

## 2 Theoretical Part

## 2.1 PAC Learnability

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

Hints:

- For  $(a) \Rightarrow (b)$  show that  $\forall \varepsilon, \delta > 0$  and  $\forall m \geq m(\varepsilon, \delta)$  it holds that  $\mathbb{E}_{S \sim \mathcal{D}^m}[L_{\mathcal{D}}(A(S))] \leq \varepsilon + \delta$ .
  - For  $(b) \Rightarrow (a)$  use Markov's inequality

$$\forall \epsilon > 0 \quad \exists M \in \mathbb{N} \quad \forall m \geq M \quad |E_{S_m} (L_0(A(S)))| < \epsilon$$

- פ"ג נס ערך ג' רצ'ר

$m\left(\frac{k}{2}, \frac{k}{2}\right) = \rho''(a)$   $\rho''(a) < \varepsilon + \delta$   $\Rightarrow \varepsilon = \delta = \frac{k}{2} > 0$   $\text{and } sk > 0$

$$E_{S \sim \sigma^m} (L_b(A(S))) = E[L_b(A(S)) \leq \frac{k}{2}] + E[L_b(A(S)) > \frac{k}{2}] \leq \frac{k}{2} \cdot P_{S \sim \sigma^m}(L_b(A(S)) \leq \frac{k}{2}) +$$

$$1 \cdot P_{S \sim \Omega^m} (L_b(A(S)) > \frac{k}{2}) \leq \frac{k}{2} \cdot 1 + 1 \cdot \frac{k}{2} = k$$

הו תומך בראויים גנויים מטענה. נורווגיה נרואה כמי שפונה לנגד עיניו.

$$P_{S \sim \Omega^m}(L_b(A(S)) > \frac{k}{2}) < \frac{1}{2}$$

INCUBATOR ה-0 ק-ה פ-מ. (היכרוי ס-מ ה-מ)

$$P(|X| \geq a) \leq \frac{E[X]}{a}$$

אנו מודים לך נינה פון פולס.

. $\exists \delta > 0$  .(  $L_\rho(A(s)) \in [0, 1]$  )  $\Rightarrow$   $\exists \epsilon > 0$  . $L_\rho(A(s)) < \epsilon$  .  
 סעיף ב' בeweis

நாட்காலிகளில் முறையின் பொருளை விட்டு விடுவதே தமிழ்நாட்காலிகளின் பொருளாகும்.

$$E\left(\underset{s \sim 0^m}{L}_0(A(s))\right) = \left|E\left(\underset{s \sim 0^m}{L}_0(A(s))\right)\right| < \varepsilon \delta$$

$\downarrow$

$$L_0(A(s)) \geq 0$$

$$\rho(L_0(\alpha(s)) > \varepsilon) \leq \frac{E(L_0(\alpha(s)))}{\varepsilon} = \frac{\varepsilon d}{\varepsilon} = d$$

לפיכך, אם  $\mu$  מוגדר ככזה,  $P(L_0(A(s)) > \varepsilon) \geq 1 - \delta$  ומכאן

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.,

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

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

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

When proving, *do not* use a VC-Dimension argument. Instead, prove the claim directly from the PAC learnability definition by showing a specific algorithm and analyzing its sample complexity.

Hint: Remember that for every  $\varepsilon > 0$  it holds that  $1 - \varepsilon \leq e^{-\varepsilon}$

כזכור נזכר בה הדרישה גפונטיה

**Definition 1.1** An hypothesis class  $\mathcal{H}$  is *PAC learnable* if there exists a function  $m_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathbb{N}$  and a learning algorithm  $\mathcal{A}$  such that:

- For every  $\varepsilon, \delta \in (0, 1)$
  - For every distribution  $\mathcal{D}$  over  $\mathcal{X}$
  - For every labeling function  $f: \mathcal{X} \rightarrow \{0, 1\}$

if the realizable assumption holds with respect to  $\mathcal{H}, \mathcal{D}, f$ , then when running the learning algorithm  $\mathcal{A}$  on  $m \geq m_{\mathcal{H}}(\epsilon, \delta)$  i.i.d samples drawn from  $\mathcal{D}$  and labeled by  $f$ , the algorithm

returns a hypothesis  $h_S = \mathcal{A}(S)$  such that:

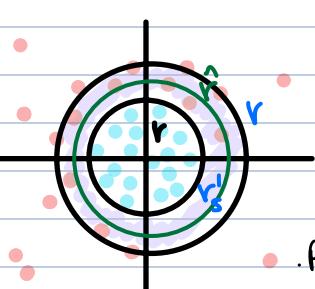
$$\mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S) \leq \varepsilon] \geq 1 - \delta$$

$$r_s = \max_{\substack{i: y_i = 1 \\ i \in [m]}} x_{ni}^2 + x_{zi}^2 = \|x\|_F h_{rs}$$

## רְבָעָה וְיָמִים

\*Cc, יתך רך לאנגלית, ניקון (קן) ניקון, אולין כויה.

