

פתרון תרגיל מספר 5 - מערכות לומדות

שם: מיכאל גרינבאום, ת.ז: 211747639

13 ביוני 2020

1. פתרון:

$$L_D(h^*) \leq \min_{h \in \mathcal{H}_k} L_D(h) + \sqrt{\frac{2 \ln\left(\frac{2|\mathcal{H}_k|}{\delta}\right)}{m}} \quad \text{(א) צ"ל: הוכחה:}$$

נגדיר $\mathcal{H} = \mathcal{H}_k$ ונגדיר $\varepsilon = \sqrt{\frac{\ln\left(\frac{2|\mathcal{H}|}{\delta}\right)}{2m}}$, ראינו בהרצאה 5 בהוכחה של המשפט היסודי של ה- IML שמתקיים (ראינו בכך ש- L_S הוא L_D בתוחלת, והוא ממוצע של משתנים מקרים חסומים ב- $0 - 1$ מאותה התפלגות ובלתי תלויים)

$$\mathbb{P}_{S \sim \mathcal{D}^m} (S \mid \exists h \in \mathcal{H} \text{ s.t. } |L_D(h) - L_S(h)| > \varepsilon) \leq 2 \cdot |\mathcal{H}| \cdot \exp(-2m \cdot \varepsilon^2)$$

נעזר בתכונה זאת ונראה כי

$$\begin{aligned} \mathbb{P}_{S \sim \mathcal{D}^m} (S \mid \exists h \in \mathcal{H} \text{ s.t. } |L_D(h) - L_S(h)| > \varepsilon) &\leq 2 \cdot |\mathcal{H}| \cdot \exp(-2m \cdot \varepsilon^2) \\ &= 2 \cdot |\mathcal{H}| \cdot \exp\left(-2m \cdot \frac{\ln\left(\frac{2|\mathcal{H}|}{\delta}\right)}{2m}\right) = 2 \cdot |\mathcal{H}| \cdot \exp\left(-\ln\left(\frac{2|\mathcal{H}|}{\delta}\right)\right) = 2 \cdot |\mathcal{H}| \cdot \frac{\delta}{2 \cdot |\mathcal{H}|} = \delta \end{aligned}$$

כלומר נקבל כי

$$\begin{aligned} \mathbb{P}_{S \sim \mathcal{D}^m} (S \mid S \text{ is } \varepsilon \text{ representative}) &= 1 - \mathbb{P}_{S \sim \mathcal{D}^m} (S \mid S \text{ is not } \varepsilon \text{ representative}) \\ &= 1 - \mathbb{P}_{S \sim \mathcal{D}^m} (S \mid \exists h \in \mathcal{H} \text{ s.t. } |L_D(h) - L_S(h)| > \varepsilon) = 1 - \delta \end{aligned}$$

תהי $S \sim \mathcal{D}^m$ כך ש- S is ε representative, תהי $h \in \mathcal{H}$, נשים לב כי

$$\begin{aligned} L_D(h^*) &\stackrel{S \text{ is } \varepsilon \text{ representative}}{\leq} L_S(h^*) + \varepsilon \\ &\stackrel{h^* \text{ minimizes the error on } S}{\leq} L_S(h) + \varepsilon \\ &\stackrel{S \text{ is } \varepsilon \text{ representative}}{\leq} L_D(h) + \varepsilon + \varepsilon = L_D(h) + 2 \cdot \varepsilon \\ &= L_D(h) + 2 \cdot \sqrt{\frac{\ln\left(\frac{2|\mathcal{H}|}{\delta}\right)}{2m}} = L_D(h) + \sqrt{\frac{2 \cdot \ln\left(\frac{2|\mathcal{H}|}{\delta}\right)}{m}} \end{aligned}$$

ולכן מהיות זה נכון לכל $h \in \mathcal{H}$ נקבל כי $L_D(h^*) \leq \min_{h \in \mathcal{H}} L_D(h) + \sqrt{\frac{2 \cdot \ln\left(\frac{2|\mathcal{H}|}{\delta}\right)}{m}}$ ולכן מתקיים

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left(L_D(h^*) \leq \min_{h \in \mathcal{H}_k} L_D(h) + \sqrt{\frac{2 \cdot \ln\left(\frac{2|\mathcal{H}_k|}{\delta}\right)}{m}} \right) \geq \mathbb{P}_{S \sim \mathcal{D}^m} (S \mid S \text{ is } \varepsilon \text{ representative}) = 1 - \delta$$

מ.ש.ל.א.☺

$$L_{\mathcal{D}}(h^*) \leq \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{2}{\alpha \cdot m} \cdot \ln\left(\frac{4k}{\delta}\right)} + \sqrt{\frac{2}{(1-\alpha) \cdot m} \cdot \ln\left(\frac{4|\mathcal{H}_j|}{\delta}\right)} \quad \text{ב) צ"ל:}$$

הוכחה:

תחילה נשים לב ש- $|V| = \alpha \cdot m$ וגם $|S| = (1-\alpha) \cdot m$,
 עתה נסתכל על שלב הולידציה, יש לנו $|\{h_1, \dots, h_k\}| = k$ היפותזות ואנחנו ממזערים את השגיאה על $|V|$, לכן לפי מה שהוכח בסעיף הקודם מתקיים כי

$$\begin{aligned} & \mathbb{P}\left(L_{\mathcal{D}}(h^*) > \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{2 \cdot \ln\left(\frac{4k}{\delta}\right)}{\alpha \cdot m}}\right) \\ &= \mathbb{P}\left(L_{\mathcal{D}}(h^*) > \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{2 \cdot \ln\left(\frac{2|\{h_1, \dots, h_k\}|}{\left(\frac{\delta}{2}\right)}\right)}{|V|}}\right) < \frac{\delta}{2} \end{aligned}$$

