

# מבוא למערכות לומדות תרגיל 5

עמית בסקין 312259013

## חלק I

## חלק תיאורתי

### 1 ואלידציה

בשאלה זו נראה מתי לפרדיגמת "בחירת מודל" יש יתרון על פני השיטה הסטנדרטית כאשר בוחרים בין  $k$  מחלקות היפותיזות אפשריות:

$$\mathcal{H}_1 \subseteq \mathcal{H}_2 \subseteq \dots \subseteq \mathcal{H}_k$$

עבור  $\mathcal{H}_k$  סופית.

תהא  $S_{all} = \{(x_i, y_i)\}_1^m$  קבוצת דגימות. כרגיל, נרצה ללמוד היפותיזה עם שגיאת הכללה קטנה ככל האפשר.

בכיתה דנו בבעיית התאמת פולינום, שם היו לנו  $k$  מחלקות היפותיזות לבחור מביניהן:

$$\mathcal{H}_1 \subseteq \mathcal{H}_2 \subseteq \dots \subseteq \mathcal{H}_k$$

בשאלה זו נשווה בין שתי שיטות לבחירת היפותיזה:

- השיטה הסטנדרטית: לבחור את ההיפותיזה הטובה ביותר ב-  $\mathcal{H}_k$  על ידי שימוש בדגימות הנתונות ובכלל ה-  $ERM$ .
- שיטת "בחירת מודל":

- ניקח איזשהי  $\alpha \in (0, 1)$  ונחלק את  $m$  הדגימות לקבוצת אימון  $S$  מגודל  $(1 - \alpha)m$  וקבוצת ואלידציה  $V$  מגודל  $\alpha m$ , עבור  $\alpha m$  מספר שלם.

- לכל מחלקת היפותיזות  $\mathcal{H}_i$  עם  $i \in [k]$ , נמצא  $h_i \in ERM_{\mathcal{H}_i}(S)$ .

- נחזיר  $h^* \in ERM_{\mathcal{H}}(V)$  כאשר  $\mathcal{H} = \{h_i\}_1^k$ .

נניח ש-  $\mathcal{H}_k$  סופית ופונקציית ההפסד חסומה מלמעלה על ידי 1.

א. צלחוסם את שגיאת ההכללה שמתקבלת בשיטה הסטנדרטית: להוכיח שלמידה- $PAC$  אגנוסיטית של  $\mathcal{H}_k$ , מקיימת את החסם הבא:

עבור  $h^* \in ERM_{\mathcal{H}_k}(S_{all})$ , עם הסתברות של  $1 - \delta$  לכל הפחות, מתקיים:

$$L_D(h^*) \leq \min_{h \in \mathcal{H}_k} L_D(h) + \sqrt{\frac{2 \ln\left(\frac{2}{\delta} |\mathcal{H}_k|\right)}{m}}$$

פיתרון:

באמצעות חסם הופדינג, הוכחנו בתרגול שלכל  $\delta \in (0, 1)$ , בהסתברות של  $1 - \delta$  לכל הפחות, ולכל היפותיזה  $h \in \mathcal{H}_k$ , מתקיים:

$$|L_{S_{all}}(h) - L_{\mathcal{D}}(h)| \leq \sqrt{\frac{\ln\left(\frac{2}{\delta}\right)}{2m}}$$

ובפרט בהסתברות של  $1 - \frac{\delta}{|\mathcal{H}_k|}$  לכל הפחות, מתקיים:

$$|L_{S_{all}}(h) - L_{\mathcal{D}}(h)| \leq \sqrt{\frac{\ln\left(\frac{2}{\delta} |\mathcal{H}_k| \frac{2}{\delta}\right)}{2m}}$$

אז בהסתברות של  $\frac{\delta}{|\mathcal{H}_k|}$  לכל היותר מתקיים:

$$|L_{S_{all}}(h) - L_{\mathcal{D}}(h)| \geq \sqrt{\frac{\ln\left(\frac{2}{\delta} |\mathcal{H}_k|\right)}{2m}}$$

ומחסם האיחוד, ההסתברות שקיימת היפותיזה  $h$  ב- $\mathcal{H}_k$  עם הא"ש לעיל, היא לכל היותר  $|\mathcal{H}_k| \frac{\delta}{|\mathcal{H}_k|}$  קרי  $\delta$ , ולכן ההסתברות שלכל היפותיזה  $h$  ב- $\mathcal{H}_k$  מתקיים:

$$|L_{S_{all}}(h) - L_{\mathcal{D}}(h)| \leq \sqrt{\frac{\ln\left(\frac{2}{\delta} |\mathcal{H}_k|\right)}{2m}}$$

היא לכל הפחות  $1 - \delta$ , ובפרט לכל  $h \in \mathcal{H}_k$  ההסתברות שמתקיים הא"ש הנ"ל היא לכל הפחות  $1 - \delta$ .