כג' (טבנין) מ' (טבנין) ה' (טבנין)



הוכיחו ש- $\sum_{\|x_i\| \leq r} \|x_i\|^2$  מוגדרת כפונקציית נורמה על  $\mathbb{R}^n$ .

$$P_{\text{large } m} (L_0(h_s) > \varepsilon) = P(\{S | \forall x \in S \quad \|x\| \leq r \quad \vee \quad \|x\| > r\}) \leq (1 - \varepsilon)^m \leq (e^{-\varepsilon})^m = e^{-\varepsilon m}$$

$$m \geq \frac{\log \frac{1}{\delta}}{\varepsilon} \iff m \geq -\frac{\log \delta}{\varepsilon} \iff -\varepsilon m \leq \log \delta \iff e^{-\varepsilon m} \leq \delta$$

$$m_H(\varepsilon, \delta) = \frac{\log \frac{1}{\delta}}{\varepsilon} \quad \text{by PAC guarantee} \quad \text{defn of } H - \ell \text{ PONIIPC PC}$$

3. 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  $m_{\mathcal{H}}(\varepsilon, \delta) \leq m_{\mathcal{H}}^{UC}(\varepsilon/2, \delta)$ .

לעומת הוכחה

**Definition 4.6.4 — Uniform Convergence Property.** A hypothesis class  $\mathcal{H}$  is said to have the *uniform convergence property* if and only if there exists a function  $m_{\mathcal{H}}^{UC} : (0, 1)^2 \rightarrow \mathbb{N}$  such that for every  $\varepsilon, \delta \in (0, 1)$  and every distribution  $\mathcal{D}$  on  $\mathcal{X} \times \mathcal{Y}$

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

**Definition 4.6.3 —  $\varepsilon$ -representative.** A training sample  $S$  is called  $\varepsilon$ -representative for  $\mathcal{D}, \mathcal{H}, \ell$  if and only if

$$\forall h \in \mathcal{H} \quad |L_S(h) - L_{\mathcal{D}}(h)| < \varepsilon$$

**Definition 4.5.5 — Agnostic PAC Learnability.** A hypothesis class  $\mathcal{H}$  is Agnostic-PAC learnable with respect to loss  $\ell : \mathcal{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, \infty)$  if there exists a function  $\tilde{m}_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathbb{N}$  and a learning algorithm  $\mathcal{A} : (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathcal{H}$  with the following property:

- For any  $\varepsilon, \delta \in (0, 1)$
- For any distribution  $\mathcal{D}$  over  $\mathcal{X} \times \mathcal{Y}$

when running the learning algorithm  $\mathcal{A}$  on  $m \geq \tilde{m}_{\mathcal{H}}(\varepsilon, \delta)$  i.i.d samples generated by  $\mathcal{D}$ , the algorithm returns a hypothesis  $h_S = \mathcal{A}(S)$  such that, with probability of at least  $1 - \delta$ :

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

רנו  $\varepsilon, \delta > 0$  ו-  $\chi \times \chi$  ספנ  $\mathcal{D}$  מילבדת,  $\mathcal{A} \in \text{ERM}_{\mathcal{H}}$  נאנו רג'יג'ס נקי. נמי  $\tilde{m}_{\mathcal{H}}(\varepsilon, \delta)$  פון, Uniform Convergence של ערך ה-  $\ell$  ה-Fe NCII מיל  $\tilde{m}_{\mathcal{H}}(\varepsilon, \delta)$  גוף

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

: מילנו מילנו נס פה ווקט

ו-  $L_0(h_S) \leq \min_{h' \in \mathcal{H}} L_0(h') + \varepsilon$  - פון  $h_S \in \arg \min_{h \in \mathcal{H}} L_S(h)$  מיל,  $h_S \in \text{ERM}_{\mathcal{H}}(S)$

$$\begin{aligned} \frac{\varepsilon}{2}\text{-representative} & \quad \frac{\varepsilon}{2}\text{-representative} \\ L_0(h_S) \leq \frac{\varepsilon}{2} + L_S(h_S) & \leq \frac{\varepsilon}{2} + L_S(h) \leq \frac{\varepsilon}{2} + L_0(h) + \frac{\varepsilon}{2} = L_0(h) + \varepsilon \end{aligned} \quad : \text{מיל}$$

ו-  $L_0(h_S) \leq \frac{\varepsilon}{2} + L_S(h_S)$  ו-  $L_S(h) \leq \frac{\varepsilon}{2} + L_0(h)$  ו-  $L_0(h) \leq \frac{\varepsilon}{2} + L_0(h)$  ו-  $L_0(h) + \frac{\varepsilon}{2} = L_0(h) + \varepsilon$

:  $P$ " $\exists$   $m \geq m_H(\varepsilon, \delta)$   $\forall \delta$   $\exists \bar{\delta}$ ,  $m_H(\varepsilon, \delta) = m_H\left(\frac{\varepsilon}{2}, \bar{\delta}\right)$   $\exists \bar{m}$

$$\Omega^m \left( \{s \mid L_0(h_s) \leq \min_{h \in H} L_0(h) + \varepsilon \} \right) \geq \Omega^m \left( \{s \in (xxy)^m \mid s \text{ is } \frac{\varepsilon}{2}\text{-representative}\} \right) \geq 1 - \delta$$

$\downarrow$   
בנוסף, נניח  
הנחתה נסיעה נרחבת ורואה

## 2.2 VC-Dimension

פערת נס כריה

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

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

נונען פק 1 גורן פק 0  
פוק פק 0 גורן פק 1

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

Hint: what is the size of the hypothesis class?

- 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 VC-dimension( $\mathcal{H}_{\text{parity}}$ ) =  $n - e$  |  
180/1 (1)

ונדרה לנו  $\mathcal{H}_{\text{parity}}$  מוגדרת על ידי סדרת האפשרויות  $\{e_1, e_2, \dots, e_n\} \subseteq \{0, 1\}^n$  ופונקציית פולינומית  $h_I$  ביחס לסדרת הנקודות  $(y_1, y_2, \dots, y_n) \in \{0, 1\}^n$  מוגדרת כ  $h_I(e_i) = y_i \bmod 2 = y_i$ . כלומר  $e_i$  מוגדרת כ  $y_i$ .