עתה, נסתכל על שלה האימון של \mathcal{H}_j , נשים לב כי מחלקת ההיפותזות היא \mathcal{H}_j ואנחנו ממזערים את השגיאה כל $|S|$ ולכן לפי מה שהוכח בסעיף הקודם נקבל כי

$$\begin{aligned} & \mathbb{P}\left(L_{\mathcal{D}}(h_j) > \min_{h \in \mathcal{H}_j} L_{\mathcal{D}}(h) + \sqrt{\frac{2 \cdot \ln\left(\frac{4|\mathcal{H}_j|}{\delta}\right)}{(1-\alpha) \cdot m}}\right) \\ &= \mathbb{P}\left(L_{\mathcal{D}}(h_j) > \min_{h \in \mathcal{H}_j} L_{\mathcal{D}}(h) + \sqrt{\frac{2 \cdot \ln\left(\frac{2|\mathcal{H}_j|}{\left(\frac{\delta}{2}\right)}\right)}{|S|}}\right) < \frac{\delta}{2} \end{aligned}$$

עתה נבחין כי מתקיים $\min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) = \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h)$ וגם מהנתון $h_j = h^*$ ולכן

$$\begin{aligned} & \mathbb{P}\left(L_{\mathcal{D}}(h^*) > \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{2}{\alpha \cdot m} \cdot \ln\left(\frac{4k}{\delta}\right)} + \sqrt{\frac{2}{(1-\alpha) \cdot m} \cdot \ln\left(\frac{4|\mathcal{H}_j|}{\delta}\right)}\right) \\ & \leq \mathbb{P}\left(\left(L_{\mathcal{D}}(h^*) > \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{2}{\alpha \cdot m} \cdot \ln\left(\frac{4k}{\delta}\right)}\right) \cup \left(L_{\mathcal{D}}(h^*) > \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{2}{(1-\alpha) \cdot m} \cdot \ln\left(\frac{4|\mathcal{H}_j|}{\delta}\right)}\right)\right) \\ & \leq \mathbb{P}\left(L_{\mathcal{D}}(h^*) > \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{2}{\alpha \cdot m} \cdot \ln\left(\frac{4k}{\delta}\right)}\right) + \mathbb{P}\left(L_{\mathcal{D}}(h^*) > \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{2}{(1-\alpha) \cdot m} \cdot \ln\left(\frac{4|\mathcal{H}_j|}{\delta}\right)}\right) \\ & < \frac{\delta}{2} + \mathbb{P}\left(L_{\mathcal{D}}(h_j) > \min_{h \in \mathcal{H}_j} L_{\mathcal{D}}(h) + \sqrt{\frac{2}{(1-\alpha) \cdot m} \cdot \ln\left(\frac{4|\mathcal{H}_j|}{\delta}\right)}\right) < \frac{\delta}{2} + \frac{\delta}{2} = \delta \end{aligned}$$

ולכן נקבל כי

$$\begin{aligned} & \mathbb{P}\left(L_{\mathcal{D}}(h^*) \leq \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{2}{\alpha \cdot m} \cdot \ln\left(\frac{4k}{\delta}\right)} + \sqrt{\frac{2}{(1-\alpha) \cdot m} \cdot \ln\left(\frac{4|\mathcal{H}_j|}{\delta}\right)}\right) \\ &= 1 - \mathbb{P}\left(L_{\mathcal{D}}(h^*) > \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{2}{\alpha \cdot m} \cdot \ln\left(\frac{4k}{\delta}\right)} + \sqrt{\frac{2}{(1-\alpha) \cdot m} \cdot \ln\left(\frac{4|\mathcal{H}_j|}{\delta}\right)}\right) \\ & \geq 1 - \delta \end{aligned}$$

כלומר קיבלנו

$$\mathbb{P} \left(L_{\mathcal{D}}(h^*) \leq \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{2}{\alpha \cdot m} \cdot \ln \left(\frac{4k}{\delta} \right)} + \sqrt{\frac{2}{(1-\alpha) \cdot m} \cdot \ln \left(\frac{4|\mathcal{H}_j|}{\delta} \right)} \right) \geq 1 - \delta$$

כנדרש

מ.ש.ל.ב.⊙

(ג) צ"ל: אי אפשר להשוות בין 2 החסמים
הנוכחה:

נסמן את החסמים ב- $\varepsilon_2 = \sqrt{\frac{2}{\alpha \cdot m} \cdot \ln \left(\frac{4k}{\delta} \right)} + \sqrt{\frac{2}{(1-\alpha) \cdot m} \cdot \ln \left(\frac{4|\mathcal{H}_j|}{\delta} \right)}$, $\varepsilon_1 = \sqrt{\frac{2 \ln \left(\frac{2|\mathcal{H}_k|}{\delta} \right)}{m}}$

נסתכל על המנה שלהם ונקבל

$$\begin{aligned} \frac{\varepsilon_2}{\varepsilon_1} &= \frac{\sqrt{\frac{2}{\alpha \cdot m} \cdot \ln \left(\frac{4k}{\delta} \right)} + \sqrt{\frac{2}{(1-\alpha) \cdot m} \cdot \ln \left(\frac{4|\mathcal{H}_j|}{\delta} \right)}}{\sqrt{\frac{2 \ln \left(\frac{2|\mathcal{H}_k|}{\delta} \right)}{m}}} = \frac{\sqrt{\frac{2}{\alpha \cdot m} \cdot \ln \left(\frac{4k}{\delta} \right)}}{\sqrt{\frac{2 \ln \left(\frac{2|\mathcal{H}_k|}{\delta} \right)}{m}}} + \frac{\sqrt{\frac{2}{(1-\alpha) \cdot m} \cdot \ln \left(\frac{4|\mathcal{H}_j|}{\delta} \right)}}{\sqrt{\frac{2 \ln \left(\frac{2|\mathcal{H}_k|}{\delta} \right)}{m}}} \\ &= \sqrt{\frac{\frac{2}{\alpha \cdot m} \cdot \ln \left(\frac{4k}{\delta} \right)}{\frac{2 \ln \left(\frac{2|\mathcal{H}_k|}{\delta} \right)}{m}}} + \sqrt{\frac{\frac{2}{(1-\alpha) \cdot m} \cdot \ln \left(\frac{4|\mathcal{H}_j|}{\delta} \right)}{\frac{2 \ln \left(\frac{2|\mathcal{H}_k|}{\delta} \right)}{m}}} = \sqrt{\frac{\ln \left(\frac{4k}{\delta} \right)}{\alpha \cdot \ln \left(\frac{2|\mathcal{H}_k|}{\delta} \right)}} + \sqrt{\frac{\ln \left(\frac{4|\mathcal{H}_j|}{\delta} \right)}{(1-\alpha) \cdot \ln \left(\frac{2|\mathcal{H}_k|}{\delta} \right)}} \end{aligned}$$