אם כן, מאחר ששגיאת ההכללה גדולה יותר באופן טיפוס, אזי שבפרט עבור  $h^*$  מתקיים בהסתברות של  $1 - \delta$  לכל הפחות הא"ש:

$$L_{\mathcal{D}}(h^*) \leq L_{S_{all}}(h^*) + \sqrt{\frac{\ln\left(\frac{2}{\delta} |\mathcal{H}_k|\right)}{2m}}$$

אבל  $h^* \in \text{ERM}_{\mathcal{H}_k}(S_{all})$  ולכן לכל  $h \in \mathcal{H}_k$ :

$$\leq L_{S_{all}}(h) + \sqrt{\frac{\ln\left(\frac{2}{\delta} |\mathcal{H}_k|\right)}{2m}}$$

ומקיום הא"ש לעיל עבור כל  $h \in \mathcal{H}_k$ :

$$\leq L_{\mathcal{D}}(h) + 2\sqrt{\frac{\ln\left(\frac{2}{\delta} |\mathcal{H}_k|\right)}{2m}}$$

נכניס את ה-2 לתוך השורש:

$$= L_{\mathcal{D}}(h) + \sqrt{\frac{2 \ln\left(\frac{2}{\delta} |\mathcal{H}_k|\right)}{m}}$$

ומאחר שזה נכון לכל  $h \in \mathcal{H}_k$ , אזי ש-

$$L_{\mathcal{D}}(h^*) \leq \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{2 \ln\left(\frac{2}{\delta} |\mathcal{H}_k|\right)}{m}}$$

כאמור בהסתברות של  $1 - \delta$  לכל הפחות. מש"ל

ב. צ.לחסום את שגיאת ההכללה על ידי שימוש בשיטת "בחירת מודל":

נסמן:  $h' = \arg \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h)$ . יהא  $j$  מינימלי עם  $h' \in \mathcal{H}_j$ .

צ.להוכיח שבהסתברות של לפחות  $1 - \delta$  מתקיים כי:

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

פיתרון:

מהסעיף הקודם, על ידי לקיחת  $\delta/2$  והפעלת המסקנה על כל  $\mathcal{H}_i$  בנפרד, נקבל שלכל  $i \in [k]$  ו-  $h_i \in ERM_{\mathcal{H}_i}(S)$  מתקיים:

$$\mathbb{P} \left( L_{\mathcal{D}}(h_i) \leq \min_{h \in \mathcal{H}_i} L_{\mathcal{D}}(h) + \sqrt{\frac{2 \ln \left( \frac{4}{\delta} |\mathcal{H}_i| \right)}{(1-\alpha)m}} \right) \geq 1 - \frac{\delta}{2}$$

ועבור  $h^*$  שנבחרה בשלב הוואלידציה:

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

(כי משווים בין  $k$  היפותיזות)

אבל  $h' = \arg \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h)$ , ולכן בהסתברות של לכל הפחות  $1 - \delta$ , מתקיים:

$$L_{\mathcal{D}}(h^*) \leq L_{\mathcal{D}}(h') + \sqrt{\frac{2 \ln \left( \frac{4}{\delta} k \right)}{\alpha m}}$$

ומקיום הא"ש הנ"ל עבור  $\mathcal{H}_j$ :

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

אבל  $\min_{h \in \mathcal{H}_j} L_{\mathcal{D}}(h) = \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h)$  ולכן:

$$= \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{2 \ln \left( \frac{4}{\delta} |\mathcal{H}_j| \right)}{(1-\alpha)m}} + \sqrt{\frac{2 \ln \left( \frac{4}{\delta} k \right)}{\alpha m}}$$

וזאת כאמור בהסתברות של  $1 - \delta$  לכל הפחות. מש"ל

ג. צ.להראות ששני החסמים אינם בני-שוואה: כלומר צ.לתאר מקרה שבו השיטה הסטנדרטית עדיפה ומקרה הפוך.

הוכחה:

עבור השיטה הסטנדרטית קיבלנו:

$$L_{\mathcal{D}}(h^*) \leq \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{2 \ln \left( \frac{2}{\delta} |\mathcal{H}_k| \right)}{m}}$$

ועבור שיטת "בחירת מודל":

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

אז נרצה לתת דוגמה לשני הבאים:

$$\min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{2 \ln \left( \frac{2}{\delta} |\mathcal{H}_k| \right)}{m}} < \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{2 \ln \left( \frac{4}{\delta} |\mathcal{H}_j| \right)}{(1-\alpha)m}} + \sqrt{\frac{2 \ln \left( \frac{4}{\delta} k \right)}{\alpha m}}$$

