

# Contents

<a href="#">1 Ex4 Answers.pdf</a>	2
<a href="#">2 adaboost.py</a>	21
<a href="#">3 ex4 tools.py</a>	24

## מבוא למערכות לומדות תרגיל 4

עמית בסקין 312259013

### 1 למידות-PAC

1. יהיו  $A$  אלגוריתם למידה ו-  $\mathcal{D}$  התפלגות כלשהי.

נניח שפונקציית ה- Loss נמצאת בטווח  $[0, 1]$ .

צ. להוכיח ששתי הטענות הבאות שקולות:

- לכל  $\epsilon, \delta > 0$  יש  $m(\epsilon, \delta)$  כך שלכל  $m \geq m(\epsilon, \delta)$  מתקיים:

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

- מתקיים:

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

הוכחה:

כיוון ראשון: נניח שלכל  $\epsilon, \delta > 0$  יש  $m(\epsilon, \delta)$  כך שלכל  $m \geq m(\epsilon, \delta)$  מתקיים:

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

ונראה ש-

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

אכן:

יהא  $m \geq m(\epsilon, \delta)$ . ניקח  $\epsilon = \delta = \frac{1}{m}$ . נסמן  $x = L_{\mathcal{D}}(A(S))$ . נתבונן בפונקציית ההצטברות  $f(x)$  שמקיימת:

$$\mathbb{P}(a \leq x \leq b) = \int_a^b f(x) dx$$

נתון שפונקציית ה-Loss נמצאת בטווח  $[0, 1]$  ולכן זהו תחום האינטגרציה. אז נקבל:

$$\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S))] = \int_0^1 x f(x) dx$$

נפצל לתחומים:

$$\begin{aligned} &= \int_0^{1/m} x f(x) dx + \int_{1/m}^1 x f(x) dx \leq \\ &\leq \frac{1}{m} \int_0^{1/m} f(x) dx + \int_{1/m}^1 f(x) dx \end{aligned}$$

אבל לכל  $a, b$  מתקיים:

$$\mathbb{P}(a \leq x \leq b) = \int_a^b f(x) dx \leq 1$$

ולכן:

$$\frac{1}{m} \int_0^{1/m} f(x) dx + \int_{1/m}^1 f(x) dx \leq \frac{1}{m} + \int_{1/m}^1 f(x) dx$$

ועתה מההנחה:

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

נובע ש-

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left[ L_{\mathcal{D}}(A(S)) > \frac{1}{m} \right] < \frac{1}{m}$$

קרי:

$$\frac{1}{m} + \int_{1/m}^1 f(x) dx < \frac{1}{m} + \frac{1}{m} = \frac{2}{m}$$

ניקח את  $m$  לאינסוף ונקבל:

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

כנדרש.

כיוון שני: נניח ש-

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

ונראה שלכל  $\epsilon, \delta > 0$  יש  $m(\epsilon, \delta)$  כך שלכל  $m \geq m(\epsilon, \delta)$  מתקיים:

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

אכן: יהיו  $\delta > 0, \epsilon$ . לפי א"ש מרקוב מתקיים:

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

עתה מההנחה ש-

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

אזי שקיים  $M$  כך שלכל  $m > M$  מתקיים ש-

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

אזי יהא  $m > M$  ונקבל:

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

קרי:

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

כלומר:

$$m(\epsilon, \delta) = M$$

מש"ל.

2. יהיו  $\mathcal{X} = \mathbb{R}^2, \mathcal{Y} = \{0, 1\}$  ו-  $\mathcal{H} = \{h_r \mid r \in \mathbb{R}_+\}$  מחלקה של היפותיזות כאשר  $h_r(x) = \mathbf{1}[\|x\|_2 \leq r]$ .

צלהוכיח ש-  $\mathcal{H}$  היא למידה-PAC וש-

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

הוכחה:

מהנחת הריאליזביליות יש  $h_r \in \mathcal{H}$  עם  $L_{\mathcal{D}}(h_r) = 0$ .

תהא  $S = \{(x_i, y_i)\}_{i=1}^m$  קבוצה של  $m$  דגימות שנדגמו באופן זהה ובלתי תלוי מהתפלגות  $\mathcal{D}$  מעל  $\mathcal{X} \times \mathcal{Y}$ .

נשים לב שמאחר ש-  $L_{\mathcal{D}}(h_r) = 0$  אזי ש-

$$\mathbb{P}_{(x,y) \sim \mathcal{D}} [\|x\|_2 > r \text{ and } y = 1] = 0$$

ניקח:

$$r_{alg} = \begin{cases} \max_{i \in [m]} \{\|x_i\|_2 \in \mathbb{R}_+ \mid y_i = 1\} & \exists i \in [m] (y_i = 1) \\ 0 & \text{else} \end{cases}$$

יהא  $r'$  ממשי אי-שלילי עם  $r' < r$  כך ש-

$$\mathbb{P}_{x \sim \mathcal{D}_X} [r' \leq \|x\|_2 \leq r] = \epsilon$$

ונשים לב שכדי שיתקיים  $L_{\mathcal{D}}(h_{r_{alg}}) \leq \epsilon$  מספיק שיתקיים ש- $r_{alg} \geq r'$  כלומר:

$$\mathbb{P}[L_{\mathcal{D}}(h_{r_{alg}}) > \epsilon] \leq \mathbb{P}[r_{alg} < r']$$

אבל  $r' < r_{alg}$  אם לכל  $(x_i, y_i) \in S$  מתקיים ש- $\|x_i\|_2 \notin [r', r]$  קרי:

$$= \mathbb{P}_{S \sim \mathcal{D}^m} [\forall i \in [m] (\|x_i\|_2 \notin [r', r])] = (1 - \epsilon)^m \leq e^{-\epsilon m}$$

ומאחר שאנחנו מניחים ש- $m > \frac{\log(1/\delta)}{\epsilon}$  אזי ש-

$$\leq e^{-\epsilon \frac{\log(1/\delta)}{\epsilon}} = e^{-\log(1/\delta)} = \frac{1}{e^{\log(1/\delta)}} = \frac{1}{1/\delta} = \delta$$

ומכאן ש- $\mathcal{H}$  הינה למידה-PAC וכן:

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

כנדרש. ■

## 2 מימד VC

3. יהיו  $\mathcal{X} = \{0, 1\}^d$ ,  $\mathcal{Y} = \{0, 1\}$  עם  $d \geq 2$ .

כל דגימה  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  מורכבת מהשמה של  $d$  משתנים בוליאניים ותויות.

לכל משתנה בוליאני  $x_k$  עם  $k \in [d]$  ישנם שני ליטרלים:  $x_k$  ו- $\overline{x_k} = 1 - x_k$ .

נגדיר מחלקה  $\mathcal{H}_{\text{con}}$ : כל היפותזה ניתנת על ידי גימומים של תת-קבוצה כלשהי מתוך  $2d$  הליטרלים הנ"ל.

צ.לחשב את מימד ה-VC של  $\mathcal{H}_{\text{con}}$ .

תזכורות:

- תהא  $C \subseteq \mathcal{X}$ . הצמצום של  $\mathcal{H}$  ל- $C$  הוא:

$$\mathcal{H}_C = \left\{ h_C \in \{0, 1\}^C \mid \exists h \in \mathcal{H} \text{ s.t. } \forall x \in C (h_C(x) = h(x)) \right\}$$

- נאמר ש- $\mathcal{H}$  מנתצת קבוצה סופית  $C \subseteq \mathcal{X}$  אם  $|\mathcal{H}_C| = 2^{|C|}$ .

- הגדרנו:

$$\text{VC dim}(\mathcal{H}) = \sup \{m \in \mathbb{N} \mid \exists C \subseteq \mathcal{X} (|C| = m \text{ and } \mathcal{H} \text{ shatters } C)\}$$

- כדי להראות ש- $\text{VC dim}(\mathcal{H}) = d$  צריך להראות ש-

- יש  $C \subseteq \mathcal{X}$  מגודל  $d$  שמנותצת על ידי  $\mathcal{H}$ , קרי  $\text{VC dim}(\mathcal{H}) \geq d$

- כל  $C \subseteq \mathcal{X}$  מגודל  $d+1$  אינה מנותצת על ידי  $\mathcal{H}$ , קרי  $\text{VC dim}(\mathcal{H}) \leq d$

פיתרון: נטען כי  $\text{VC dim}(\mathcal{H}) = d$ . הוכחה:

- $\text{VC dim}(\mathcal{H}) \geq d$ : ניקח את קבוצת הוקטורים הסטנדרטיים ב- $\mathbb{R}^d$ , קרי:  $C = \{e_i \in \mathbb{R}^d \mid i \in [d]\}$  ונראה שניתן לנתץ אותה על ידי  $\mathcal{H}$ :

אכן, יהא  $y \in \{0, 1\}^d$  וקטור תוויות עבור הוקטורים ב- $C$ . אז ניקח:

$$h = \bigwedge_{y_i=0} \overline{x_i}$$

ראשית לכל  $e_i \in C$  ולכל  $j \neq i$  מתקיים שהכניסה ה- $j$  ב- $e_i$  היא אפס, ומאחר שכל המשתנים בגימום של  $h$  הם ליטרליים שליליים אזי שההשמה של כניסה  $j$  ב- $e_i$  בתוך משתנה  $\overline{x_j}$  בגימום של  $h$  (אם הוא מופיע), נותן 1.