≤ מקרה שבו החסם הראשון יותר טוב, נבחר $\mathcal{H}_1 = \dots = \mathcal{H}_k = \mathcal{H}$ (כש \mathcal{H} בגודל קבוע ללא תלות ב- k) ואז נקבל כי

$$\begin{aligned} \frac{\varepsilon_2}{\varepsilon_1} &= \sqrt{\frac{\ln \left(\frac{4k}{\delta} \right)}{\alpha \cdot \ln \left(\frac{2|\mathcal{H}_k|}{\delta} \right)}} + \sqrt{\frac{\ln \left(\frac{4|\mathcal{H}_j|}{\delta} \right)}{(1-\alpha) \cdot \ln \left(\frac{2|\mathcal{H}_k|}{\delta} \right)}} \\ &= \underbrace{\sqrt{\frac{1}{\alpha \cdot \ln \left(\frac{2|\mathcal{H}|}{\delta} \right)}}}_{\text{constant in k}} \cdot \underbrace{\sqrt{\ln \left(\frac{4k}{\delta} \right)}}_{\text{goes to infinity when } k \rightarrow \infty} + \underbrace{\sqrt{\frac{\ln \left(\frac{4|\mathcal{H}_j|}{\delta} \right)}{(1-\alpha) \cdot \ln \left(\frac{2|\mathcal{H}|}{\delta} \right)}}}_{\text{constant in k}} \end{aligned}$$

כלומר נקבל כי $\lim_{k \rightarrow \infty} \frac{\varepsilon_2}{\varepsilon_1} = \infty$ כלומר ε_1 חסם יותר טוב מ- ε_2 במקרה הזה

⇒ מקרה שבו החסם השני יותר טוב, נבחר $|\mathcal{H}_i| = e^{i^2}$ לכל $i \in [k]$ ונקבל כי

$$\begin{aligned} \frac{\varepsilon_2}{\varepsilon_1} &= \sqrt{\frac{\ln \left(\frac{4k}{\delta} \right)}{\alpha \cdot \ln \left(\frac{2 \cdot e^{k^2}}{\delta} \right)}} + \sqrt{\frac{\ln \left(\frac{4e^{j^2}}{\delta} \right)}{(1-\alpha) \cdot \ln \left(\frac{2e^{k^2}}{\delta} \right)}} \\ &= \sqrt{\frac{\ln \left(\frac{4k}{\delta} \right)}{\alpha \cdot \left[\ln \left(\frac{2}{\delta} \right) + \ln(e^{k^2}) \right]}} + \sqrt{\frac{\ln \left(\frac{4}{\delta} \right) + \ln(e^{j^2})}{(1-\alpha) \cdot \left[\ln \left(\frac{2}{\delta} \right) + \ln(e^{k^2}) \right]}} \\ &= \sqrt{\frac{\ln \left(\frac{4k}{\delta} \right)}{\alpha \cdot \left[\ln \left(\frac{2}{\delta} \right) + k^2 \right]}} + \sqrt{\frac{\ln \left(\frac{4}{\delta} \right) + j^2}{(1-\alpha) \cdot \left[\ln \left(\frac{2}{\delta} \right) + k^2 \right]}} \\ &= \underbrace{\sqrt{\frac{\ln \left(\frac{4k}{\delta} \right)}{\alpha \cdot \left[\ln \left(\frac{2}{\delta} \right) + k^2 \right]}}}_{\text{goes to 0 when } k \rightarrow \infty \text{ because } \frac{\ln(k)}{k^2} \rightarrow 0} + \underbrace{\sqrt{\frac{\ln \left(\frac{4}{\delta} \right) + j^2}{(1-\alpha) \cdot \left[\ln \left(\frac{2}{\delta} \right) + k^2 \right]}}}_{\text{goes to 0 when } k \rightarrow \infty \text{ because } \frac{j^2}{k^2} \rightarrow 0 \text{ for big enough k, j will stay fixed}} \end{aligned}$$

כלומר נקבל כי $\lim_{k \rightarrow \infty} \frac{\varepsilon_2}{\varepsilon_1} = 0$ כלומר ε_2 חסם יותר טוב מ- ε_1 במקרה הזה

מ.ש.ל.ג. ☺

2. פתרון:

$$\hat{w}_\lambda^{ridge} = \frac{\hat{w}^{LS}}{1+\lambda} \quad \text{צ"ל:} \quad \text{הוכחה:}$$