- לכן סדרת הנקודות מוגדרת על ידי סדרת הנקודות  $\{e_1, e_2, \dots, e_n\}$ .

**Claim 2.1** Let  $\mathcal{H}$  be a finite hypothesis class then  $VC\text{-Dim}(\mathcal{H}) \leq \log_2 |\mathcal{H}|$ .

ונדרה לנו  $\mathcal{H}_{\text{parity}}$  מוגדרת על ידי סדרת האפשרויות  $\{e_1, e_2, \dots, e_n\} \subseteq \{0, 1\}^n$  ופונקציית פולינומית  $h_I$  ביחס לסדרת הנקודות  $(y_1, y_2, \dots, y_n) \in \{0, 1\}^n$  מוגדרת כ  $h_I(e_i) = y_i \bmod 2 = y_i$ . כלומר  $e_i$  מוגדרת כ  $y_i$ .

.VC-dimension( $\mathcal{H}_{\text{parity}}$ )  $\leq \log 2^n = n - e$  |  
פונקציית פולינומית מוגדרת על ידי סדרת הנקודות  $\{e_1, e_2, \dots, e_n\}$ .

.VC-dimension( $\mathcal{H}_{\text{parity}}$ ) =  $n$  |  
פונקציית פולינומית מוגדרת על ידי סדרת הנקודות  $\{e_1, e_2, \dots, e_n\}$ .

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

$c \subseteq \mathcal{X}$  פ"ג  $VC\text{-dim}(\mathcal{H}_2) < VC\text{-dim}(\mathcal{H}_1)$  פ"ג  $\mathcal{H}_1 \subseteq \mathcal{H}_2$  גורף הינה  $|c| > VC\text{-dim}(\mathcal{H}_2)$  -0 פ"ג  $|\mathcal{H}_1| = 2^c$   $\mathcal{H}_1 \subseteq \mathcal{H}_2$  פ"ג  $|\mathcal{H}_2| \geq |\mathcal{H}_1| = 2^c$   $\Leftarrow |\mathcal{H}_1| = 2^c$   $c \in \mathbb{N}$  פ"ג  $|\mathcal{H}_2| \geq |\mathcal{H}_1| = 2^c$   $c \in \mathbb{N}$

## 2.3 Regularization

## Based on Lecture 7 and Recitation 9

1. In the following question we will show that although the Ridge estimator is biased it can achieve lower MSE compared to the LS estimator.

Let  $\mathbf{X} \in \mathbb{R}^{m \times d}$  be a **constant** design matrix,  $\mathbf{y} \in \mathbb{R}^m$  a response vector, and assume that  $\mathbf{X}^\top \mathbf{X}$  is invertible. Denote  $\hat{\mathbf{w}}$  the LS solution and  $\hat{\mathbf{w}}_\lambda$  the ridge solution for the regularization parameter  $\lambda \geq 0$  (where  $\hat{\mathbf{w}}_0 \equiv \hat{\mathbf{w}}$ )

- Assume the linear model is correct, namely  $\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\varepsilon}$  where  $\varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$ .
  - Recall that in this case LS estimator is unbiased:  $\mathbb{E}[\hat{\mathbf{w}}] = \mathbf{w}$ .

- (a) Show that  $\hat{\mathbf{w}}_\lambda = A_\lambda \hat{\mathbf{w}}$  where  $A_\lambda := (\mathbf{X}^\top \mathbf{X} + \lambda I_d)^{-1}(\mathbf{X}^\top \mathbf{y})$

(b) From the above, conclude that for any  $\lambda > 0$  the ridge estimator is a biased estimator of  $\mathbf{w}$ . That is, show that for any  $\lambda > 0$   $\mathbb{E}[\hat{\mathbf{w}}_\lambda] \neq \mathbf{w}$ .

(c) Show that:  $\text{Var}(\hat{\mathbf{w}}_\lambda) = \sigma^2 A_\lambda (\mathbf{X}^\top \mathbf{X})^{-1} A_\lambda^\top$ , for  $\sigma^2$  the variance of the assumed noise.

*Hint::* Recall that for a constant matrix  $B$  and a random vector  $\mathbf{z}$  it holds that  $\text{Var}(B\mathbf{z}) = B \cdot \text{Var}(\mathbf{z}) \cdot B^\top$  and that  $\text{Var}(\hat{\mathbf{w}}) = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}$ .

(d) Derive explicit expressions for the (squared) bias and variance of  $\hat{\mathbf{w}}_\lambda$  as a function of  $\lambda$ , i.e. write a bias-variance decomposition for the mean square error of  $\hat{\mathbf{w}}_\lambda$ .

Hint: recall that for the multivariate case the MSE defined to be:

$$MSE(\hat{\mathbf{y}}) = \mathbb{E} [\|\hat{\mathbf{y}} - \mathbf{y}\|^2] = \mathbb{E} [(\hat{\mathbf{y}} - \mathbf{y})^\top (\hat{\mathbf{y}} - \mathbf{y})]$$

where  $y$  is the true value, and  $\hat{y}$  is our estimation.

- (e) The derivative of the MSE with respect to  $\lambda$  is negative at  $\lambda = 0$ . We omit here the proof. Conclude that, if the linear model is correct, a little Ridge regularization helps to reduce the MSE.