עתה, אם  $y_i = 1$  אז  $\overline{x_i}$  לא משתתף בגימום ובסה"כ מתקבל גימום של אחדות ולכן  $h(e_i) = 1$ . מצד שני אם  $y_i = 0$  אז  $\overline{x_i}$  משתתף בגימום ומקבל השמה של 0 ובסה"כ מתקבל גימום של אחדות ואפס ולכן  $h(e_i) = 0$ .

- $\text{VC dim}(\mathcal{H}) \leq d$ : תהא קבוצה  $C$  עם  $|C| = d+1$ . נניח בשלילה שניתן להתאים עבור  $C$  כל וקטור תוויות  $y$ . נמספר את הוקטורים ב-

$$C = \{c_i\}_1^{d+1}. \text{ אז בפרט לכל } c_i \text{ יש } h_C \text{ כך ש- } h_C^i(c_i) = 0 \text{ ולכל } j \neq i \text{ מתקיים } h_C^i(c_j) = 1.$$

לכל היפותיזה  $h_C^i$  כנ"ל, נבחר ליטרל  $z_i$  עליו  $c_i$  מקבל 0, קיים כזה כי  $h_C^i(c_i)$ .

נתבונן בקבוצת הליטרלים שקיבלנו:  $\{z_i\}_1^{d+1}$ . מעקרון שובך היונים יש שני ליטרלים  $z_i, z_j$  שמתייחסים לאותו משתנה. נחלק למקרים:

- אם  $z_i = z_j$  אז  $c_i, c_j$  שניהם מאפסים את  $z_i$  בסתירה לכך שאחד מהם אמור לקבל תיוג 0 והשני 1.

- אם  $z_i = \overline{z_j}$  כ"כ  $z_i = x_k$  ו- $z_j = \overline{x_k}$ . יהא  $l \notin \{i, j\}$ . אז  $l$  מקבל 1 גם  $x_k$  וגם על  $\overline{x_k}$  בסתירה.

מש"ל

### 3 PAC-אגנוסטיות

4. צ.להוכיח שאם ל- $\mathcal{H}$  יש את תכונת ההתפלגות האחידה עם פונקציה

$$m_{\mathcal{H}}^{UC}: (0, 1)^2 \rightarrow \mathbb{N}$$

אז  $\mathcal{H}$  היא למידה-PAC-אגנוסטית עם סיבוכיות של

$$m_{\mathcal{H}}(\epsilon, \delta) \leq m_{\mathcal{H}}^{UC}\left(\frac{\epsilon}{2}, \delta\right)$$

תזכורת:

- קבוצת אימון  $S_m$  נקראת  $\epsilon$ -מייצגת עבור  $\mathcal{D}, \mathcal{H}, L$  אם לכל  $h \in \mathcal{H}$  מתקיים:  $|L_S(h) - L_{\mathcal{D}}(h)| < \epsilon$ .

- נאמר שלמחלקת היפותיזות  $\mathcal{H}$  יש את תכונת ההתכנסות האחידה אם יש  $m_{\mathcal{H}}^{UC}: (0, 1)^2 \rightarrow \mathbb{N}$  כך שלכל  $\epsilon > 0$  ו- $\delta < 1$  ולכל התפלגות  $\mathcal{D}$

על  $\mathcal{X} \times \mathcal{Y}$  מתקיים שאם  $m \geq m_{\mathcal{H}}^{UC}(\epsilon, \delta)$  אז  $m \geq m_{\mathcal{H}}^{UC}(\epsilon, \delta)$   $\mathcal{D}^m(\{S_m \in (\mathcal{X} \times \mathcal{Y})^m \mid S \text{ is } \epsilon\text{-representative}\}) \geq 1 - \delta$

- מחלקת היפותיזות  $\mathcal{H}$  נקראת למידה-PAC-אגנוסטית ביחס לפונקציית Loss  $L: \mathcal{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, \infty)$  אם קיימת פונקציה  $\tilde{m}_{\mathcal{H}}: (0, 1)^2 \rightarrow \mathbb{N}$  ואלגוריתם למידה  $\mathcal{A}_m: (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathcal{H}$  עם התכונה הבאה: לכל  $\epsilon, \delta \in (0, 1)$  ולכל התפלגות  $\mathcal{D}$  מעל  $\mathcal{X} \times \mathcal{Y}$  ולכל  $m \geq \tilde{m}_{\mathcal{H}}(\epsilon, \delta)$  מתקיים ש-  $\mathbb{P}_{S \sim \mathcal{D}^m} \left( \left\{ S_m \mid L_{\mathcal{D}}(h_{S_m}) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \epsilon \right\} \right) \geq 1 - \delta$  כאשר  $S_m = \{(x_i, y_i)\}_1^m$  דגימות שנדגמו בצורה זהה ובלתי-תלויה מ-  $\mathcal{D}$ , וכן  $h_{S_m} = \mathcal{A}_m(S_m)$ .

למה: תהא  $S_m$  קבוצת אימון  $-\frac{\epsilon}{2}$ -מייצגת עבור  $\mathcal{D}, \mathcal{H}, L$ . תהא  $h_{S_m}$  פלט כלשהו של  $ERM_{\mathcal{H}}(S_m)$ , קרי  $h_{S_m} \in \arg \min_{h \in \mathcal{H}} L_{S_m}(h)$ . אז  $L_{\mathcal{D}}(h_{S_m}) < \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \epsilon$ .

הוכחה: מכך ש-  $S_m$  היא  $-\frac{\epsilon}{2}$ -מייצגת עבור  $\mathcal{D}, \mathcal{H}, L$ , אזי שלכל  $h \in \mathcal{H}$  מתקיים:

$$|L_{S_m}(h) - L_{\mathcal{D}}(h)| < \frac{\epsilon}{2} \implies \begin{cases} L_{S_m}(h) - L_{\mathcal{D}}(h) < \frac{\epsilon}{2} & \star_1 \\ L_{\mathcal{D}}(h) - L_{S_m}(h) < \frac{\epsilon}{2} & \star_2 \end{cases}$$

ובפרט זה מתקיים עבור  $h = h_{S_m}$ , אז מ-  $\star_2$  נקבל:

$$L_{\mathcal{D}}(h_{S_m}) - L_{S_m}(h_{S_m}) < \frac{\epsilon}{2}$$

קרי:

$$L_{\mathcal{D}}(h_{S_m}) < L_{S_m}(h_{S_m}) + \frac{\epsilon}{2}$$

עתה ניקח  $h' \in \arg \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h)$ , ומאחר שב-  $h'$  אילוף המינימום חזק יותר, אזי ש-

$$\leq L_{S_m}(h') + \frac{\epsilon}{2}$$

ומ-  $\star_1$  ב-  $\epsilon/2$ -רפרסנטטיביות:

$$\leq L_{\mathcal{D}}(h') + \epsilon$$

כנדרש.

הוכחת השאלה המרכזית:

ל-  $\mathcal{H}$  יש את תכונת ההתכנסות האחידה, קרי יש  $m_{\mathcal{H}}^{UC}: (0, 1)^2 \rightarrow \mathbb{N}$  כך שלכל  $\epsilon > 0$  ו-  $\delta < 1$  ולכל התפלגות  $\mathcal{D}$  על  $\mathcal{X} \times \mathcal{Y}$  מתקיים

$$\mathbb{P}_{S \sim \mathcal{D}^m}(\{S_m \in (\mathcal{X} \times \mathcal{Y})^m \mid S \text{ is } \epsilon\text{-representative}\}) \geq 1 - \delta \text{ אז } m \geq m_{\mathcal{H}}^{UC}(\epsilon, \delta).$$

אז יהיו  $\epsilon, \delta \in (0, 1)$  וניקח  $m \geq m_{\mathcal{H}}^{UC}(\frac{\epsilon}{2}, \delta)$ .