ו-

$$\min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{2 \ln \left( \frac{4}{\delta} |\mathcal{H}_j| \right)}{(1-\alpha)m}} + \sqrt{\frac{2 \ln \left( \frac{4}{\delta} k \right)}{\alpha m}} < \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{2 \ln \left( \frac{2}{\delta} |\mathcal{H}_k| \right)}{m}}$$

קרי:

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

- אם מתקיים  $j = k$  ו-  $\alpha = \frac{1}{2}$ , אז עבור הא"ש הראשון מתקיים:

$$\begin{aligned} \sqrt{\frac{2 \ln \left( \frac{4}{\delta} |\mathcal{H}_j| \right)}{0.5m}} + \sqrt{\frac{2 \ln \left( \frac{4}{\delta} k \right)}{0.5m}} &= \\ = 2\sqrt{\frac{\ln \left( \frac{4}{\delta} |\mathcal{H}_k| \right)}{m}} + 2\sqrt{\frac{\ln \left( \frac{4}{\delta} k \right)}{m}} \end{aligned}$$

אבל:

$$2\sqrt{\ln \left( \frac{4}{\delta} |\mathcal{H}_k| \right)} + 2\sqrt{\ln \left( \frac{4}{\delta} k \right)} > \sqrt{2 \ln \left( \frac{2}{\delta} |\mathcal{H}_k| \right)}$$

ולכן מתקיים הא"ש הראשון במקרה הזה.

- עבור הא"ש השני:

אם  $|\mathcal{H}_j| < e^8 \delta$  ו-  $k, |\mathcal{H}_k| > \frac{1}{2} \delta e^{50}$ ,  $\alpha = 0.5$ , אז:

$$\begin{aligned} \sqrt{\frac{2 \ln \left( \frac{4}{\delta} |\mathcal{H}_j| \right)}{(1-\alpha)m}} + \sqrt{\frac{2 \ln \left( \frac{4}{\delta} k \right)}{\alpha m}} &< \\ < 2\sqrt{\frac{8}{0.5m}} = 2\sqrt{\frac{16}{m}} = \frac{8}{\sqrt{m}} \end{aligned}$$

ואילו:

$$\sqrt{\frac{2 \ln \left( \frac{2}{\delta} |\mathcal{H}_k| \right)}{m}} > \sqrt{\frac{2 \ln (e^{50})}{m}} = \frac{10}{\sqrt{m}}$$

ולכן מתקיים הא"ש השני במקרה הזה.

במקרה המיוחד של מטריצת תכנון אורתוגונלית,  $XX^t = Id$ , ניתן לגזור נוסחאות סגורות עבור Ridge, Lasso ו-Best Subset, כפונקציה של  $LS$ .

תהא  $\lambda > 0$ . נסמן ב-  $\hat{w}_\lambda^{LS}$ ,  $\hat{w}_\lambda^{ridge}$ ,  $\hat{w}_\lambda^{lasso}$ ,  $\hat{w}_\lambda^{subset}$  את הפתרונות עבור רגרסיה לינארית סטנדרטית, Ridge, Lasso ו-Best Subset בהתאמה.

ניזכר בפונקציות הבאות:

$$\eta_\lambda^{soft}(x) = \begin{cases} x - \lambda & x \geq 0 \\ 0 & |x| < \lambda \\ x + \lambda & x \leq -\lambda \end{cases}$$

$$\eta_\lambda^{hard}(x) = \mathbf{1}[|x| \geq \lambda] \cdot x$$

צ. להוכיח ש-

א.

$$\hat{w}_\lambda^{ridge} = \frac{\hat{w}^{LS}}{1 + \lambda}$$

ב.

$$\hat{w}_\lambda^{subset} = \eta_{\sqrt{\lambda}}^{hard}(\hat{w}^{LS})$$

הוכחה:

א. בתרגול ראינו ש-

$$\hat{w}_\lambda^{ridge} = (XX^t + \lambda I)^{-1} Xy$$

ומכך ש-  $XX^t = I$ :

$$\hat{w}_\lambda^{ridge} = (I + \lambda I)^{-1} Xy = \frac{1}{1 + \lambda} Xy$$

אבל:

$$\hat{w}^{LS} = (XX^t)^{-1} Xy = Xy$$

ולכן:

$$\hat{w}_\lambda^{ridge} = \frac{\hat{w}^{LS}}{1 + \lambda}$$

כנדרש.