נזכר שבמקרה שלנו מתקיים $\hat{w}^{LS} = (X \cdot X^T)^{-1} \cdot X \cdot y = X \cdot y$ נזכר כי

$$f_{l_2}(w) = \|y - X^T \cdot w\|^2 + \lambda \cdot \|w\|^2$$

לכן

$$\begin{aligned} J_w(f_{l_2}) &= J_w(\|y - X^T \cdot w\|^2 + \lambda \cdot \|w\|^2) = J_w(\|y - X^T \cdot w\|^2) + J_w(\lambda \cdot \|w\|^2) \\ &= J_{y-X^T \cdot w}(\|y - X^T \cdot w\|^2) \cdot J_w(y - X^T \cdot w) + 2 \cdot \lambda \cdot w^T \\ &= 2 \cdot (y - X^T \cdot w)^T \cdot [-X^T] + 2 \cdot \lambda \cdot w^T \\ &= 2 \cdot [w^T X X^T - y^T X^T + \lambda w^T] = 2 \cdot [X \cdot X^T \cdot w - X \cdot y + \lambda w]^T \\ &= 2 \cdot [I_n \cdot w - \hat{w}^{LS} + \lambda w]^T = 2 \cdot [(1 + \lambda) \cdot w - \hat{w}^{LS}]^T \end{aligned}$$

לכן אנחנו רוצים ש- $\nabla(f_{l_2}) = 0$

$$0 = \nabla(f_{l_2}) = J_w(f_{l_2})^T = [2 \cdot [(1 + \lambda) \cdot w - \hat{w}^{LS}]^T]^T = 2 \cdot [(1 + \lambda) \cdot w - \hat{w}^{LS}]$$

$$\Rightarrow 2 \cdot (1 + \lambda) \cdot w = 2 \cdot \hat{w}^{LS} \Rightarrow \boxed{w = \frac{1}{1 + \lambda} \cdot \hat{w}^{LS}}$$

כלומר כל נקודת קיצון w של f_{l_2} חייבת לקיים $w = \frac{1}{1+\lambda} \cdot \hat{w}^{LS}$

ולכן $\hat{w}_\lambda^{ridge} = \frac{1}{1 + \lambda} \cdot \hat{w}^{LS}$ כי הוא נקודת מינימום של f_{l_2} ובפרט נקודת קיצון, כנדרש

מ.ש.ל.א. ☺

$$\hat{w}_\lambda^{subset} = \eta_{\sqrt{\lambda}}^{hard}(\hat{w}^{LS}) \quad \text{צ"ל:} \quad \text{הוכחה:}$$

נזכר שבמקרה שלנו מתקיים $\hat{w}^{LS} = (X \cdot X^T)^{-1} \cdot X \cdot y = X \cdot y$ נזכר כי

$$\begin{aligned} f_{l_0}(w) &= \|y - X^T \cdot w\|^2 + \lambda \cdot \|w\|_0 = \|X^{-1} \cdot [X \cdot y - X \cdot X^T \cdot w]\|^2 + \lambda \cdot \|w\|_0 \\ &= (X^{-1} \cdot [X \cdot y - X \cdot X^T \cdot w])^T \cdot (X^{-1} \cdot [X \cdot y - X \cdot X^T \cdot w]) + \lambda \cdot \|w\|_0 \\ &= [X \cdot y - X \cdot X^T \cdot w]^T \cdot [X^{-1}]^T X^{-1} \cdot [X \cdot y - X \cdot X^T \cdot w] + \lambda \cdot \|w\|_0 \\ &= [X \cdot y - X \cdot X^T \cdot w]^T \cdot (X \cdot X^T)^{-1} \cdot [X \cdot y - X \cdot X^T \cdot w] + \lambda \cdot \|w\|_0 \\ &= [X \cdot y - X \cdot X^T \cdot w]^T \cdot I_n \cdot [X \cdot y - X \cdot X^T \cdot w] + \lambda \cdot \|w\|_0 \\ &= \|X \cdot y - X \cdot X^T \cdot w\|^2 + \lambda \cdot \|w\|_0 = \|\hat{w}^{LS} - X \cdot X^T \cdot w\|^2 + \lambda \cdot \|w\|_0 \\ &= \sum_{i=1}^{|w|} (\hat{w}_i^{LS} - w_i)^2 + \lambda \cdot 1_{w_i \neq 0} \end{aligned}$$

נרצה למצוא מזעור ל- f_{l_0} , נשים לב שנוכל למזער כל אינדקס בנפרד.

יהי $i \in [|\hat{w}_i^{LS}|]$ נחלק ל- 2 מקרים:

i. אם $\hat{w}_i \neq 0$ אז $|\hat{w}_i^{LS}| < \sqrt{\lambda}$ אם $\hat{w}_i \neq 0$

$$[f_{l_0}(\hat{w})]_i = \left(\hat{w}_i^{LS} - \hat{w}_i\right)^2 + \lambda \cdot 1_{w_i \neq 0} \geq \lambda \cdot 1_{w_i \neq 0} = \lambda$$

ואם נבחר $\hat{w}_i = 0$ נקבל

$$[f_{l_0}(\hat{w})]_i = \left(\hat{w}_i^{LS} - \hat{w}_i\right)^2 + \lambda \cdot 1_{w_i \neq 0} = \left(\hat{w}_i^{LS} - 0\right)^2 + \lambda \cdot 0 = \left|\hat{w}_i^{LS}\right|^2 < \sqrt{\lambda}^2 = \lambda$$

כלומר המינימום מתקבל כאשר $\hat{w}_i = 0$ ונשים לב כי $\hat{w}_i = 0 = \eta_{\sqrt{\lambda}}^{hard}(\hat{w}^{LS})$ כנדרש

ii. אחרת $\hat{w}_i = 0$ אם $|\hat{w}_i^{LS}| \geq \sqrt{\lambda}$ אז

$$[f_{l_0}(\hat{w})]_i = \left(\hat{w}_i^{LS} - \hat{w}_i\right)^2 + \lambda \cdot 1_{w_i \neq 0} \geq \left(\hat{w}_i^{LS} - 0\right)^2 = \left|\hat{w}_i^{LS}\right|^2 = \sqrt{\lambda}^2 = \lambda$$