מהלמה אנחנו יודעים שאם  $S_m$  היא  $-\frac{\epsilon}{2}$ -מייצגת אז  $L_{\mathcal{D}}(h_{S_m}) < \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \epsilon$

אז ממונוטוניות פונקציית ההסתברות נקבל ש-

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left[ L_{\mathcal{D}} \left( h_{S_m} \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \epsilon \right) \right] \geq \mathbb{P}^m \left( \left\{ S_m \in (\mathcal{X} \times \mathcal{Y})^m \mid S \text{ is } \frac{\epsilon}{2}\text{-representative} \right\} \right)$$

ומקיום תכונת ההתכנסות האחידה:

$$\geq 1 - \delta$$

כלומר קיבלנו את  $m_{\mathcal{H}}^{UC}: (0,1)^2 \rightarrow \mathbb{N}$  מקיים תכונת ההתכנסות האחידה, ואלגוריתם הלמידה שלנו

$$\mathcal{A}_m: (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathcal{H}$$

הוא  $h_{S_m}$  והוא מקיים שלכל  $\epsilon, \delta \in (0,1)$  ולכל  $\mathcal{X} \times \mathcal{Y}$  מעל  $\mathcal{D}$  ולכל  $m \geq \tilde{m}_{\mathcal{H}}(\epsilon, \delta)$  מתקיים ש-

$$\mathcal{D}^m \left( \left\{ S_m \mid L_{\mathcal{D}}(h_{S_m}) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \epsilon \right\} \right) \geq 1 - \delta$$

כאשר  $S_m = \{(x_i, y_i)\}_1^m$  דגימות שנדגמו בצורה זהה ובלתי-תלויה מ- $\mathcal{D}$ , וזוהי בדיוק ההגדרה של למידות-PAC-אגנוסטית. מש"ל ■

5. תהא  $\mathcal{H}$  מחלקת היפותיזות מעל  $Z = \mathcal{X} \times \{\pm 1\}$ .

נתבונן בפונקציית ההפסד 0-1.

נניח שיש פונקציה  $m_{\mathcal{H}}(\epsilon, \delta)$  כך שלכל התפלגות  $\mathcal{D}$  מעל  $Z$  יש אלגוריתם  $A$  עם התכונה הבאה:

כאשר מריצים את  $A$  על  $m \geq m_{\mathcal{H}}(\epsilon, \delta)$  דגימות שנדגמות בהתפלגות זהה ובאופן בלתי-תלוי שנלקחות על ידי  $\mathcal{D}$ , אז מובטח ש- $A$

יחזיר עם הסתברות של לפחות  $1 - \delta$ , היפותיזה  $h_S: \mathcal{X} \rightarrow \{\pm 1\}$  עם  $L_{\mathcal{D}}(h_S) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \epsilon$

צ.לבדוק האם  $\mathcal{H}$  היא למידה-PAC-אגנוסטית.

#### בחרתי לענות על שאלה 4.

## 4 מונוטוניות

6. תהא  $\mathcal{H}$  מחלקת היפותיזות עבור משימת קלסיפקציה בינארית.

נניח ש- $\mathcal{H}$  היא למידה-PAC ושהסיבוכיות ניתנת על ידי  $m_{\mathcal{H}}(\cdot, \cdot)$ .

צ.להראות ש- $m_{\mathcal{H}}$  היא מונוטונית לא-עולה בכל אחד מהפרמטרים שלה, קרי בהינתן  $\delta \in (0,1)$  ו- $0 < \epsilon_1 \leq \epsilon_2 < 1$  מתקיים ש-

$$m_{\mathcal{H}}(\epsilon_1, \delta) \geq m_{\mathcal{H}}(\epsilon_2, \delta)$$

ובהינתן  $\epsilon \in (0,1)$  ו- $0 < \delta_1 \leq \delta_2 < 1$  מתקיים ש-

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

#### בחרתי לענות על שאלה 7.

7. יהיו  $\mathcal{H}_1, \mathcal{H}_2$  מחלקות של קלסיפקציה בינארית כך ש- $\mathcal{H}_1 \subseteq \mathcal{H}_2$ . צ.להראות ש-

$$\text{VC}(\mathcal{H}_1) \leq \text{VC}(\mathcal{H}_2)$$

תזכורות:

- תהא  $C \subseteq \mathcal{X}$  הצמצום של  $\mathcal{H}$  ל- $C$  הוא:

$$\mathcal{H}_C = \left\{ h_C \in \{0,1\}^C \mid \exists h \in \mathcal{H} \text{ s.t. } \forall x \in C (h_C(x) = h(x)) \right\}$$



• נאמר ש- $\mathcal{H}$  מנתצת קבוצה סופית  $C \subseteq \mathcal{X}$  אם  $|\mathcal{H}_C| = 2^{|C|}$ .

• הגדרנו:

$$\text{VC dim}(\mathcal{H}) = \sup \{m \in \mathbb{N} \mid \exists C \subseteq \mathcal{X} (|C| = m \text{ and } \mathcal{H} \text{ shatters } C)\}$$

• כדי להראות ש- $\text{VC dim}(\mathcal{H}) = d$  צריך להראות ש-

- יש  $C \subseteq \mathcal{X}$  מגודל  $d$  שמנותצת על ידי  $\mathcal{H}$ , קרי  $\text{VC dim}(\mathcal{H}) \geq d$

- כל  $C \subseteq \mathcal{X}$  מגודל  $d+1$  אינה מנותצת על ידי  $\mathcal{H}$ , קרי  $\text{VC dim}(\mathcal{H}) \leq d$

הוכחת הטענה:

$$\text{VC dim}(\mathcal{H}_1) = \sup \{m \in \mathbb{N} \mid \exists C \subseteq \mathcal{X} (|C| = m \text{ and } \mathcal{H}_1 \text{ shatters } C)\} =$$

$$= \sup \{m \in \mathbb{N} \mid \exists C \subseteq \mathcal{X} (|C| = m \text{ and } |(\mathcal{H}_1)_C| = 2^m)\} =$$

$$= \sup \left\{ m \in \mathbb{N} \mid \exists C \subseteq \mathcal{X} \left( |C| = m \text{ and } \left| \left\{ h_C \in \{0,1\}^C \mid \exists h \in \mathcal{H}_1 \text{ s.t. } \forall x \in C (h_C(x) = h(x)) \right\} \right| = 2^m \right) \right\}$$

ומאחר ש- $\mathcal{H}_1 \subseteq \mathcal{H}_2$ :

$$\leq \sup \left\{ m \in \mathbb{N} \mid \exists C \subseteq \mathcal{X} \left( |C| = m \text{ and } \left| \left\{ h_C \in \{0,1\}^C \mid \exists h \in \mathcal{H}_2 \text{ s.t. } \forall x \in C (h_C(x) = h(x)) \right\} \right| = 2^m \right) \right\} =$$

$$= \sup \{m \in \mathbb{N} \mid \exists C \subseteq \mathcal{X} (|C| = m \text{ and } |(\mathcal{H}_2)_C| = 2^m)\} =$$

$$= \sup \{m \in \mathbb{N} \mid \exists C \subseteq \mathcal{X} (|C| = m \text{ and } \mathcal{H}_2 \text{ shatters } C)\} =$$

$$= \text{VC dim}(\mathcal{H}_2)$$

כנדרש.

## 5 טענה תיאורתית

8. יהא  $\mathcal{X}$  מרחב דגימות ו- $\mathcal{Y} = \{\pm 1\}$ .

תהא  $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$  מחלקת היפותיזות.

עבור  $C \subseteq \mathcal{X}$ , ניזכר בסימון  $\mathcal{H}_C$  עבור הצמצום של  $\mathcal{H}$  לקבוצה  $C$ .

נגדיר את הפונקציה  $\tau_m(\mathcal{H}) : \mathbb{N} \rightarrow \mathbb{N}$  שמתאימה ל- $\mathcal{H}$  על ידי:

$$\tau_{\mathcal{H}}(m) := \max \{|\mathcal{H}_C| \mid C \subseteq \mathcal{X} \text{ and } |C| = m\}$$

א. צלהסביר מהי המשמעות של הפונקציה  $\tau_{\mathcal{H}}$ . הסבר:

ראשית, עבור  $|C| = m$  תמיד מתקיים ש-  $|\mathcal{H}_C| \leq 2^m$ , כי מספר התיגים המקסימלי הוא  $2^m$ , קרי 0 או 1 לכל איבר ב-  $C$ . ולכן  $\tau_{\mathcal{H}}(m) \leq 2^m$  לכל  $m$ .

עתה, אם  $\tau_{\mathcal{H}}(m) = 2^m$  אז  $\text{VC dim}(\mathcal{H}) \geq m$  כי  $\tau_{\mathcal{H}}(m) = 2^m$  גורר שקיימת קבוצה מגודל  $m$  שניתנת לניתוח. כמו כן, אם  $\tau_{\mathcal{H}}(m) < 2^m$  אז  $\text{VC dim}(\mathcal{H}) < m$  כי זה שלא קיימת קבוצה מגודל  $m$  שניתנת לניתוח.

אם כן, במידה ש-  $\text{VC dim}(\mathcal{H}) < m$ , הפונקציה  $\tau_{\mathcal{H}}(m)$  יכולה לתת מושג עד כמה  $m$  רחוק מלהיות מימד ה-  $\text{VC}$  של  $\mathcal{H}$ . אם ההפרש בין  $\tau_{\mathcal{H}}(m)$  הוא "גדול" אז מימד ה-  $\text{VC}$  של  $\mathcal{H}$  "הרבה" יותר קטן מ-  $m$ , ואם ההפרש קטן אז מימד ה-  $\text{VC}$  של  $\mathcal{H}$  הוא "כמעט"  $m$ . ובמילים אחרות המספר  $\tau_{\mathcal{H}}(m)$  נותן את מספר ההיפותיזות המקסימלי האפשרי עבור תת-קבוצה מגודל  $m$  של קבוצת המדגם.