ב. בדומה לטענה 2 בתרגול, נכתוב:

$$\begin{aligned} f_{l_0}(w) &= \|y - X^t w\|_2^2 + \lambda \|w\|_0 = \\ &= \|y\|_2^2 - 2y^t X^t w + w^t X X^t w + \lambda \|w\|_0 = \\ &= \|y\|_2^2 + (w^t - 2\hat{w}^t) w + \lambda \|w\|_0 = \\ &= \|y\|_2^2 + \sum_1^d (w_i^2 - 2\hat{w}_i w_i + \lambda \|w_i\|_0) \end{aligned}$$

אם  $w_i = 0$  אז הנסכם מתאפס ואם  $w_i \neq 0$  אז מתקבל בנסכם:

$$w_i^2 - 2\hat{w}_i w_i + \lambda$$

נגזור לפי  $w_i$  ונשווה לאפס:

$$w_i = \hat{w}_i$$

קרי:

$$\hat{w}_i^2 - 2\hat{w}_i^2 + \lambda = -\hat{w}_i^2 + \lambda$$

וזה קטן מאפס אם  $|\hat{w}_i| \geq \sqrt{\lambda}$ , ולכן  $\hat{w}_i \neq 0$  אם  $|\hat{w}_i| \geq \sqrt{\lambda}$ .

ובסה"כ קיבלנו:

$$\hat{w}_\lambda^{subset} = \eta_{\sqrt{2\lambda}}^{hard}(\hat{w}^{LS})$$

### 3 רגולריזציה

בשאלה זו נראה שלמרות ש-Ridge מוסיפה  $bias$  להערכה של  $w$ , היא עדיין מקטינה את ה- $MSE$ :

תהא  $X \in \mathbb{R}_{d \times m}$  מטריצת רגרסיה קבועה כך ש- $XX^t$  הפיכה.

יהא

$$\hat{w}(\lambda) = \arg \min_w \left( \|y - X^t w\|_2^2 + \lambda \|w\|_2^2 \right)$$

הפיתרון של Ridge וכן יהא  $\hat{w}(\lambda = 0) \equiv \hat{w}$  הפיתרון של  $LS$ .

• נניח שהמודל הלינארי הוא נכון, קרי:

$$y = X^t \hat{w} + \epsilon$$

כאשר  $\epsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$

• נזכר כי במקרה זה מתקיים  $\mathbb{E}[\hat{w}] = \hat{w}$

א. צ. להראות שמתקיים כי

$$\hat{w}(\lambda) = A_\lambda \hat{w}$$

כאשר

$$A_\lambda = (XX^t + \lambda Id)^{-1} XX^t$$

הוכחה:

$$A_\lambda \hat{w} = (XX^t + \lambda Id)^{-1} XX^t \cdot (XX^t)^{-1} Xy = (XX^t + \lambda Id)^{-1} Xy$$

ואכן ראינו בהרצאה כי:

$$Xy = (XX^t + \lambda I) \hat{w}(\lambda)$$

ומאחר ש-  $(XX^t + \lambda I)$  הפיכה אזי שקיבלנו את הנדרש.

ב. צ. להסיק מהסעיף הקודם הפיתרון של Ridge הוא עם  $bias$  לכל  $\lambda > 0$ , קרי להראות שאם  $\lambda > 0$  אז  $\mathbb{E}[\hat{w}(\lambda)] \neq \hat{w}$ .

הוכחה:

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

ג. צ. להראות ש-

$$\text{Var}(\hat{w}(\lambda)) = \sigma^2 A_\lambda (XX^t)^{-1} A_\lambda^t$$

הוכחה:

נשתמש בעובדה שעבור מטריצה  $B$  שאינה אקראית ווקטור אקראי  $z$ , מתקיים:

$$\text{Var}(Bz) = B \text{Var}(z) B^t$$

$$\text{Var}(\hat{w}) = \sigma^2 (XX^t)^{-1}$$

קרי:

$$\text{Var}(\hat{w}(\lambda)) = \text{Var}(A_\lambda \hat{w}) = A_\lambda \text{Var}(\hat{w}) A_\lambda^t =$$

$$= A_\lambda \sigma^2 (X X^t)^{-1} A_\lambda^t = \sigma^2 A_\lambda (X X^t)^{-1} A_\lambda^t$$

ד. צלגזור ביטויים מפורשים עבור ה-  $bias$  (בריבוע) והשונות של  $\hat{w}(\lambda)$ , כפונקציה של  $\lambda$ , קרי, לכתוב פירוק של bias-variance עבור ה-  $MSE$  של  $\hat{w}(\lambda)$ .

כמו כן צלהראות באמצעות גזירה ש-

$$\frac{d}{d\lambda} MSE(\lambda) |_{\lambda=0} = \frac{d}{d\lambda} bias^2(\lambda) |_{\lambda=0} + \frac{d}{d\lambda} Var(\lambda) |_{\lambda=0} < 0$$