ענה נשים לב שעבור $\hat{w}_i \neq 0$ נקבל

$$[f_{l_0}(\hat{w})]_i = \left(\hat{w}_i^{LS} - \hat{w}_i\right)^2 + \lambda \cdot 1_{w_i \neq 0} = \left(\hat{w}_i^{LS} - \hat{w}_i\right)^2 + \lambda$$

המינימום יתקבל כאשר $\hat{w}_i = \hat{w}_i^{LS}$ ובמקרה זה נקבל

$$[f_{l_0}(\hat{w})]_i = \left(\hat{w}_i^{LS} - \hat{w}_i\right)^2 + \lambda \cdot 1_{w_i \neq 0} = \left(\hat{w}_i^{LS} - \hat{w}_i^{LS}\right)^2 + \lambda = \lambda$$

כלומר המינימום מתקבל כאשר $\hat{w}_i = \hat{w}_i^{LS}$ ונשים לב כי $\hat{w}_i = \hat{w}_i^{LS} = \eta_{\sqrt{\lambda}}^{hard}(\hat{w}^{LS})$ כנדרש

מ.ש.ל.ב. ☺

3. פתרון:

$$A_\lambda = (X \cdot X^T + \lambda \cdot I_n)^{-1} \cdot (X \cdot X^T)$$

צ"ל: **הוכחה:**

ראינו בתרגול שמתקיים $\hat{w}(\lambda) = (X \cdot X^T + \lambda \cdot I_n)^{-1} \cdot X \cdot y$ ולכן

$$\begin{aligned} \hat{w}(\lambda) &= (X \cdot X^T + \lambda \cdot I_n)^{-1} \cdot X \cdot y \\ &= (X \cdot X^T + \lambda \cdot I_n)^{-1} \cdot \left[(X \cdot X^T) \cdot (X \cdot X^T)^{-1} \right] \cdot X \cdot y \\ &= (X \cdot X^T + \lambda \cdot I_n)^{-1} \cdot (X \cdot X^T) \cdot \left[(X \cdot X^T)^{-1} \cdot X \cdot y \right] \\ &= (X \cdot X^T + \lambda \cdot I_n)^{-1} \cdot (X \cdot X^T) \cdot \hat{w} = A_\lambda \cdot \hat{w} \end{aligned}$$

מ.ש.ל.א. ☺

$$\mathbb{E}[\hat{w}(\lambda)]$$

צ"ל: **הוכחה:**

ראינו בתרגול שמתקיים

$$\mathbb{E}[\hat{w}(\lambda)] = \mathbb{E}[A_\lambda \cdot \hat{w}] = A_\lambda \cdot \mathbb{E}[\hat{w}] = A_\lambda \cdot w = (X \cdot X^T + \lambda \cdot I_n)^{-1} \cdot (X \cdot X^T) \cdot w$$

נשים לב שעבור $\lambda > 0$ מתקיים

$$\begin{aligned}
 w - \mathbb{E}[\hat{w}(\lambda)] &= w - (X \cdot X^T + \lambda \cdot I_n)^{-1} \cdot (X \cdot X^T) \cdot w \\
 &= w \cdot \left(I_n - (X \cdot X^T + \lambda \cdot I_n)^{-1} \cdot (X \cdot X^T) \right) \\
 &= w \cdot \left((X \cdot X^T + \lambda \cdot I_n)^{-1} \cdot (X \cdot X^T + \lambda \cdot I_n) - (X \cdot X^T + \lambda \cdot I_n)^{-1} \cdot (X \cdot X^T) \right) \\
 &= w \cdot (X \cdot X^T + \lambda \cdot I_n)^{-1} \cdot ((X \cdot X^T + \lambda \cdot I_n) - (X \cdot X^T)) \\
 &= w \cdot (X \cdot X^T + \lambda \cdot I_n)^{-1} \cdot (\lambda \cdot I_n) \neq 0 \\
 &\Rightarrow \boxed{w \neq \mathbb{E}[\hat{w}(\lambda)]}
 \end{aligned}$$

כלומר קיבלנו כי $\boxed{\mathbb{E}[\hat{w}(\lambda)] = (X \cdot X^T + \lambda \cdot I_n)^{-1} \cdot (X \cdot X^T) \cdot w}$ ולכל $\lambda > 0$ מתקיים $\mathbb{E}[\hat{w}(\lambda)] \neq w$ כנדרש
 מ.ש.ל.ב. ☺

(ג) צ"ל: $Var[\hat{w}(\lambda)]$
 הוכחה:

$$\begin{aligned}
 Var[\hat{w}(\lambda)] &= Var[A_\lambda \cdot \hat{w}] \stackrel{A_\lambda \text{ constant non random matrix}}{=} A_\lambda \cdot Var[\hat{w}] \cdot A_\lambda^T \\
 &= A_\lambda \cdot \sigma^2 \cdot (X \cdot X^T)^{-1} \cdot A_\lambda^T = \sigma^2 \cdot A_\lambda \cdot (X \cdot X^T)^{-1} \cdot A_\lambda^T
 \end{aligned}$$

כלומר קיבלנו כי

$$\boxed{Var[\hat{w}(\lambda)] = \sigma^2 \cdot A_\lambda \cdot (X \cdot X^T)^{-1} \cdot A_\lambda^T}$$

מ.ש.ל.ב. ☺

(ד) צ"ל: $MSE(\lambda)$
 הוכחה:
 בהרצאה ראינו ש-

$$MSE(\hat{y}) = \|y - \hat{y}\|^2 = \|y - \mathbb{E}[\hat{y}]\|^2 + \sum_{i=1}^{|y|} Var(\hat{y}_i)$$