ב. נניח ש-  $\text{VC dim}(\mathcal{H}) = \infty$ . צלמצוא ביטוי עבור הערך של  $\tau_{\mathcal{H}}(m)$  עבור  $m \in \mathbb{N}$ . פיתרון:

המשמעות של  $\text{VC dim}(\mathcal{H}) = \infty$  היא ש-  $\mathcal{H}$  מנתצת כל קבוצה מגודל  $m$  לכל  $m$  טבעי ולכן, בהסתמך על הסעיף הקודם, לכל  $m$  מתקיים ש-  $\tau_{\mathcal{H}}(m) = 2^m$ .

ג. נניח ש-  $\text{VC dim}(\mathcal{H}) = d$ . צלמצוא ביטוי עבור הערך של  $\tau_{\mathcal{H}}(m)$  עבור  $m \leq d$ . פיתרון:

$d$  הוא הגודל המקסימלי של קבוצה ש-  $\mathcal{H}$  מנתצת, אז מכך ש-  $\text{VC dim}(\mathcal{H}) = d$  נובע שקיימת קבוצה  $C \subseteq \mathcal{X}$  עם  $|C| = d$  כך ש-  $|\mathcal{H}_C| = 2^d$ . ואז אם  $m \leq d$  אזי שכל קבוצה  $C'$  עם  $|C'| = m$  ו-  $C' \subseteq C$ , מתקיים ש-  $|\mathcal{H}_{C'}| = 2^m$  שהרי:

$$\mathcal{H}_C = \left\{ h_C \in \{0, 1\}^C \mid \exists h \in \mathcal{H} \text{ s.t. } \forall x \in C (h_C(x) = h(x)) \right\}$$

קרי אם יש  $h \in \mathcal{H}$  כך ש- לכל  $x \in C$  מתקיים ש-  $h_C(x) = h(x)$ , אז בפרט לכל  $x \in C' \subseteq C$   $h_{C'}(x) = h(x)$  ולכן אם מתקבלת כל קלסיפיקציה אפשרית עבור  $C$  אז בהכרח גם עבור  $C'$  ומכאן ש-

$$\tau_{\mathcal{H}}(m) = 2^m$$

מש"ל

ד. נניח ש-  $\text{VC dim}(\mathcal{H}) = d$ . יהא  $m > d$ . נראה ש-

$$\tau_{\mathcal{H}}(m) \leq \left( \frac{em}{d} \right)^d$$

עבור  $e$  בסיס הלוגריתם הטבעי. נעשה זאת בשלבים הבאים:

i. צלהראות באינדוקציה שלכל  $C \subseteq \mathcal{X}$  סופית, מתקיים כי:

$$|\mathcal{H}_C| \leq |\{B \subseteq C \mid \mathcal{H} \text{ shatters } B\}|$$

הוכחה:

באינדוקציה על  $|C| = m$ :

בסיס האינדוקציה: עבור  $m = 1$ : יש איבר אחד ב-  $C$ .

אם ניתן לתייג אותו רק באפס או רק באחד אז  $|\mathcal{H}_C| = 1$  ואז הקבוצה היחידה שמוכלת ב-  $C$  ושמונתצת היא הקבוצה הריקה ולכן גם צד ימין שווה לאחד. אם ניתן לתייג את האיבר ב-  $C$  גם באפס וגם באחד אז צד שמאל בא"ש שווה לשתיים ואז ניתן לנתץ גם את  $C$  עצמה וגם את הקבוצה הריקה ולכן גם צד ימין הוא 2. מש"ל בסיס האינדוקציה.

צעד: נניח שהטענה נכונה לכל  $k < m$  עבור  $m > 1$  ונוכיח עבור  $m$ : נמספר את האיברים ב- $C$ :  $C = \{c_i\}_1^m$  וכן נגדיר  $C' = \{c_i\}_2^m$ .

כמו כן נגדיר:

$$Y_0 = \{(y_2, \dots, y_m) \mid (0, y_2, \dots, y_m) \in \mathcal{H}_C \vee (1, y_2, \dots, y_m) \in \mathcal{H}_C\}$$

$$Y_1 = \{(y_2, \dots, y_m) \mid (0, y_2, \dots, y_m) \in \mathcal{H}_C \wedge (1, y_2, \dots, y_m) \in \mathcal{H}_C\}$$

נשים לב שמתקיים כי  $Y_1 \subseteq Y_0$  כי לכל  $(y_2, \dots, y_m)$  עם  $(0, y_2, \dots, y_m) \in \mathcal{H}_C \wedge (1, y_2, \dots, y_m) \in \mathcal{H}_C$  בפרט מתקיים ש- $(0, y_2, \dots, y_m) \in \mathcal{H}_C \vee (1, y_2, \dots, y_m) \in \mathcal{H}_C$ . כמו כן, לכל  $(y_2, \dots, y_m) \in Y_1$  מתקיים ש- $(0, y_2, \dots, y_m) \in \mathcal{H}_C$  וגם  $(1, y_2, \dots, y_m) \in \mathcal{H}_C$  ולכל  $(y_2, \dots, y_m) \in Y_0 \setminus Y_1$  מתקיים ש- $(0, y_2, \dots, y_m) \in \mathcal{H}_C$  XOR  $(1, y_2, \dots, y_m) \in \mathcal{H}_C$  ולכן  $|\mathcal{H}_C| = 2|Y_1| + |Y_0 \setminus Y_1| = |Y_0| + |Y_1|$ .

עבור היפותיזה  $h_C \in \mathcal{H}_C$  שונתנת את התווית  $y_i$  עבור  $c_i$ , נסמן:  $h_C = (y_1, y_2, \dots, y_m)$  ובאופן דומה נסמן גם עבור  $h_{C'} \in \mathcal{H}_{C'}$  אם כן, על פי ההגדרה מתקיים:

$$\begin{aligned} \mathcal{H}_{C'} &= \left\{ h_{C'} \in \{0, 1\}^{C'} \mid \exists h_C \in \mathcal{H}_C \text{ s.t. } \forall x \in C' (h_{C'}(x) = h(x)) \right\} = \\ &= \left\{ (y_2, \dots, y_m) \in \{0, 1\}^{C'} \mid (0, y_2, \dots, y_m) \in \mathcal{H}_C \vee (1, y_2, \dots, y_m) \in \mathcal{H}_C \right\} = Y_0 \end{aligned}$$

אז מהנחת האינדוקציה:

$$|Y_0| = |\mathcal{H}_{C'}| \leq |\{B \subseteq C' \mid \mathcal{H} \text{ shutters } B\}|$$

אבל  $C' = C \setminus \{c_1\}$  ולכן:

$$= |\{B \subseteq C \mid \mathcal{H} \text{ shutters } B \text{ and } c_1 \notin B\}|$$

נגדיר  $\mathcal{H}' \subseteq \mathcal{H}$  על ידי:

$$\mathcal{H}' = \{h \in \mathcal{H} \mid \exists h' \in \mathcal{H} \text{ s.t. } (1 - h'(c_1), h'(c_2), \dots, h'(c_m)) = (h(c_1), h(c_2), \dots, h(c_m))\}$$

ונשים לב ש-

$$\begin{aligned} \mathcal{H}'_{C'} &= \left\{ h_{C'} \in \{0, 1\}^{C'} \mid \exists h' \in \mathcal{H}' \text{ s.t. } \forall x \in C' (h'_{C'}(x) = h'(x)) \right\} = \\ &= \left\{ h_{C'} \in \{0, 1\}^{C'} \mid \exists h \in \mathcal{H}' \text{ s.t. } \exists h' \in \mathcal{H} \text{ s.t. } (1 - h'(c_1), h'(c_2), \dots, h'(c_m)) = (h(c_1), h(c_2), \dots, h(c_m)) \text{ and } \forall x \in C' (h'_{C'}(x) = h(x)) \right\} = \\ &= \left\{ (y_2, \dots, y_m) \in \{0, 1\}^{C'} \mid (\overline{y_1}, y_2, \dots, y_m) \in \mathcal{H}_C \wedge (y_1, y_2, \dots, y_m) \in \mathcal{H}_C \right\} = Y_1 \end{aligned}$$

אז מהנחת האינדוקציה נקבל ש-

$$|Y_1| = |\mathcal{H}'_{C'}| \leq |\{B \subseteq C' \mid \mathcal{H}' \text{ shutters } B\}|$$

אבל לכל  $(y_2, \dots, y_m) \in \mathcal{H}'_{C'}$  מתקיים ש-  $(\overline{y_1}, y_2, \dots, y_m) \in \mathcal{H}_C \wedge (y_1, y_2, \dots, y_m) \in \mathcal{H}_C$  ולכן אם  $\mathcal{H}'$  מנתצת את  $B$  אז היא מנתצת גם את  $B \cup \{c_1\}$ , קרי:

$$= |\{B \subseteq C' \mid \mathcal{H}' \text{ shutters } B \cup \{c_1\}\}| =$$

$$= |\{B \subseteq C' \mid \mathcal{H}' \text{ shutters } B \text{ and } c_1 \in B\}|$$