פיתרון:

$$MSE(\lambda) = (A_\lambda \hat{w} - \hat{w}) (A_\lambda \hat{w} - \hat{w})^t + \sigma^2 A_\lambda (X X^t)^{-1} A_\lambda^t =$$

$$= (A_\lambda - I) \hat{w} ((A_\lambda - I) \hat{w})^t + \sigma^2 A_\lambda (X X^t)^{-1} A_\lambda^t =$$

$$= (A_\lambda - I) \hat{w} \hat{w}^t (A_\lambda - I)^t + \sigma^2 A_\lambda (X X^t)^{-1} A_\lambda^t =$$

$$= (A_\lambda - I) \hat{w} \hat{w}^t (A_\lambda - I) + \sigma^2 A_\lambda (X X^t)^{-1} A_\lambda^t =$$

$$= \hat{w} \hat{w}^t \left( (X X^t + \lambda Id)^{-1} X X^t - I \right)^2$$

$$+ \sigma^2 (X X^t + \lambda Id)^{-1} X X^t (X X^t)^{-1} X X^t \left( (X X^t + \lambda Id)^{-1} \right)^t =$$

$$= \hat{w} \hat{w}^t \left( (X X^t + \lambda Id)^{-1} X X^t - I \right)^2 +$$

$$+ \sigma^2 (X X^t + \lambda Id)^{-1} X X^t \left( (X X^t + \lambda Id)^{-1} \right)$$

ועתה:

$$\frac{d}{d\lambda} \hat{w} \hat{w}^t \left( (X X^t + \lambda Id)^{-1} X X^t - I \right)^2$$

$$= -2 \hat{w} \hat{w}^t (X X^t + \lambda Id)^{-1} \frac{d}{d\lambda} (-X X^t + \lambda Id) (X X^t + \lambda Id)^{-1} \left( (X X^t + \lambda Id)^{-1} X X^t - I \right) =$$

$$= -2 \hat{w} \hat{w}^t (X X^t + \lambda Id)^{-2} \left( (X X^t + \lambda Id)^{-1} X X^t - I \right) =$$



ולכן:

$$\begin{aligned} & \frac{d}{d\lambda} \text{bias}^2(\lambda) |_{\lambda=0} = \\ & = -2\hat{w}\hat{w}^t (XX^t)^{-2} \left( (XX^t)^{-1} XX^t - I \right) = \\ & = -2\hat{w}\hat{w}^t (XX^t)^{-2} (I - I) = 0 \end{aligned}$$

ומצד שני:

$$\begin{aligned} & \left( \frac{d}{d\lambda} \sigma^2 (XX^t + \lambda Id)^{-1} XX^t \left( (XX^t + \lambda Id)^{-1} \right) \right) |_{\lambda=0} = \\ & = -\sigma^2 (XX^t + \lambda Id)^{-1} \frac{d}{d\lambda} [\sigma^2 (XX^t + \lambda Id)] \sigma^2 (XX^t + \lambda Id)^{-1} XX^t \left( (XX^t + \lambda Id)^{-1} \right) - \\ & -\sigma^2 (XX^t + \lambda Id)^{-1} XX^t \left( (XX^t + \lambda Id)^{-1} \right) \frac{d}{d\lambda} \left[ (XX^t)^{-1} (XX^t + \lambda Id) \right] XX^t \left( (XX^t + \lambda Id)^{-1} \right) = \\ & = -\sigma^2 (XX^t + \lambda Id)^{-1} \sigma^2 \sigma^2 (XX^t + \lambda Id)^{-1} XX^t (XX^t + \lambda Id)^{-1} - \\ & -\sigma^2 (XX^t + \lambda Id)^{-1} XX^t (XX^t + \lambda Id)^{-1} (XX^t)^{-1} XX^t (XX^t + \lambda Id)^{-1} = \\ & = -\sigma^6 (XX^t + \lambda Id)^{-2} XX^t (XX^t + \lambda Id)^{-1} - \\ & -\sigma^2 (XX^t + \lambda Id)^{-1} XX^t (XX^t + \lambda Id)^{-2} \end{aligned}$$

ולכן:

$$\begin{aligned} & \frac{d}{d\lambda} \text{Var}(\lambda) |_{\lambda=0} = \\ & = -\sigma^6 (XX^t)^{-2} XX^t (XX^t)^{-1} - \\ & -\sigma^2 (XX^t)^{-1} XX^t (XX^t)^{-2} = \end{aligned}$$

$$= -\sigma^6 (XX^t)^{-2} - \sigma^2 (XX^t)^{-2} =$$

$$= -\sigma^6 \left( (XX^t)^{-1} \right)^2 - \sigma^2 \left( (XX^t)^{-1} \right)^2$$

ובסה"כ קיבלנו:

$$\frac{d}{d\lambda} \text{MSE}(\lambda) |_{\lambda=0} = \frac{d}{d\lambda} \text{bias}^2(\lambda) |_{\lambda=0} + \frac{d}{d\lambda} \text{Var}(\lambda) |_{\lambda=0} =$$

$$= \frac{d}{d\lambda} \text{Var}(\lambda) |_{\lambda=0} = -\sigma^6 \left( (XX^t)^{-1} \right)^2 - \sigma^2 \left( (XX^t)^{-1} \right)^2 < 0$$

כנדרש.