נשים לב ש- $Var(\hat{w}(\lambda))$ זאת מטריצת וריאס ולכן $Var(\hat{y}_i)$ במקרה שלנו הוא באלכסון במקום i, i כלומר ב- $[Var(\hat{w}(\lambda))]_{i,i}$, נציב ונקבל

$$\begin{aligned}
 MSE(\hat{w}(\lambda)) &= \|w - \mathbb{E}[\hat{w}(\lambda)]\|^2 + \sum_{i=1}^{|w|} [Var(\hat{w}(\lambda))]_{i,i} \\
 &= \|w - A_\lambda \cdot w\|^2 + Tr(Var(\hat{w}(\lambda))) \\
 &= \|(I_n - A_\lambda) \cdot w\|^2 + Tr\left(\sigma^2 \cdot A_\lambda \cdot (X \cdot X^T)^{-1} \cdot A_\lambda^T\right) \\
 &= [(I_n - A_\lambda) \cdot w]^T \cdot [(I_n - A_\lambda) \cdot w] + Tr\left(\sigma^2 \cdot A_\lambda \cdot (X \cdot X^T)^{-1} \cdot A_\lambda^T\right) \\
 &= w^T \cdot (I_n - A_\lambda)^T \cdot (I_n - A_\lambda) \cdot w + Tr\left(\sigma^2 \cdot A_\lambda \cdot (X \cdot X^T)^{-1} \cdot A_\lambda^T\right) \\
 &= w^T \cdot [I_n - A_\lambda^T - A_\lambda + A_\lambda^T \cdot A_\lambda] \cdot w + Tr\left(\sigma^2 \cdot A_\lambda \cdot (X \cdot X^T)^{-1} \cdot A_\lambda^T\right)
 \end{aligned}$$

עתה נגזור

$$\begin{aligned}
\frac{\partial A_\lambda}{\partial \lambda} &= \frac{\partial (X \cdot X^T + \lambda \cdot I_n)^{-1} \cdot (X \cdot X^T)}{\partial \lambda} = \frac{\partial (X \cdot X^T + \lambda \cdot I_n)^{-1}}{\partial \lambda} \cdot (X \cdot X^T) \\
&= -1 \cdot (X \cdot X^T + \lambda \cdot I_n)^{-2} \cdot \frac{\partial (X \cdot X^T + \lambda \cdot I_n)}{\partial \lambda} \cdot (X \cdot X^T) \\
&= -1 \cdot (X \cdot X^T + \lambda \cdot I_n)^{-2} \cdot I_n \cdot (X \cdot X^T) \\
&= -(X \cdot X^T + \lambda \cdot I_n)^{-2} \cdot (X \cdot X^T) \\
&= -(X \cdot X^T + \lambda \cdot I_n)^{-1} \cdot (X \cdot X^T + \lambda \cdot I_n)^{-1} \cdot (X \cdot X^T) \\
&= -(X \cdot X^T + \lambda \cdot I_n)^{-1} \cdot A_\lambda
\end{aligned}$$

נשים לב כי $\lambda = 0$ נציב $A_0 = (X \cdot X^T + 0 \cdot I_n)^{-1} \cdot (X \cdot X^T) = (X \cdot X^T)^{-1} \cdot (X \cdot X^T) = I_n$

$$\frac{\partial A_\lambda}{\partial \lambda} \big|_{\lambda=0} = -(X \cdot X^T + 0 \cdot I_n)^{-1} \cdot A_0 = -(X \cdot X^T)^{-1} \cdot I_n = -(X \cdot X^T)^{-1}$$

עתה נשתמש בנגזרת הזאת כדי לגזור את הביטויים שקיבלנו

$$\begin{aligned}
\frac{\partial Tr(\sigma^2 \cdot A_\lambda \cdot (X \cdot X^T)^{-1} \cdot A_\lambda^T)}{\partial \lambda} \big|_{\lambda=0} &= Tr \left(\frac{\partial \sigma^2 \cdot A_\lambda \cdot (X \cdot X^T)^{-1} \cdot A_\lambda^T}{\partial \lambda} \right) \big|_{\lambda=0} \\
&= Tr \left[\sigma^2 \cdot \frac{\partial A_\lambda \cdot (X \cdot X^T)^{-1} \cdot A_\lambda^T}{\partial \lambda} \right] \big|_{\lambda=0} = \sigma^2 \cdot Tr \left[\frac{\partial A_\lambda \cdot (X \cdot X^T)^{-1} \cdot A_\lambda^T}{\partial \lambda} \right] \big|_{\lambda=0} \\
&= \sigma^2 \cdot Tr \left[\frac{\partial A_\lambda}{\partial \lambda} \cdot (X \cdot X^T)^{-1} \cdot A_\lambda^T + A_\lambda \cdot (X \cdot X^T)^{-1} \cdot \frac{\partial A_\lambda^T}{\partial \lambda} \right] \big|_{\lambda=0} \\
&= \sigma^2 \cdot Tr \left[-(X \cdot X^T)^{-1} \cdot (X \cdot X^T)^{-1} \cdot A_0^T + A_0 \cdot (X \cdot X^T)^{-1} \cdot -(X \cdot X^T)^{-1} \right] \\
&= \sigma^2 \cdot Tr \left[-(X \cdot X^T)^{-2} - (X \cdot X^T)^{-2} \right] = -2 \cdot \sigma^2 \cdot Tr \left[(X \cdot X^T)^{-2} \right]
\end{aligned}$$

נשים לב כי $(X \cdot X^T)^{-2}$ היא מטריצה סימטרית ולכן הטרייס שלה אי שלילי והייתה הפיכה הוא לא 0, כלומר נקבל כי

$$\frac{\partial Tr(\sigma^2 \cdot A_\lambda \cdot (X \cdot X^T)^{-1} \cdot A_\lambda^T)}{\partial \lambda} \big|_{\lambda=0} < 0$$