ומאחר ש-  $\mathcal{H}' \subseteq \mathcal{H}$ :

$$\leq |\{B \subseteq C' \mid \mathcal{H} \text{ shutters } B \text{ and } c_1 \in B\}|$$

ולבסוף:

$$|\mathcal{H}_C| = |Y_0| + |Y_1| \leq$$

$$\leq |\{B \subseteq C \mid \mathcal{H} \text{ shutters } B \text{ and } c_1 \notin B\}| + |\{B \subseteq C' \mid \mathcal{H} \text{ shutters } B \text{ and } c_1 \in B\}| =$$

$$= |\{B \subseteq C \mid \mathcal{H} \text{ shutters } B\}|$$

■ כנדרש.

ii. צ.להסביר את המשמעות של אי-השוויון הנ"ל. הסבר:

אם יש ב-  $\mathcal{H}_C$   $m$  היפותיזות אז יש ב-  $C$  לפחות  $m$  תת-קבוצות ש-  $\mathcal{H}$  מנתצת אותן.

iii. צ.להראות שלכל קבוצה סופית  $C \subseteq \mathcal{X}$ , מתקיים:

$$|\{B \subseteq C \mid \mathcal{H} \text{ shutters } B\}| \leq \sum_{k=0}^d \binom{m}{k}$$

הוכחה:

נתון ש-  $\text{VCdim}(\mathcal{H}) = d$  ולכן  $\mathcal{H}$  יכולה לנתץ רק קבוצות שהן מגודל  $d$  ומטה. אז כל הקבוצות ב-  $C$  ש-  $\mathcal{H}$  יכולה לנתץ הן תת-קבוצות מגודל שהוא קטן או שווה ל-  $d$ . אז לכל  $k \leq d$ , יש  $\binom{m}{k}$  תת-קבוצות ב-  $C$  שהן מגודל  $k$ , ולכן מתקיים הא"ש הנ"ל.

iv. צ.להשתמש באי-השוויון הבא:

$$\sum_{k=0}^d \binom{m}{k} \leq \left(\frac{em}{d}\right)^d$$

כדי לסיים את ההוכחה ש-

$$\tau_{\mathcal{H}}(m) \leq \left(\frac{em}{d}\right)^d$$

הוכחה:

$$\tau_{\mathcal{H}}(m) := \max \{ |\mathcal{H}_C| \mid C \subseteq \mathcal{X} \text{ and } |C| = m \}$$

ומאחר שהטענות שהוכחנו לעיל נכונות לכל  $C$  מגודל  $m$ :

$$\leq \max \{ |\{B \subseteq C \mid \mathcal{H} \text{ shatters } B\}| \mid C \subseteq \mathcal{X} \text{ and } |C| = m \} \leq$$

$$\leq \sum_{k=0}^d \binom{m}{k}$$

ומאחר ש-  $m > d + 1$ :

$$\leq \left(\frac{em}{d}\right)^d$$

■ כנדרש.

ה. נניח ש-  $m = d$ . צלבדוק האם מתקיים הא"ש  $\tau_{\mathcal{H}}(m) \leq \left(\frac{em}{d}\right)^d$ . אם כן, האם הוא הדוק? פיתרון:

הא"ש נשמר כי:

$$\tau_{\mathcal{H}}(m) = \tau_{\mathcal{H}}(d) = 2^d \leq e^d = \left(\frac{ed}{d}\right)^d = \left(\frac{em}{d}\right)^d$$

וככל ש-  $d$  גדול יותר, כך החסם יותר הדוק כי ככל ש-  $d$  גדל כך קטן ההפרש  $e^d - 2^d = (e - 2)^d \approx (0.718)^d$ .

ו. צ.לאפיין במילים את ההתנהגות של  $\tau_{\mathcal{H}}(m)$  עבור  $m \leq \text{VC dim}(\mathcal{H})$  ועבור  $m > \text{VC dim}(\mathcal{H})$ , ולהשתמש בו כדי להציע הגדרה אלטרנטיבית עבור  $\text{VC dim}(\mathcal{H})$ . פיתרון:

עבור  $m > \text{VC dim}(\mathcal{H})$  מתקיים ש-  $\tau_{\mathcal{H}}(m) < 2^m$ , קרי מתקבלת התנהגות פולינומיאלית ועבור  $m \geq \text{VC dim}(\mathcal{H})$  מתקיים  $\tau_{\mathcal{H}}(m) = 2^m$  קרי מתקבלת התנהגות אקספוננציאלית.

מכאן נובע שאם גרף הפונקציה של  $\tau_{\mathcal{H}}(m)$  הוא אקספוננציאלי ב-  $m$  תמיד, אז  $\mathcal{H}$  אינה למידה-PAC כי אז  $\text{VC dim}(\mathcal{H}) = \infty$ , ואילו אם הגרף של  $\tau_{\mathcal{H}}(m)$  הינו פולינומי החל ממקום מסוים אזי ש-  $\mathcal{H}$  הינה למידה-PAC כי מתקיים ש-  $\text{VC dim}(\mathcal{H}) < \infty$ .

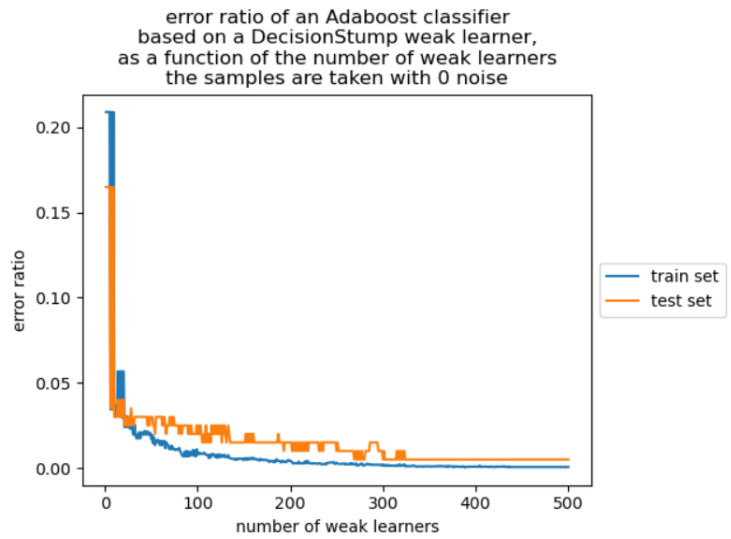
אז נוכל להגדיר:

$$\text{VC dim}(\mathcal{H}) = \min \{ m \in \mathbb{N} \mid \tau_{\mathcal{H}}(m) < 2^m \}$$

ואז אם נסמן ב-  $\widehat{\text{VC dim}}(\mathcal{H})$  את מימד ה-VC לפי ההגדרה החדשה, אז נקבל ש-  $\widehat{\text{VC dim}}(\mathcal{H}) = \text{VC dim}(\mathcal{H}) + 1$  כאשר  $\text{VC dim}(\mathcal{H}) < \infty$  והוא מימד ה-VC לפי ההגדרה המקורית:  $\widehat{\text{VC dim}}(\mathcal{H}) = d$  אזי שכל קבוצה מגודל  $d$  ומעלה לא ניתנת לניתוח ואילו  $d - 1$  ניתנת לניתוח ולכן  $\text{VC dim}(\mathcal{H}) = d$ . מש"ל

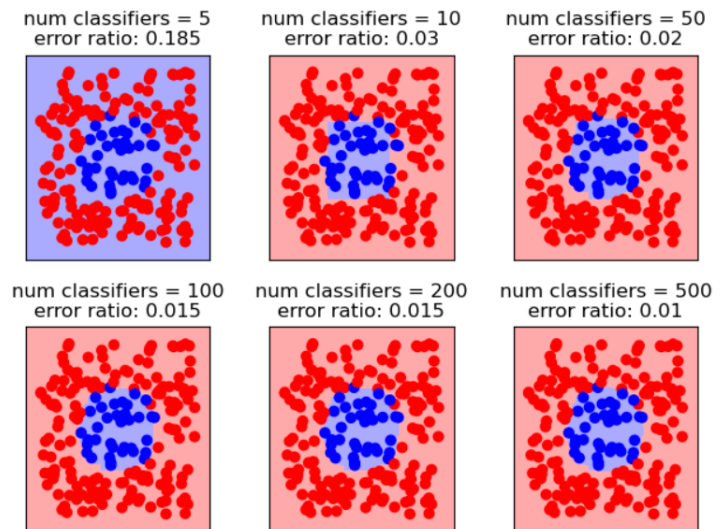
## 6 להפריד את הבלתי ניתן להפרדה - Adaboost

.10



.11

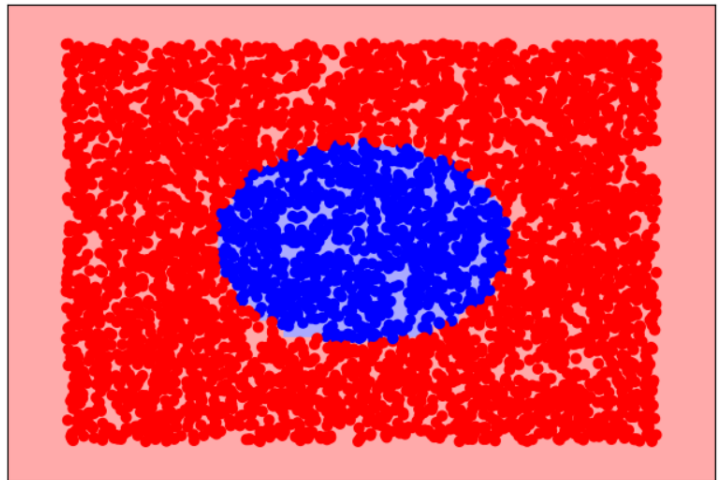
An Adaboost classifier based on a DecisionStump weak learner  
using different amounts of weak learners  
over the test data which consists of 200 samples  
with zero noise



.12. בשאלה הקודמת ניתן לראות ש-  $T = 500$  ממזער את השגיאה עבור ה- test set.