ה. צ. להסיק שאם המודל הלינארי הוא נכון, אז מעט רגורליזציה של *Ridge*, מסייעת להורדת ה- $MSE$ . הוכחה:  
מהסעיף הקודם מתקבל שבסביבת  $\lambda = 0$ , פונקציית ה- $MSE$  יורדת, ומכאן שיש  $\lambda$  כך ש- $MSE(\lambda) < MSE(0)$ .

## חלק II

## חלק מעשי

### 4 k-Fold Cross Validation על התאמת פולינום

יהיו  $\mathcal{X}, \mathcal{Y} \subseteq \mathbb{R}$ . נניח שהקשר בין  $\mathcal{X}$  ו- $\mathcal{Y}$  ניתן לתיאור על ידי פולינום מגדרה  $d \in [15]$ .

לכל  $d \in [15]$ , נסמן ב- $\mathcal{H}_d$  את מחלקת הפולינומים מדרגה  $d$ .

המשימה היא לאמן כל אחת מהמחלקות על קבוצת האימון ולבצע וואלידציה על 15 ההיפותיזות שהתקבלו, בשביל לבחור את הפלט הסופי.

לבסוף, נבחן את הביצועים של הפרדיקטור שהתקבל, על קבוצת המבחן.

א. צ. לייצר דאטא באופן הבא:

i. ניקח  $\mathcal{X} = [-3.2, 2.2]$  ונדגום מתוכו בהתפלגות אחידה.

ii. נקבע את הקשר בין  $y$  ו- $x$  על ידי:

$$y(x) = f(x) + \epsilon$$

כאשר:

$$f(x) = \prod_{i=-2}^3 (x+i)$$

ו- $\epsilon \sim \mathcal{N}(0, \sigma^2)$ .

iii. ניקח  $\sigma = 1$  ונייצר  $m = 1,500$  דגימות.

iv. נחלק את  $m$  הדגימות לשתי קבוצות: 1,000 דגימות עבור קבוצת האימון וקבוצת הוואלידציה, נסמנה  $D$ , ו- 500 עבור קבוצת המבחן, נסמנה  $T$ .

ב. נחלק את  $D$  לשתי קבוצות של 500 דגימות כל אחת: קבוצה אחת עבור קבוצת האימון  $S$  וקבוצה שנייה עבור קבוצת הוואלידציה  $V$ .

נאמן כל אחת מהמחלקות באמצעות  $S$  בשביל לקבל היפותיזה  $h_d$ , לכל  $d \in [15]$ , שממזער את ה- $loss$  על פני  $S$ .

ג. צלבצע וואלידציה בעבור  $\{h_i\}_1^{15}$  בשביל לקבל  $h^*$  שממזער את ה- $loss$  על פני קבוצת הוואלידציה,  $V$ :

תוצאה:

The best degree with two folds for polynomial regression on the given data is: 7

כלומר, הרעש משפיע והחלוקה של קבוצת האימון לשני חלקים אינה מספיקה על מנת להתגבר עליו. אך התוצאה הטובה ביותר מתקבלת עבור פולינום ממעלה שבע בעוד שהפולינום המקורי הוא ממעלה חמש.

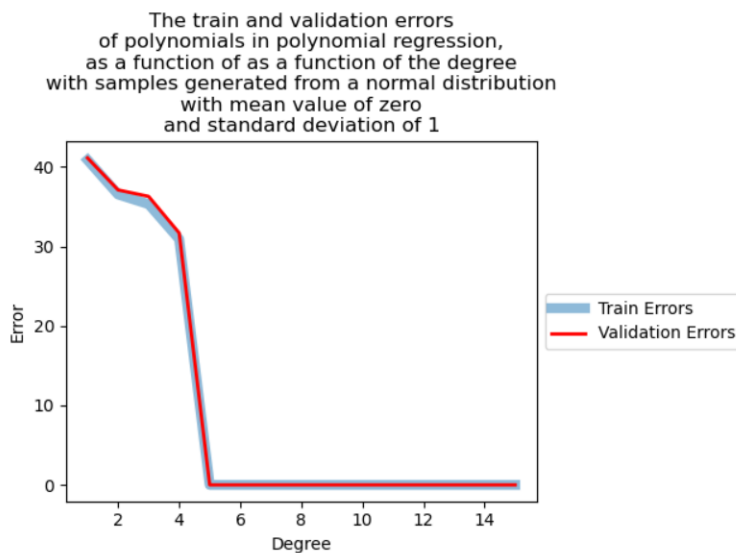
ד. סעיפים א עד ג הם k-fold cross validation עם  $k = 2$ .

צלבצע k-fold cross validation אבל הפעם עם  $k = 5$  על הקבוצה  $D$ .