עתה נגזור את החלק השני של הביטוי

$$\begin{aligned}
\frac{\partial w^T \cdot [I_n - A_\lambda^T - A_\lambda + A_\lambda^T \cdot A_\lambda] \cdot w}{\partial \lambda} \big|_{\lambda=0} &= w^T \cdot \frac{\partial [I_n - A_\lambda^T - A_\lambda + A_\lambda^T \cdot A_\lambda]}{\partial \lambda} \cdot w \big|_{\lambda=0} \\
&= w^T \cdot \left[\frac{\partial I_n}{\partial \lambda} - \frac{\partial A_\lambda^T}{\partial \lambda} - \frac{\partial A_\lambda}{\partial \lambda} + \frac{\partial A_\lambda^T \cdot A_\lambda}{\partial \lambda} \right] \cdot w \big|_{\lambda=0} \\
&= w^T \cdot \left[\frac{\partial I_n}{\partial \lambda} - \frac{\partial A_\lambda}{\partial \lambda} - \frac{\partial A_\lambda^T}{\partial \lambda} + \frac{\partial A_\lambda^T}{\partial \lambda} \cdot A_\lambda + A_\lambda^T \cdot \frac{\partial A_\lambda}{\partial \lambda} \right] \cdot w \big|_{\lambda=0} \\
&= w^T \cdot \left[0 - \left(-(X \cdot X^T)^{-1} \right) - \left(-(X \cdot X^T)^{-1} \right)^T + \left(-(X \cdot X^T)^{-1} \right)^T \cdot A_0 + A_0^T \cdot \left(-(X \cdot X^T)^{-1} \right) \right] \cdot w \\
&= w^T \cdot \left[0 + (X \cdot X^T)^{-1} + \left[(X \cdot X^T)^{-1} \right]^T - \left[(X \cdot X^T)^{-1} \right]^T \cdot I_n - I_n \cdot (X \cdot X^T)^{-1} \right] \cdot w \\
&= w^T \cdot [0] \cdot w = 0
\end{aligned}$$

כלומר קיבלנו כי

$$\begin{aligned} \frac{\partial MSE(\hat{w}(\lambda))}{d\lambda} \Big|_{\lambda=0} &= \frac{\partial \left[w^T \cdot [I_n - A_\lambda^T - A_\lambda + A_\lambda^T \cdot A_\lambda] \cdot w + Tr \left(\sigma^2 \cdot A_\lambda \cdot (X \cdot X^T)^{-1} \cdot A_\lambda^T \right) \right]}{\partial \lambda} \Big|_{\lambda=0} \\ &= \frac{\partial \left[w^T \cdot [I_n - A_\lambda^T - A_\lambda + A_\lambda^T \cdot A_\lambda] \cdot w \right]}{\partial \lambda} \Big|_{\lambda=0} + \frac{\partial \left[Tr \left(\sigma^2 \cdot A_\lambda \cdot (X \cdot X^T)^{-1} \cdot A_\lambda^T \right) \right]}{\partial \lambda} \Big|_{\lambda=0} \\ &= 0 + \frac{\partial \left[Tr \left(\sigma^2 \cdot A_\lambda \cdot (X \cdot X^T)^{-1} \cdot A_\lambda^T \right) \right]}{\partial \lambda} \Big|_{\lambda=0} = \frac{\partial \left[Tr \left(\sigma^2 \cdot A_\lambda \cdot (X \cdot X^T)^{-1} \cdot A_\lambda^T \right) \right]}{\partial \lambda} \Big|_{\lambda=0} < 0 \end{aligned}$$

מ.ש.ל.ד.☺

(ה) צ"ל: לאמבדה ישפר טיפה אם המודל צודק

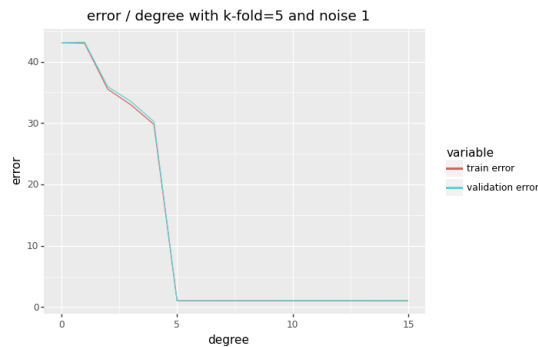
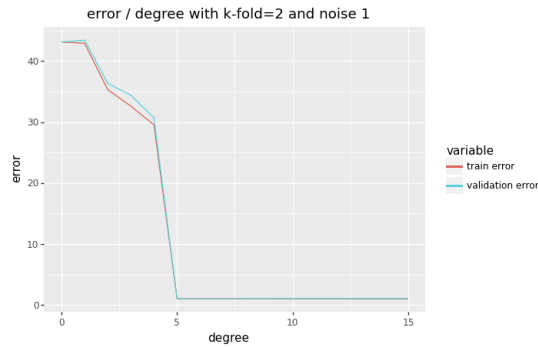
הוכחה:

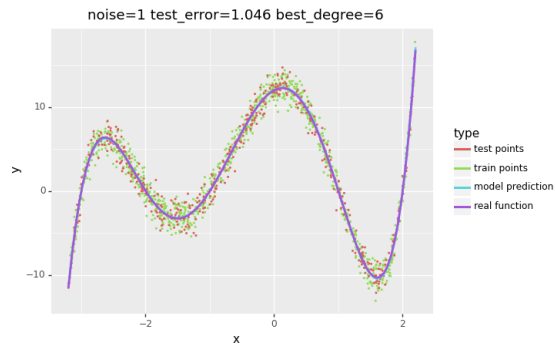
נשים לב שהנגזרת של ה- MSE שלילי עבור $\lambda = 0$, כלומר קיימת $[0, a]$ כאשר $a > 0$ שבה פונקצית השגיאה יורדת, ולכל $\lambda \in (0, a)$ יתקיים $MSE(\hat{w}(\lambda)) < MSE(\hat{w}(0))$, כלומר טיפה רגולריזציה רק תשפר את המודל.

