights - Noise=0.4



ראשית, נשים לב ש-*Adaboost* מתקשה מאוד במקרה זה לאמוד את התפלגות הדגימות, שכן בגרף הראשון אנו רואים שהשגיאה גדולה מאוד ביחס למקרה בו אין רעש, גם לאחר 250 איטרציות, אפילו עבור ה-*Train set*. שנית, נוכל לראות כי כאשר אנו משתמשים במספר גדול של *weak learners*, קרי במקרה זה 250, האומד יתאים את עצמו ל-*Train set* ויבצע *Overfitting* עליו. בנוסף, כפי שנוכל לראות מהגרף השני, ככל שמספר זה יהיה גדול יותר כך יעלה ה-*variance*, שכן יהיה שינוי גדול עבור כל שינוי בדגימות, לעומת ירידה ב-*bias* שכן ההתאמה ל-*Train set* מתהדקת.