ה. צלשרטט את השגיאות על פני קבוצת האימון וקבוצת הוואלידציה, כממוצע על פני  $k$ -החלקים, של הפולינומים מדרגות  $d \in [15]$ , שכפונקציה של  $d$ . צלבדוק עבור איזו דרגה מתקבלת השגיאה הנמוכה ביותר, נסמנה  $d^*$ .

פיתרון:

הגרף להלן:



The best degree with cross validation for polynomial regression on the given data is: 5

כלומר התוצאה הטובה ביותר היא עבור פולינום מדרגה 5. כמו כן יש הבדל מזערי בין השגיאה על קבוצת המבחן ובין השגיאה על קבוצת הוואלידציה, קרי יש מעט מאוד *overfit*, כלומר למדנו מתוך מידע מספיק גדול.

ו. צלבצע  $ERM_{\mathcal{H}_{d^*}}$  על הקבוצה  $D$ : למצוא פולינום מדרגה  $d^*$  בעל שגיאה מינימלית על הדגימות הללו, נסמנו ב-  $h^*$ .

ז. צלבחון את הביצועים של  $h^*$  על קבוצת המבחן  $T$ : לחשב את השגיאה של  $h^*$  על פני  $T$ . צלבדוק האם השגיאה שונה בהרבה מהשגיאה שמצאנו בסעיף הקודם.

פיתרון:

The train error is: 0.9813278322975978

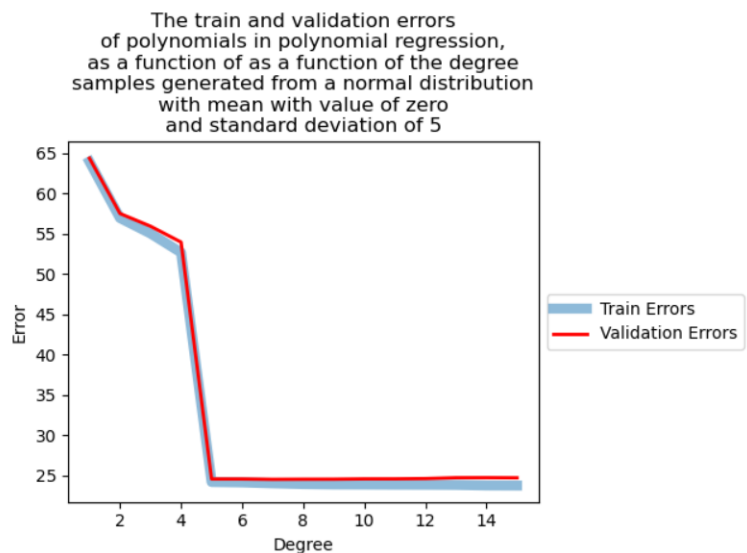
The test error is: 0.9726685025411308

The difference between the train error and the test error is: 0.008659329756466994

כפי שניתן לראות, ההבדל הוא מזערי, כלומר ה- *overfit* הינו קטן ביותר, כלומר ה- Cross Validation עושה את העבודה.

ח. צ. לחזור על השלבים הקודמים עבור  $\sigma = 5$  ולבדוק מה השתנה.

פיתרון:



The best degree with two folds for polynomial regression on the given data is: 6

The best degree with cross validation for polynomial regression on the given data is: 7

The train error is: 24.09464342161084

The test error is: 21.74461657093929

The difference between the train error and the test error is: 2.3500268506715507

הפעם הרעש יותר משמעותי, ולכן השגיאה גדלה, וכן ה- Cross Validation לא מצליח להתגבר עליו.

## 5 k-fold ורגולריזציה

בשאלה זו נערוך השוואה בין רגולריזציות *Lasso* ו- *Ridge*. נשתמש בדאטא שמעריך את מידת הגלוקוז בדם של מטופלי סוכרת.

א. נטען את הדאטא.

ב. רגולריזציה היא שימושית כאשר קבוצת האימון קטנה. אז ניצור את קבוצת האימון להיות  $m = 50$  הדגימות הראשונות, והשאר יהיו את קבוצת המבחן.

ג. צ. לבצע את הבאים פעמיים: פעם אחת עבור *Ridge* ופעם שנייה עבור *Lasso*:

i. צ. לערוך k-fold cross-validation מעל קבוצת האימון עם  $k = 5$ , תוך שימוש בערכים שונים עבור פרמטר הרגולריזציה  $\lambda$ : צ. לחקור את הטווח של הערכים האפשריים.

ii. צ. לציין איזה ערכים של  $\lambda$  בחרתי לבדיקה על פני קבוצת הוואלידציה, ולהסביר מדוע.