(a)

$$\hat{\omega} = (\mathbf{X}^\top \mathbf{X} + \lambda I_d)^{-1} \mathbf{X}^\top \mathbf{y}$$

$$E[\hat{\omega}_x] = E[A_x \hat{\omega}] \stackrel{?}{=} A_x E[\hat{\omega}] \stackrel{?}{=} A_x \omega$$

ר'ג'ז'נ'ז'י'ז' א'ג'ז'מ'נ'ט' ה'ז'

רנור. מינימום נרחב בפונקציית האנרגיה. אם  $\lambda > 0$ , אז  $A_\lambda w \neq w$  ולכן  $I = (x^T x)^{-1}(x^T x) = A_\lambda \leftrightarrow \lambda = 0$ . אם  $\lambda < 0$ , אז  $A_\lambda \neq I$ .

$$\text{Var}(\hat{\omega}_\lambda) = \text{Var}(A_\lambda \hat{\omega}) = A_\lambda \text{Var}(\hat{\omega}) A_\lambda^T = A_\lambda \sigma^2 (X^T X)^{-1} A_\lambda^T = \sigma^2 A_\lambda (X^T X)^{-1} A_\lambda^T \quad (\text{C})$$

*אנו מוכיחים כי*  $\sigma^2$

- (d) Derive explicit expressions for the (squared) bias and variance of  $\hat{\mathbf{w}}_\lambda$  as a function of  $\lambda$ , i.e. write a bias-variance decomposition for the mean square error of  $\hat{\mathbf{w}}_\lambda$ .

Hint: recall that for the multivariate case the MSE defined to be:

$$MSE(\hat{\mathbf{y}}) = \mathbb{E} [\|\hat{\mathbf{y}} - \mathbf{y}\|^2] = \mathbb{E} [(\hat{\mathbf{y}} - \mathbf{y})^\top (\hat{\mathbf{y}} - \mathbf{y})]$$

:(ד) כו�ה נגזרת נורית חנוקה הינה נמי לא נקי

$$\begin{aligned} \mathbb{E}[(\hat{\mathbf{y}} - \mathbf{y})^2] &= \mathbb{E}[(\hat{\mathbf{y}} - \bar{\mathbf{y}}) + (\bar{\mathbf{y}} - \mathbf{y})]^2 \\ MSE &= \mathbb{E}[(\hat{\mathbf{y}} - \bar{\mathbf{y}})^2] + 2(\bar{\mathbf{y}} - \mathbf{y})\mathbb{E}[\hat{\mathbf{y}} - \bar{\mathbf{y}}] + (\bar{\mathbf{y}} - \mathbf{y})^2 \\ &= \mathbb{E}[(\hat{\mathbf{y}} - \bar{\mathbf{y}})^2] + (\bar{\mathbf{y}} - \mathbf{y})^2 \\ &= \text{var}[\hat{\mathbf{y}}] + \text{bias}^2[\hat{\mathbf{y}}] \end{aligned}$$

In words,

$$MSE = \text{Variance} + \text{bias}^2$$

$$MSE = E(\|(\hat{\mathbf{w}}_\lambda - E(\hat{\mathbf{w}}_\lambda)) + (E(\hat{\mathbf{w}}_\lambda) - \mathbf{w})\|^2) =$$

↓  
נפערת נורית,  $\|u-v\|^2 = \sum_i (u_i - v_i)^2 = \sum_i u_i^2 - 2u_i v_i + v_i^2 = \|u\|^2 - 2\langle u, v \rangle + \|v\|^2$

$$\begin{aligned} &= E(\|\hat{\mathbf{w}}_\lambda - E(\hat{\mathbf{w}}_\lambda)\|^2) + 2(E(\hat{\mathbf{w}}_\lambda) - \mathbf{w})^\top E(\hat{\mathbf{w}}_\lambda - E(\hat{\mathbf{w}}_\lambda)) + E(\|E(\hat{\mathbf{w}}_\lambda) - \mathbf{w}\|^2) = \\ &= E(\|\hat{\mathbf{w}}_\lambda - E(\hat{\mathbf{w}}_\lambda)\|^2) + \|E(\hat{\mathbf{w}}_\lambda) - \mathbf{w}\|^2 = \text{Var}_w(\lambda) + \text{Bias}_w^2(\lambda) \quad E(\hat{\mathbf{w}}_\lambda - E(\hat{\mathbf{w}}_\lambda)) = E(\hat{\mathbf{w}}_\lambda) - E(\hat{\mathbf{w}}_\lambda) = 0 \\ &\quad \text{প্রয়োগ } w, E(\hat{\mathbf{w}}_\lambda) \end{aligned}$$

$$\text{Bias}_w^2(\lambda) \downarrow \frac{(E(\hat{\mathbf{w}}_\lambda) - \mathbf{w})^\top (E(\hat{\mathbf{w}}_\lambda) - \mathbf{w})}{\|A_\lambda \mathbf{w} - \mathbf{w}\|^2} = \frac{\|E(\hat{\mathbf{w}}_\lambda) - \mathbf{w}\|^2}{\|A_\lambda \mathbf{w} - \mathbf{w}\|^2} = \frac{\|(A_\lambda - I)\mathbf{w}\|^2}{\|A_\lambda \mathbf{w} - \mathbf{w}\|^2}$$

$$\text{Var}_w(\lambda) \downarrow \frac{E(\|\hat{\mathbf{w}}_\lambda - E(\hat{\mathbf{w}}_\lambda)\|^2)}{\|A_\lambda \mathbf{w} - \mathbf{w}\|^2} = E\left(\sum_i (\hat{\mathbf{w}}_{\lambda,i} - E(\hat{\mathbf{w}}_{\lambda,i}))^2\right) = \sum_i E(\hat{\mathbf{w}}_{\lambda,i} - E(\hat{\mathbf{w}}_{\lambda,i}))^2 =$$

$$\sum_i \text{Var}(\hat{\mathbf{w}}_{\lambda,i}) = \text{Trace}(\text{Var}(\hat{\mathbf{w}}_\lambda)) \downarrow \text{Trace}(\sigma^2 A_\lambda (\mathbf{x}^\top \mathbf{x})^{-1} A_\lambda^\top) = \sigma^2 \text{Trace}(A_\lambda (\mathbf{x}^\top \mathbf{x})^{-1} A_\lambda^\top)$$