An Adaboost classifier based on a DecisionStump weak learner  
over the training data which consists of 5,000 samples  
with zero noise

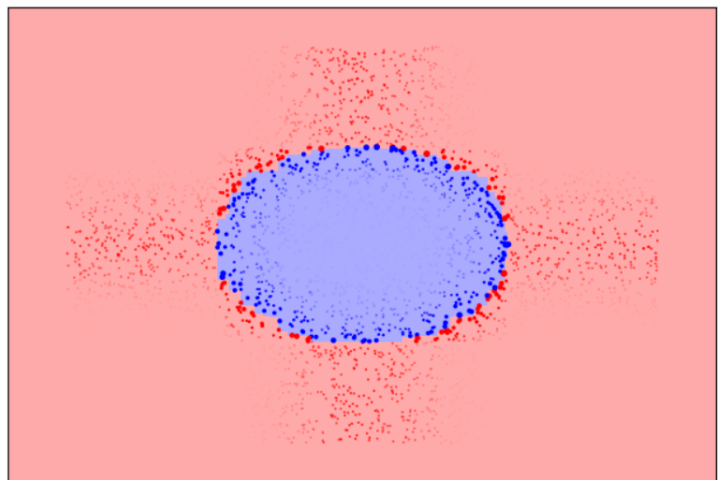
num classifiers = 500  
error ratio = 0.0



.13

The training set of an Adaboost classifier  
with size proportional to the transpose of the weights  
of the samples in the last iteration of the training  
with 5,000 training samples  
with zero noise

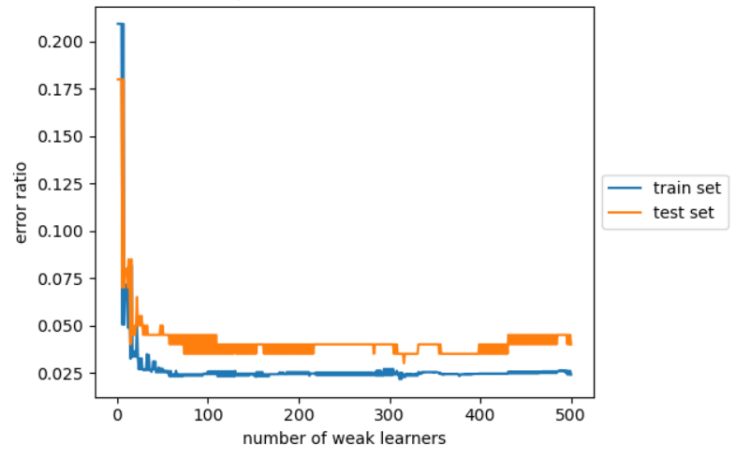
num classifiers = 500  
error ratio = 0.0002



בשוליים של העיגול הנקודות גדולות יותר כי האלגוריתם טעה עליהן הרבה פעמים במהלך תהליך הלמידה ולכן הוא נתן להן משקל גדול יותר מלנקודות באזורים אחרים. וזאת משום שקשה לסווג את השוליים באמצעות חלוקות אופקיות ואנכיות בלבד. כמו כן, הנקודות מחוץ לעיגול קטנות יותר כי האלגוריתם טעה עליהן פחות באופן יחסי, ולכן נתן להן משקל קטן יותר. וזאת משום שקל יותר לסווג אותן על ידי חלוקות אופקיות ואנכיות.

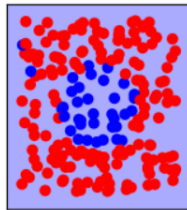
.14

error ratio of an Adaboost classifier  
based on a DecisionStump weak learner,  
as a function of the number of weak learners  
the samples are taken with 0.01 noise

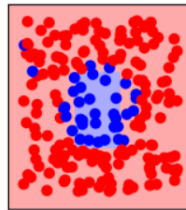


An Adaboost classifier based on a DecisionStump weak learner  
using different amounts of weak learners  
over the test data which consists of 200 samples  
with 0.01 noise

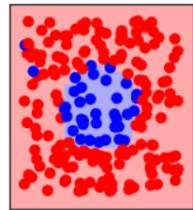
num classifiers = 5  
error ratio = 0.18



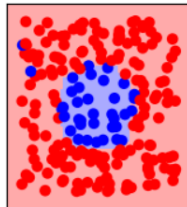
num classifiers = 10  
error ratio = 0.08



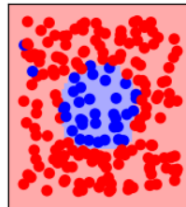
num classifiers = 50  
error ratio = 0.05



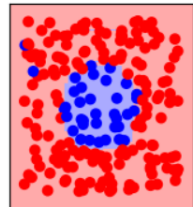
num classifiers = 100  
error ratio = 0.035



num classifiers = 200  
error ratio = 0.035



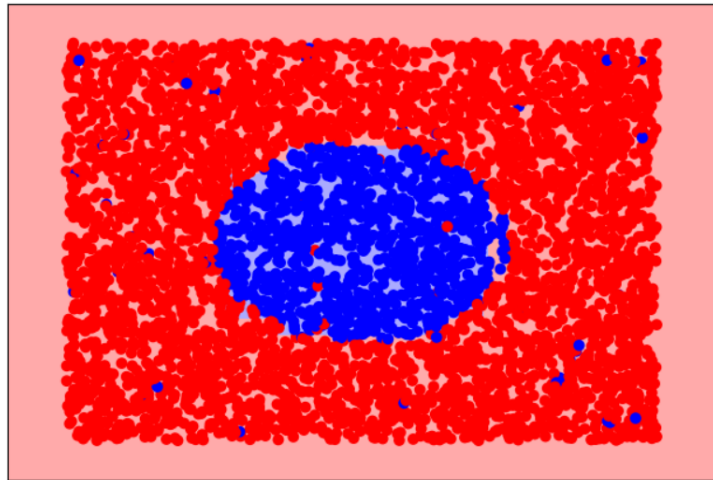
num classifiers = 500  
error ratio = 0.04





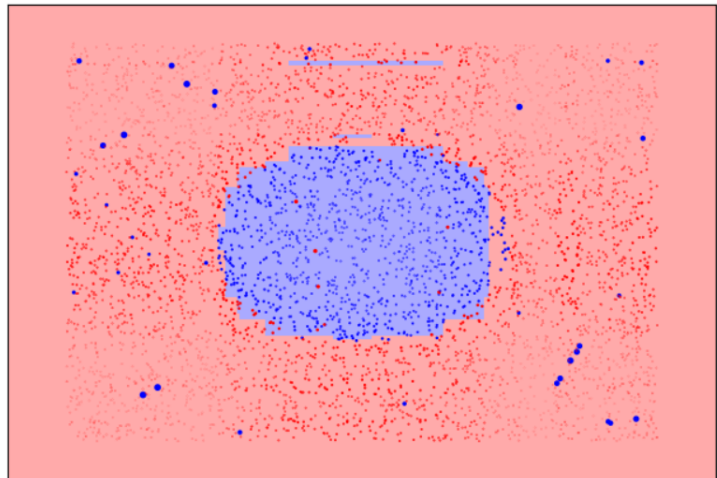
An Adaboost classifier based on a DecisionStump weak learner  
over the training data which consists of 5,000 samples  
with 0.01 noise

num classifiers = 100  
error ratio = 0.0256

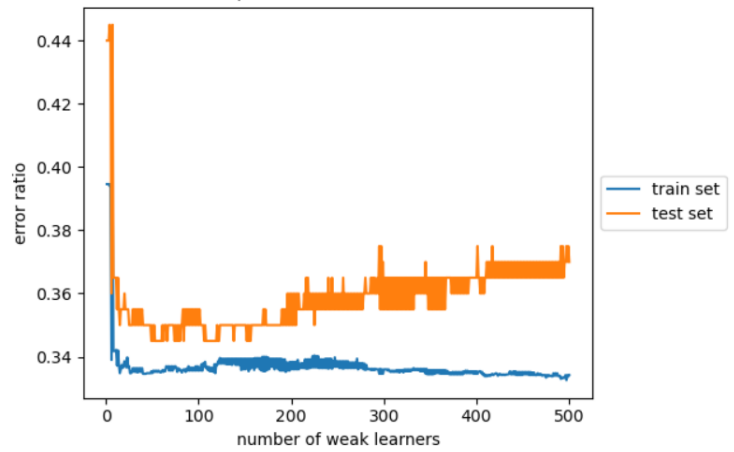


The training set of an Adaboost classifier  
with size proportional to the transpose of the weights  
of the samples in the last iteration of the training  
with 5,000 training samples  
with 0.01 noise