מ.ש.ל.ה.☺

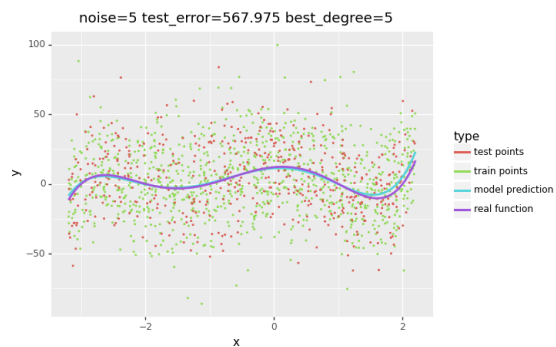
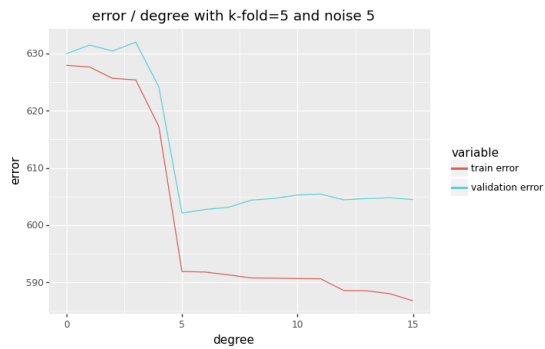
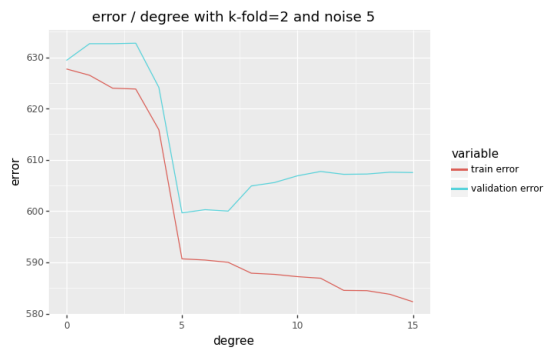
4. צ"ל: גרפים ומסקנות

הוכחה:





אפשר לשים לב שעם טיפה רעש, הוא מצליח ללמוד את הפולינום הנכון גם לולידציה וגם לאימון ואף מהדרגה הכמעט נכונה עם שגיאה אפסית, אפשר לראות כמה הפולינומים דומים בגרף השלישי שצורף.

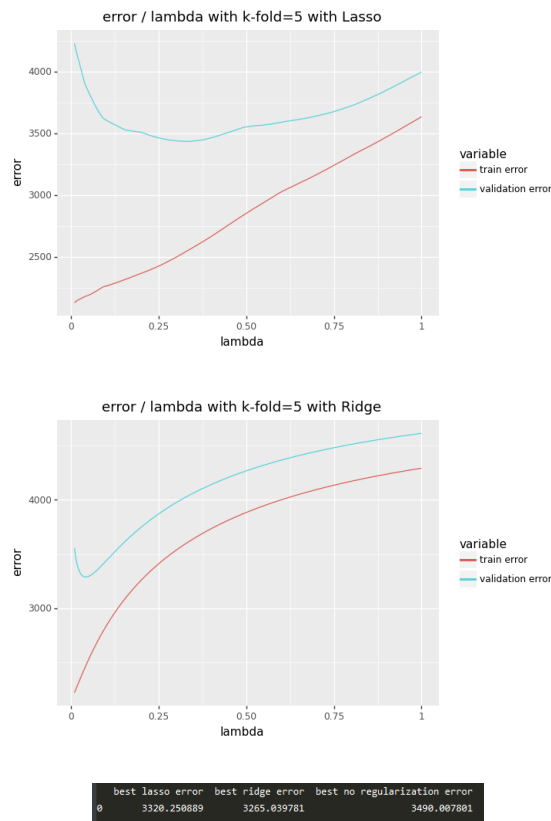


לאחר שהוספנו הרבה רעש, אפשר לראות שעבור דרגות פולינום גבוהות, הוא מתחיל ללמוד את הרעש והשגיאה על ה-*validation* עולה (כי הוא מאבד את ה-*generalization*) בזמן שהאימון יורד כי יש לו מחלקת היפותזות יותר מורכבת לתאר את הדגימות.

באותו הזמן אפשר לראות את הגרף השלישי המצורף ולראות שעל אף שהיה רעש מטורף כמו שאפשר לראות, הוא עדיין הצליח ללמוד פולינום מאוד קרוב לפולינום המקורי עם הדרגה האמיתית של הפולינום. יחסית לטיפה רעש, השגיאה גדולה מאוד, אבל ביחס למידע הניתן, זה שגיאה מדהימה.

מ.ש.ל. ☺

5. צ"ל: גרפים ומסקנות הוכחה:



תחילה אפשר לראות שגם הרידג' וגם הלאסו שיפרו את השגיאה על פני רגרסיה ללא רגולריזציה והאחד הכי טוב הוא ה-*ridge*, רוב הסיכויים בגלל שכל הפרמטרים חשובים ולא נרצה את התכונה של *lasso* שמאפסת פרמטרים. אני בהתחלה הסתכלתי על פרמטרים λ בקטע [1, 500] בקפיצות של 10, ושמתי לב שהשגיאה הכי קטנה על פני ה-*validation* היא דווקא בהתחלה בהתחלה, ולאחר חיפוש מספר פעמים, הבחנתי שהשגיאה הכי טובה על פני ה-*validation* היא בקטע [0, 1], אז החלטתי לבדוק בקטע [0.01, 1] עם 10000 דגימות בדרך ומצאתי את המינימום שחיפשתי מלכתחילה. אפשר לראות שלשניהם המינימום הוא עם λ נמוך, שזה הגיוני כי לפי מה שהוכחנו בשאלה 3.4, כשהמודל יכול לתאר את המחלקה עד כדי רעש, אז טיפה רגולריזציה יכולה רק לעזור.

מ.ש.ל. ☺