: এ যদি সেই

$$MSE(\hat{\mathbf{w}}_\lambda) = \text{Var}_w(\lambda) + \text{Bias}_w^2(\lambda) = \sigma^2 \text{Trace}(A_\lambda (\mathbf{x}^\top \mathbf{x})^{-1} A_\lambda^\top) + \|(A_\lambda - I)\mathbf{w}\|^2$$

- (e) The derivative of the MSE with respect to  $\lambda$  is negative at  $\lambda = 0$ . We omit here the proof. Conclude that, if the linear model is correct, a little Ridge regularization helps to reduce the MSE.

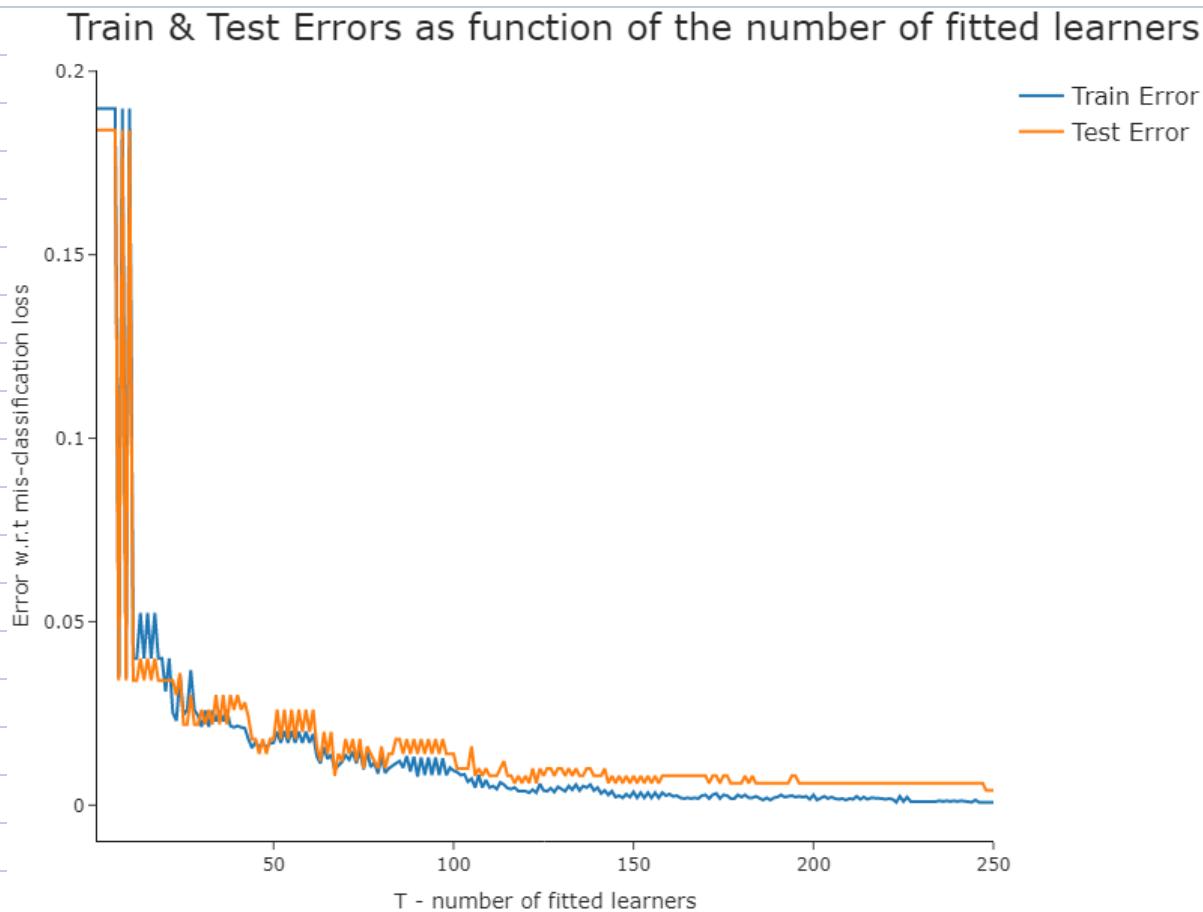
הוכיחו בזאת שORTHOGONALITY גוררת את ה- $\lambda$  ו- $\hat{\omega}$  מינימיזם את MSE. מינימיזם את MSE מושג על ידי  $\hat{\omega} = \hat{\omega}_0 - \lambda I$ . מינימיזם את MSE מושג על ידי  $\hat{\omega} = \hat{\omega}_0 - \lambda I$ . מינימיזם את MSE מושג על ידי  $\hat{\omega} = \hat{\omega}_0 - \lambda I$ .