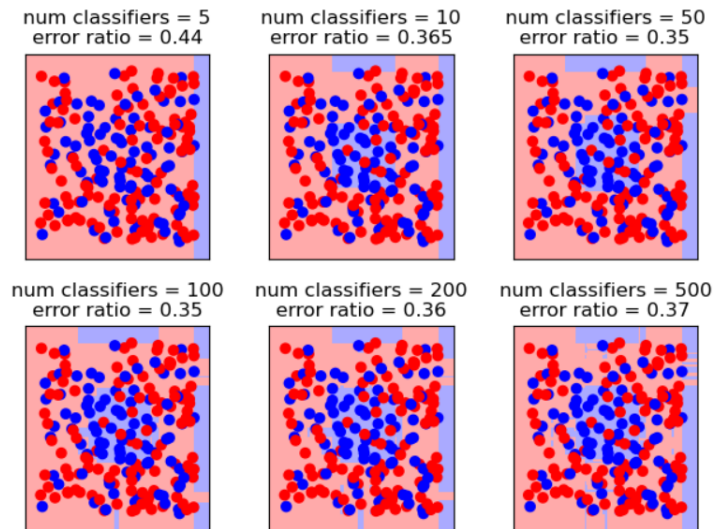
num classifiers = 500  
error ratio = 0.024



error ratio of an Adaboost classifier  
based on a DecisionStump weak learner,  
as a function of the number of weak learners  
the samples are taken with 0.4 noise

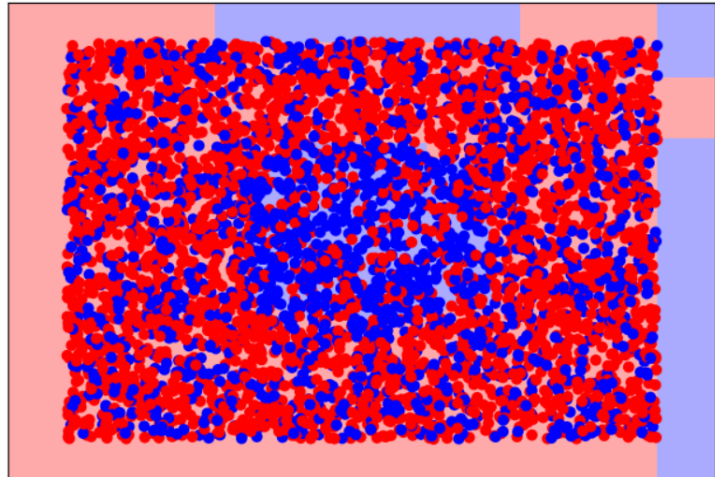


An Adaboost classifier based on a DecisionStump weak learner  
using different amounts of weak learners  
over the test data which consists of 200 samples  
with 0.4 noise



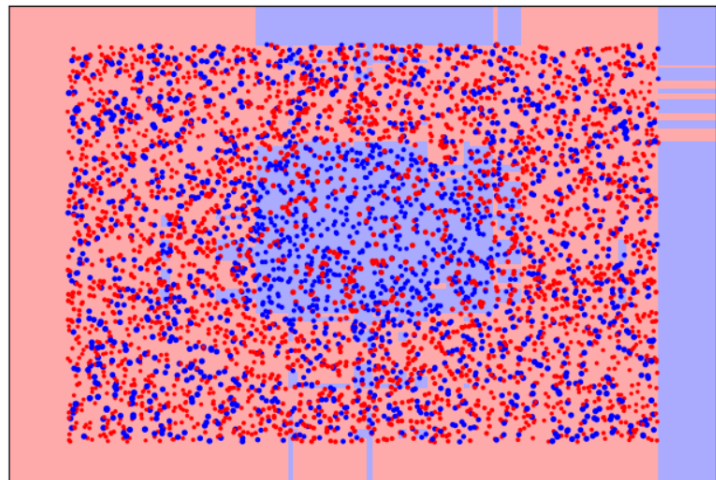
An Adaboost classifier based on a DecisionStump weak learner  
over the training data which consists of 5,000 samples  
with 0.4 noise

num classifiers = 50  
error ratio = 0.3348



The training set of an Adaboost classifier  
with size proportional to the transpose of the weights  
of the samples in the last iteration of the training  
with 5,000 training samples  
with 0.4 noise

num classifiers = 500  
error ratio = 0.3342



#### הסבר כללי לשינויים:

כשהרעש יחסית קטן, קרי 0.01, ניתן לראות שהשגיאה בסיווג היא יחסית קטנה ומספר המסווגים האופטימלי הוא עדיין המקסימלי מבין האפשרויות, וכשהרעש יחסית גדול, קרי 0.4, ניתן לראות שהשגיאה גדלה משמעותית, וכך גם מספר המסווגים האופטימלי מבין האפשרויות, נעשה קטן יותר באופן משמעותי.

#### הסבר של השינויים ביחס לשאלה 10:

ככל שהדאטא נתון עם יותר רעש, כך bias קטן יותר, קרי שימוש במספר גדול של לומדים חלשים, גורם להתאמת-יתר חמורה יותר, קרי variance גדול יותר. לכן ככל שהרעש גדל כך ניתן לראות שהשגיאה קבוצת המבחן גדלה, והפער בין מידת השגיאה בקבוצת האימון ובין השגיאה בקבוצת המבחן, גדל.

#### הסבר של השינויים ביחס לשאלה 12:

בהסתמך על ההסבר שנתתי לעיל עבור השינויים ביחס לשאלה 10, נובע שנרצה להשתמש בפחות לומדים חלשים על מנת להקטין את ה- variance, ואכן כשהרעש יחסית גדול, קרי 0.4, נקבל שגיאה מינימלית עבור 50 מסווגים ולא עבור 500 כמו במקרה ללא הרעש.