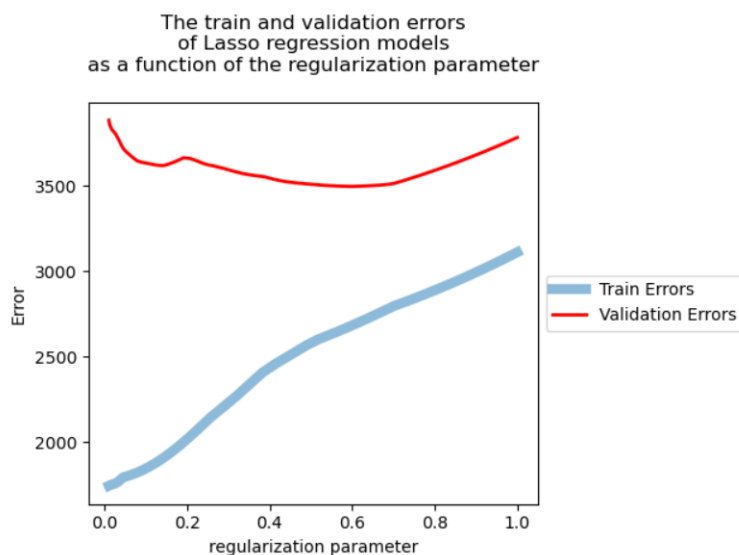
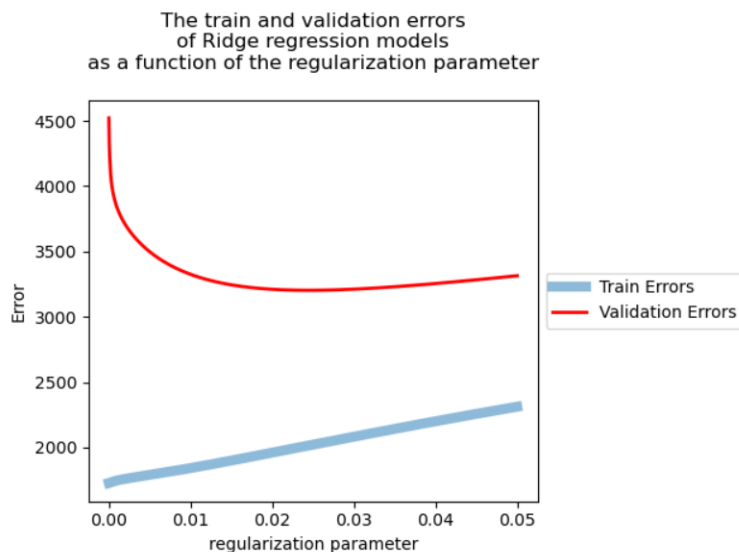
תשובה:

עבור *Ridge* בחרתי בערכים בין 0 ל- 0.05 על מנת שיהיה ניתן לראות במדויק את האזור שבו השגיאה על ה- Validation Set עובר מירידה לעלייה, שם פרמטר הרגולריזציה אופטימלי (באזור ה- 0.025).

עבור  $Lasso$  בחרתי בערכים בין 0.01 ל-1, ראשית על מנת שיהיה ניתן לראות שערכים שקרובים מדי לאפס אינם מטיבים עמו ושנית על מנת שיהיה ניתן לראות במדויק את האזור שבו השגיאה על ה- Validation Set עובר מירידה לעלייה, שם פרמטר הרגורליזציה אופטימלי (באזור ה- 0.6).

ד. צלשרטט את השגיאות על פני קבוצות האימון והואלידציה (ממוצע השגיאות על פני  $k$ -החלקים) כפונקציה של  $\lambda$ , על פני הטווח שבחרתי.

פיתרון:



ה. צלמצוא את ה-  $\lambda$  הטובה ביותר עבור  $Ridge$  וכן-ל עבור  $Lasso$ .

תשובה:

Best regularization parameter for Ridge: 0.02442442442442426  
Best regularization parameter for Lasso: 0.5966666666666667

ו. צלחשב את השגיאה של ההיפותיזות הבאות, על פני קבוצת המבחן:

i. ההיפותיזה הטובה ביותר של  $Ridge$ .

ii. ההיפותיזה הטובה ביותר של  $Lasso$ .

iii. רגרסיה לינארית.

תשובה:

The test error for the best Ridge model is: 3245.464541008676

The test error for a linear regression model is: 3612.249688324901

The test error for the best Lasso model is: 3640.7959659482776

ז. צ.לבדוק לאיזו היפותיזה יש את השגיאה הקטנה ביותר, ולבדוק האם הרגולריזציה עזרה לשפר את התוצאות על פני קבוצת המבחן, בהשוואה לתוצאות שהתקבלו ללא רגולריזציה על אותו הדאטא.

תשובה:

ל-  $Ridge$  יש את השגיאה הקטנה ביותר, כלומר הרגולריזציה עזרה לשפר את התוצאות של פני קבוצת המבחן במקרה של  $Ridge$ , אך במקרה של  $Lasso$  השגיאה גדלה, קרי הרגולריזציה לא תרמה במקרה של  $Lasso$ .