MSE  $\hat{\omega}$

### 3 Practical Part

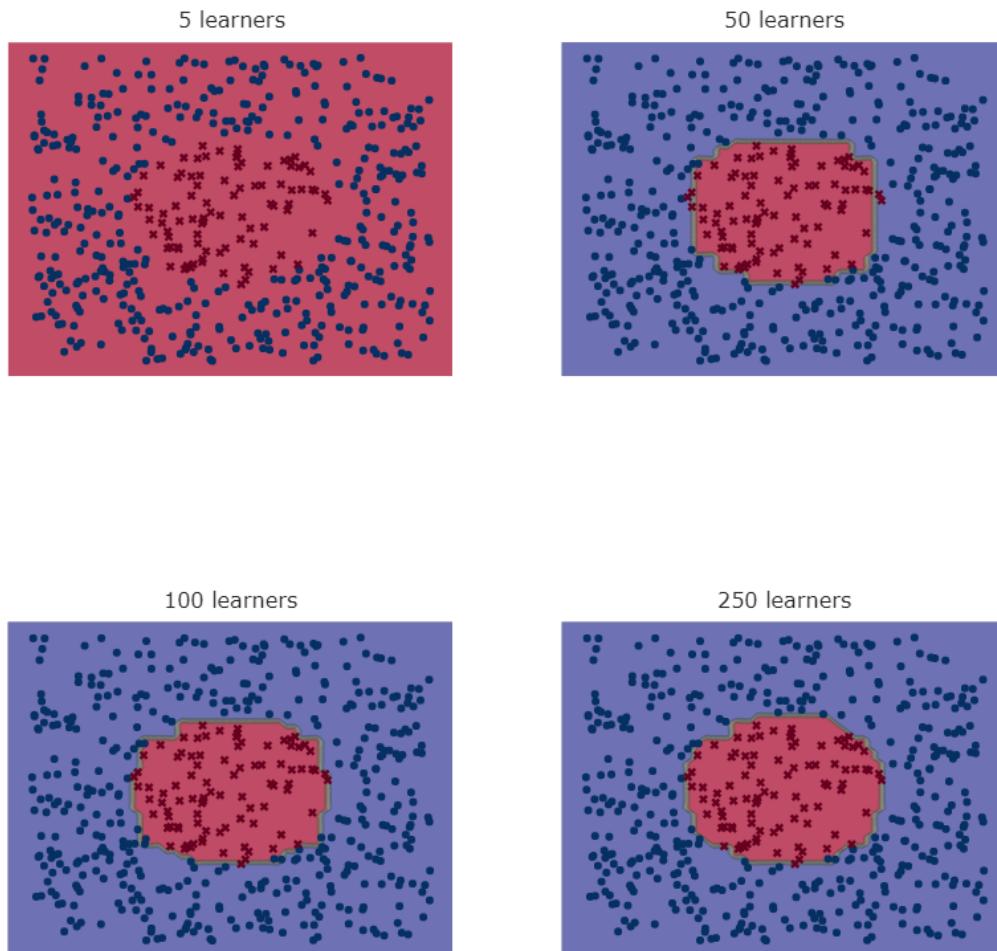
### 3.1 Boosting - Separate the inseparable

1. Use the provided `generate_data` function to generate 5000 train samples and 500 test samples, both with no noise (i.e. `noise_ratio = 0`). Train an Adaboost ensemble of size 250 (i.e passing 250 as the number of iterations) using your implementation of the `DecisionStump` as weak learner. Plot, in a single figure, the training- and test errors as a function of the number of fitted learners (i.e. use the `partial_loss` function for  $t = 1, \dots, 250$ ). Explain your results.



2. Using the previously fitted ensemble, plot the decision boundary obtained by using the ensemble up to iteration 5, 50, 100 and 250. Use your implementation of the `partial_predict` to obtain the predictions of the ensemble up to the specified size. In each of these plots also add the test set (colored and/or shaped by the actual labels). Explain your results.

Decision boundaries using the ensemble up to different number of iterations

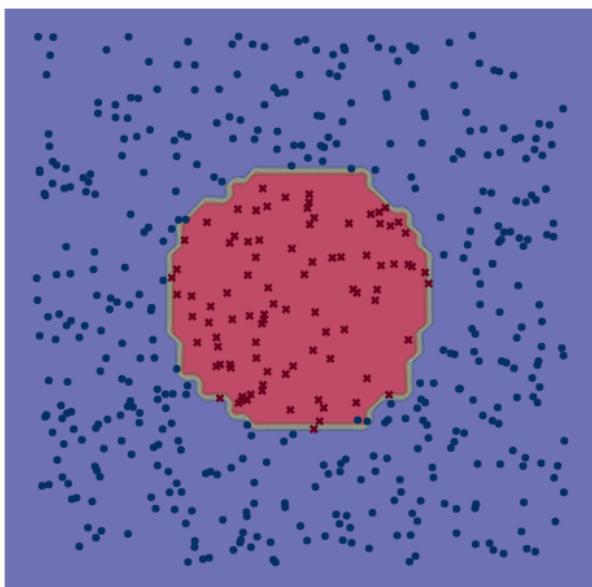


הערכה מודולית של פונקציית ניטרוליזה (ב-ELISA) מושגת באמצעות אוסף של אינטראקציות בין אנטיגן ו-抗体ים. אוסף זה מכונה **ensemble**. בפועל, אוסף אנטיגנים ייחודי ל-ELISA מושג באמצעות אינטראקציה בין אנטיגן ו-抗体ים.