## 2 adaboost.py

```
1  """
2  =====
3      Introduction to Machine Learning (67577)
4  =====
5
6  Skeleton for the AdaBoost classifier.
7
8  Author: Gad Zalcberg
9  Date: February, 2019
10
11  """
12  from ex4_tools import *
13
14
15  TRAIN_SAMPLES_AMOUNT = 5000
16  TEST_SAMPLES_AMOUNT = 200
17  DEFAULT_NOISE = 0
18  Q14_NOISES = [0.01, 0.4]
19  DEFAULT_CLASSIFIERS_AMOUNT = 500
20  CLASSIFIERS_AMOUNTS_LST = [5, 10, 50, 100, 200, 500]
21
22
23  class AdaBoost(object):
24
25      def __init__(self, WL, T):
26          """
27          Parameters
28          -----
29          WL : the class of the base weak learner
30          T : the number of base learners to learn
31          """
32          self.WL = WL
33          self.T = T
34          self.h = [None]*T # list of base learners
35          self.w = np.zeros(T) # weights
36
37      def train(self, X, y):
38          """
39          Parameters
40          -----
41          X : samples, shape=(num_samples, num_features)
42          y : labels, shape=(num_samples)
43          Train this classifier over the sample (X,y)
44          After finish the training return the weights of the samples in the
45          last iteration.
46          """
47          m = len(X)
48          D = np.ones(m) / m
49          for i in range(self.T):
50              h = self.WL(D, X, y)
51              h.train(D, X, y)
52              y_pred = h.predict(X)
53              self.h[i] = h
54              zero_indices = np.where(y_pred == 0)[0]
55              y_pred[zero_indices] = -1
56              compare = (y != y_pred).astype(int)
57              epsilon = D @ compare
58              epsilon_inverse = 1 / epsilon
59              in_log = epsilon_inverse - 1
```

```

60         curr_log = np.log(in_log)
61         w = 0.5 * curr_log
62         self.w[i] = w
63         y_and_y_pred_product = np.multiply(y, y_pred)
64         in_exp = -w * y_and_y_pred_product
65         exp_arr = np.exp(in_exp)
66         D = np.multiply(D, exp_arr)
67         D /= np.sum(D)
68     return D
69
70     def predict(self, X, max_t):
71         """
72         Parameters
73         -----
74         X : samples, shape=(num_samples, num_features)
75         :param max_t: integer < self.T: the number of classifiers to use for
76         the classification
77         :return: y_hat : a prediction vector for X. shape=(num_samples)
78         Predict only with max_t weak learners,
79         """
80         lst = [self.w[i] * self.h[i].predict(X) for i in range(max_t)]
81         sign = np.sign(np.sum(lst, axis=0))
82         return sign
83
84     def error(self, X, y, max_t):
85         """
86         Parameters
87         -----
88         X : samples, shape=(num_samples, num_features)
89         y : labels, shape=(num_samples)
90         :param max_t: integer < self.T: the number of classifiers to use for
91         the classification
92         :return: error : the ratio of the wrong predictions when predict only
93         with max_t weak learners (float)
94         """
95         pred = self.predict(X, max_t)
96         return np.sum(pred != y) / len(y)
97
98
99     class ExeSolver:
100     def __init__(self, noise):
101         self.noise = noise
102         self.train_samples = generate_data(TRAIN_SAMPLES_AMOUNT, noise)
103         self.X_train, self.y_train = self.train_samples[0], \
104                                     self.train_samples[1]
105         self.adaboost = AdaBoost(DecisionStump, DEFAULT_CLASSIFIERS_AMOUNT)
106         self.D = self.adaboost.train(self.X_train, self.y_train)
107         self.test_samples = generate_data(TEST_SAMPLES_AMOUNT, noise)
108         self.X_test = self.test_samples[0]
109         self.y_test = self.test_samples[1]
110
111
112     def q10(exe_solver):
113         T = range(1, DEFAULT_CLASSIFIERS_AMOUNT+1)
114         train_errors = \
115             [exe_solver.adaboost.error(
116                 exe_solver.X_train, exe_solver.y_train, i) for i in T]
117         test_errors = \
118             [exe_solver.adaboost.error(
119                 exe_solver.X_test, exe_solver.y_test, i) for i in T]
120         plt.plot(T, train_errors, label='train set')
121         plt.plot(T, test_errors, label='test set')
122         plt.legend(loc='center left', bbox_to_anchor=(1, 0.5))
123         plt.title('error ratio of an Adaboost classifier\n'
124                 'based on a DecisionStump weak learner,\n'
125                 'as a function of the number of weak learners\n'
126                 'the samples are taken with {} noise'.format(exe_solver.noise))
127         plt.xlabel('number of weak learners')

```

```

128     plt.ylabel('error ratio')
129     plt.show()
130
131
132     def q11(exe_solver):
133         counter = 1
134         errors_and_classifiers_amounts = []
135         for t in CLASSIFIERS_AMOUNTS_LST:
136             plt.subplot(2, 3, counter)
137             error = decision_boundaries(exe_solver.adaboost, exe_solver.X_test,
138                                       exe_solver.y_test, t)
139             errors_and_classifiers_amounts.append((error, t))
140             counter += 1
141         plt.show()
142         return min(errors_and_classifiers_amounts, key=lambda x: x[0])[1]
143
144
145     def q12(exe_solver):
146         amount = q11(exe_solver)
147         decision_boundaries(
148             exe_solver.adaboost, exe_solver.X_train, exe_solver.y_train, amount)
149         plt.show()
150
151
152     def q13(exe_solver):
153         new_D = (exe_solver.D / np.max(exe_solver.D) * 10).T
154         decision_boundaries(exe_solver.adaboost, exe_solver.X_train,
155                             exe_solver.y_train, DEFAULT_CLASSIFIERS_AMOUNT,
156                             weights=new_D)
157         plt.show()
158
159
160     def run_questions(noise):
161         exe_solver = ExeSolver(noise)
162         q10(exe_solver)
163         q12(exe_solver)
164         q13(exe_solver)
165
166
167     def run_zero_noise():
168         run_questions(0)
169
170
171     def q14():
172         for noise in Q14_NOISES:
173             run_questions(noise)
174
175
176     # run questions 10-13 with zero noise
177     # run_zero_noise()
178
179     # run question 14
180     # q14()

```

## 3 ex4 tools.py

```
1  """
2  =====
3      Introduction to Machine Learning (67577)
4  =====
5
6  This module provides some useful tools for Ex4.
7
8  Author: Gad Zalcberg
9  Date: February, 2019
10
11 """
12 import numpy as np
13 import matplotlib.pyplot as plt
14 from matplotlib.colors import ListedColormap
15 from itertools import product
16 from matplotlib.pyplot import imread
17 import os
18 from sklearn.model_selection import train_test_split
19
20
21 cm = ListedColormap(['#AAAAFF', '#FFAAAA'])
22 cm_bright = ListedColormap(['#0000FF', '#FF0000'])
23
24
25 def find_threshold(D, X, y, sign, j):
26     """
27     Finds the best threshold.
28     D = distribution
29     S = (X, y) the data
30     """
31     # sort the data so that x1 <= x2 <= ... <= xm
32     sort_idx = np.argsort(X[:, j])
33     X, y, D = X[sort_idx], y[sort_idx], D[sort_idx]
34
35     thetas = np.concatenate([[ -np.inf], (X[1:, j] + X[:-1, j]) / 2, [np.inf]])
36     minimal_theta_loss = np.sum(D[y == sign]) # loss of the smallest
37     # possible theta
38     losses = np.append(minimal_theta_loss, minimal_theta_loss -
39                        np.cumsum(D * (y * sign)))
40     min_loss_idx = np.argmin(losses)
41     return losses[min_loss_idx], thetas[min_loss_idx]
42
43
44 class DecisionStump(object):
45     """
46     Decision stump classifier for 2D samples
47     """
48
49     def __init__(self, D, X, y):
50         self.theta = 0
51         self.j = 0
52         self.sign = 0
53         self.train(D, X, y)
54
55     def train(self, D, X, y):
56         """
57         Train the classifier over the sample (X,y) w.r.t. the weights D over X
58         Parameters
59         -----
```



```

60     D : weights over the sample
61     X : samples, shape=(num_samples, num_features)
62     y : labels, shape=(num_samples)
63     """
64     loss_star, theta_star = np.inf, np.inf
65     for sign, j in product([-1, 1], range(X.shape[1])):
66         loss, theta = find_threshold(D, X, y, sign, j)
67         if loss < loss_star:
68             self.sign, self.theta, self.j = sign, theta, j
69         loss_star = loss
70
71     def predict(self, X):
72         """
73         Parameters
74         -----
75         X : shape=(num_samples, num_features)
76         Returns
77         -----
78         y_hat : a prediction vector for X shape=(num_samples)
79         """
80         y_hat = self.sign * ((X[:, self.j] <= self.theta) * 2 - 1)
81         return y_hat
82
83
84 def decision_boundaries(classifier, X, y, num_classifiers=1, weights=None):
85     """
86     Plot the decision boundaries of a binary classifiers over  $X \subset \mathbb{R}^2$ 
87
88     Parameters
89     -----
90     classifier : a binary classifier, implements classifier.predict(X)
91     X : samples, shape=(num_samples, 2)
92     y : labels, shape=(num_samples)
93     title_str : optional title
94     weights : weights for plotting X
95     """
96     cm = ListedColormap(['#AAAAFF', '#FFAAAA'])
97     cm_bright = ListedColormap(['#0000FF', '#FF0000'])
98     h = .003 # step size in the mesh
99     # Plot the decision boundary.
100     x_min, x_max = X[:, 0].min() - .2, X[:, 0].max() + .2
101     y_min, y_max = X[:, 1].min() - .2, X[:, 1].max() + .2
102     xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
103     Z = classifier.predict(np.c_[xx.ravel(), yy.ravel()], num_classifiers)
104     # Put the result into a color plot
105     Z = Z.reshape(xx.shape)
106     plt.pcolormesh(xx, yy, Z, cmap=cm)
107     # Plot also the training points
108     if weights is not None:
109         plt.scatter(X[:, 0], X[:, 1], c=y, s=weights, cmap=cm_bright)
110     else:
111         plt.scatter(X[:, 0], X[:, 1], c=y, cmap=cm_bright)
112     plt.xlim(xx.min(), xx.max())
113     plt.ylim(yy.min(), yy.max())
114     plt.xticks([])
115     plt.yticks([])
116     error = classifier.error(X, y, num_classifiers)
117     plt.title('num classifiers = {} \n error ratio = {}'.format(
118         num_classifiers, error))
119     plt.draw()
120     return error
121
122
123 def generate_data(num_samples, noise_ratio):
124     """
125     generate samples X with shape: (num_samples, 2) and labels y with shape (num_samples).
126     num_samples: the number of samples to generate
127     noise_ratio: invert the label for this ratio of the samples

```

```

128     '''
129     X = np.random.rand(num_samples, 2) * 2 - 1
130     radius = 0.5 ** 2
131     in_circle = np.sum(X ** 2, axis=1) < radius
132     y = np.ones(num_samples)
133     y[in_circle] = -1
134     y[np.random.choice(num_samples, int(noise_ratio * num_samples))] *= -1
135
136     return X, y

```