3. Using the test evaluation scores per number of learners, which ensemble size achieved the lowest test error? Plot the decision surface of that ensemble as well as the test set data points. In the plot's title provide the ensemble size and its accuracy.

3

Best Ensemble with Size: 238 & Accuracy: 1.0





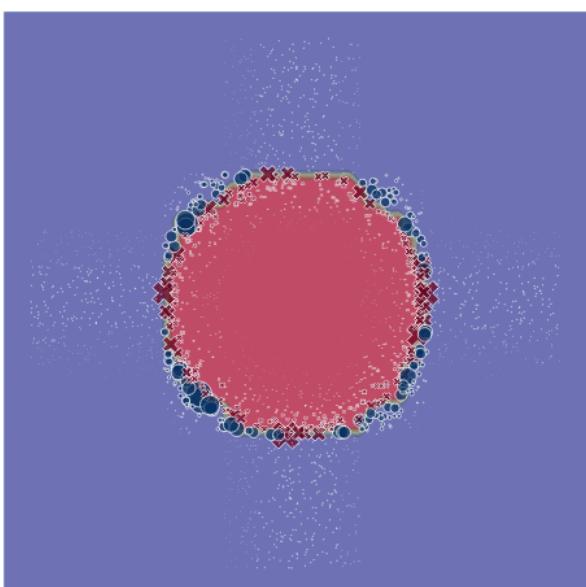
4. Using the previously fitted ensemble, use the weights of the last iteration (i.e.  $\mathcal{D}$  at  $T$ ) to plot the *training set* with a point size proportional to its weight and color (and/or shape) indicating its label.

4

- As previously, plot the decision surface (using the full ensemble).
  - As the weights are of very small numbers normalize and transform as follow:  $D = D/np.max(D) * 5$

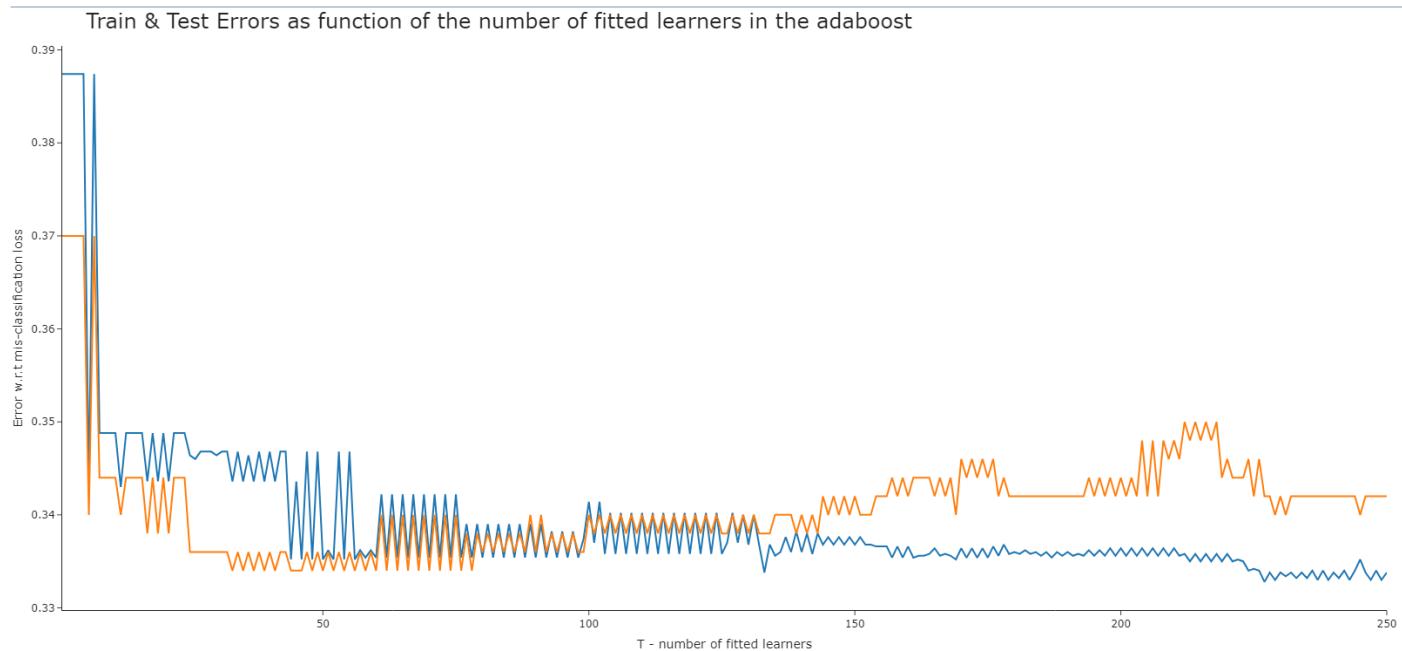
Explain your results, and specifically explain which samples are "easier" and which are more "challenging" for the classifier.

### Final AdaBoost Sample Distribution

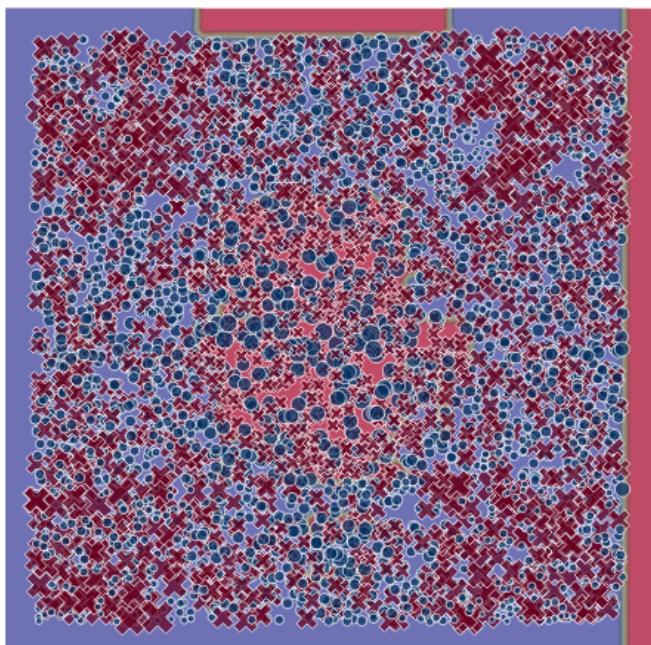


5. Repeat the steps above (while avoiding code repetition) for train- and test sets generated with noise levels of 0.4. Show graphs as in question (1) and question (4). Explain the results. In your answer explain what is seen in the plot of the loss in terms of the bias-variance tradeoff.

.5



Final AdaBoost Sample Distribution

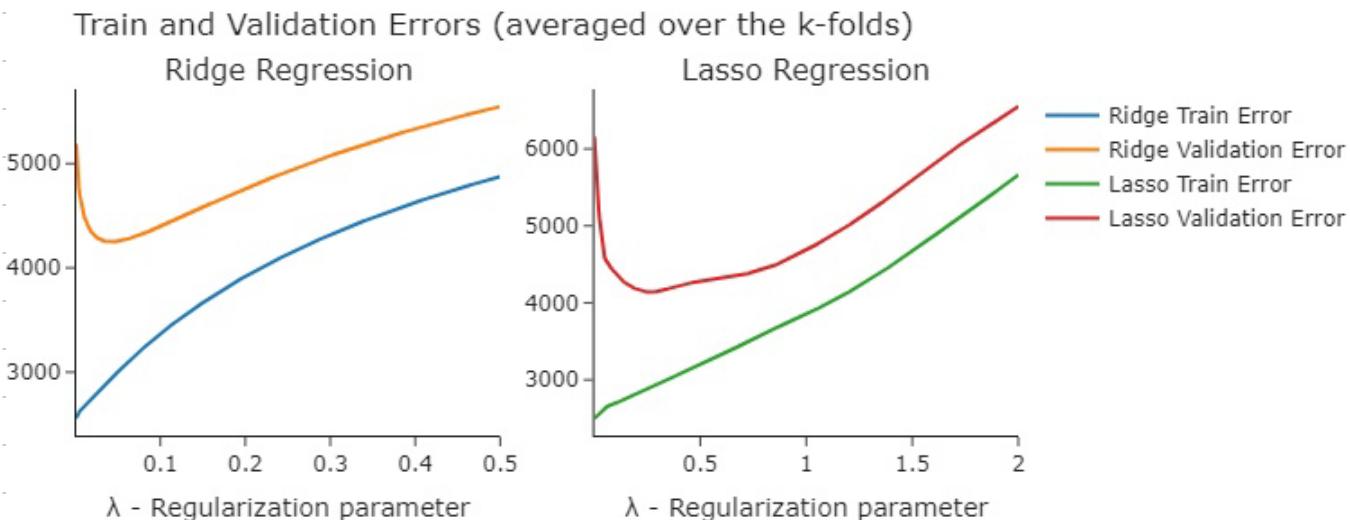


### 3.2 Choosing Regularization Parameters Using Cross Validation

2. For both Ridge (with intercept) and Lasso regularizations, run 5-Fold Cross-Validation using different values of the regularization parameter  $\lambda$ . Explore possible ranges of values of  $\lambda$  (say, take `n_evaluations=500` equally spaced values in a range of your choice). Which range of values is meaningful for the given data for each of the algorithms?

Then, over these selected ranges, for each of the two algorithms, plot the train- and validation errors as a function of the tested regularization parameter value. Explain your results. Address both the differences between the train- and validation errors as well as the differences in the plots for the two algorithms.

הערך הנקרא Ridge נקרא Ridge וערך הנקרא Lasso נקרא Lasso.



המודלים נבדק ב集 validation set - ה- validation set מוגדר כsubset של המuestים. ה- training set מוגדר כsubset של המuestים לא כולל המuestים בvalidation set. בvalidation set מוגדרות ה-  $\theta$  ו-  $J(\theta)$ . ב- test set מוגדרות ה-  $\theta$  ו-  $J(\theta)$ .

הנוגע ל- $\lambda$  נקבע על ידי שיקול Ridge - Lasso  $\lambda$  מינימום פונקציית האמצעים המבוקש. אוסף נתונים אחדים נבחר ומשתמשים בו לשליחת נתונים לא-TRAINING ל-

3. Report which regularization parameter values achieved the best validation errors for the Ridge and Lasso regularizations. Then fit the entire test set using these values for Ridge, Lasso and Least Squares regressions. Report the test errors of each of the fitted models.
- Use your own previously implemented Least Squares algorithm.

```
Best hyperparameter (regularization values):
```

```
Ridge Regression: 0.042
```

```
Lasso Regression: 0.2613907815631263
```

```
Mean Square Error for different Regression Models:
```

```
Ridge Regression: 3128.948548647699
```

```
Lasso Regression: 3188.714771629627
```

```
Least Squares Regression: 3425.32159754408
```

(over test set)