

Contents

1	Ex5 Answers.pdf	2
2	ex5.py	16

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 αm , עבור 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:

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

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 $\frac{\delta}{|\mathcal{H}_k|}$ קרי δ , 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 טיפוסי, 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- $h_i \in \text{ERM}_{\mathcal{H}_i}(S)$ 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-השוואה: 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- $k, |\mathcal{H}_k| > \frac{1}{2} \delta e^{50}$, ו- $\alpha = 0.5$, 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.

2 תכנון אורתוגונלי

במקרה המיוחד של מטריצת תכנון אורתוגונלית, קרי $XX^t = Id$, 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:

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

זה 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 רגיוולריזציה

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

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. צלהראות 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 λ , 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:

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

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 רגיורליזציה של $Ridge$, מסייעת להורדת ה- MSE . הוכחה: מהסעיף הקודם מתקבל שבסביבת $\lambda = 0$, פונקציית ה- MSE יורדת, ומכאן שיש λ 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 ואלידציה 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. אכן התוצאה הטובה ביותר מתקבלת עבור פולינום ממעלה שבע בעוד שהפולינום המקורי הוא ממעלה חמש.

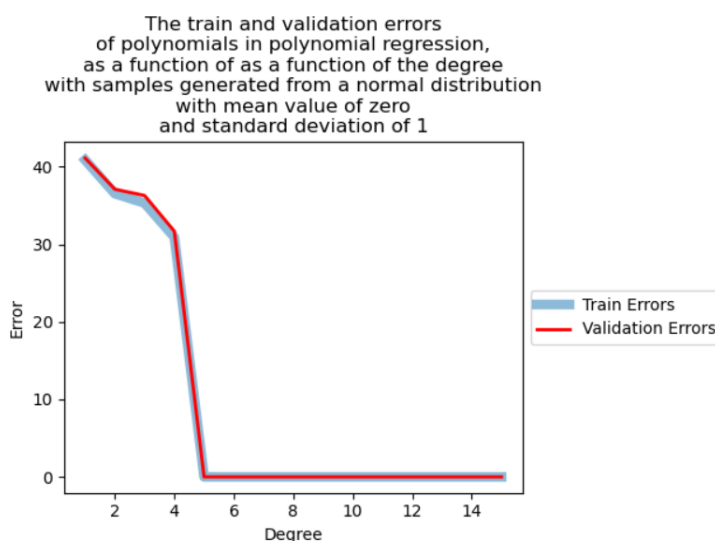
ד. סעיפים א עד ג הם k-fold cross validation עם $k = 2$.

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 $ERM_{H_{d^*}}$ 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

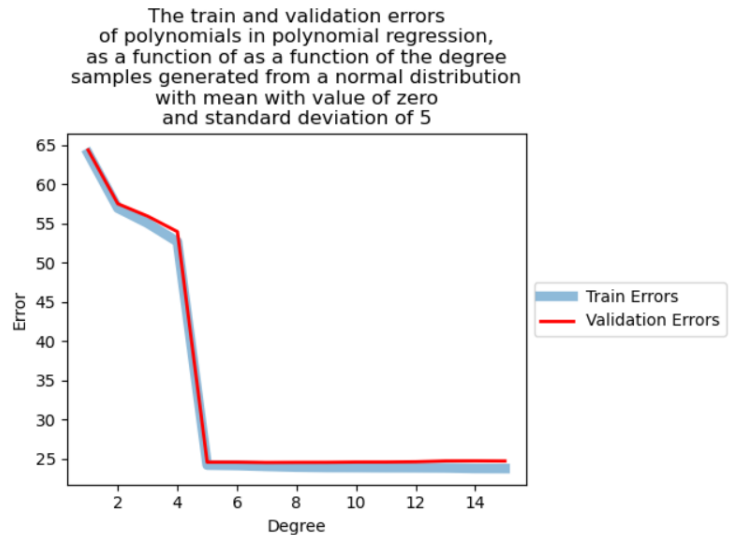
The test error is: 0.9726685025411308

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. צלחזור 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 ורגיולריזציה

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

א. נטען את הדאטא.

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

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

i. צלערוך k-fold cross-validation מעל קבוצת האימון עם $k = 5$, תוך שימוש בערכים שונים עבור פרמטר הרגיולריזציה λ : צלחקור את הטווח של הערכים האפשריים.

ii. צלציין איזה ערכים של λ בחרתי לבדיקה על פני קבוצת הוואלידציה, ולהסביר מדוע.

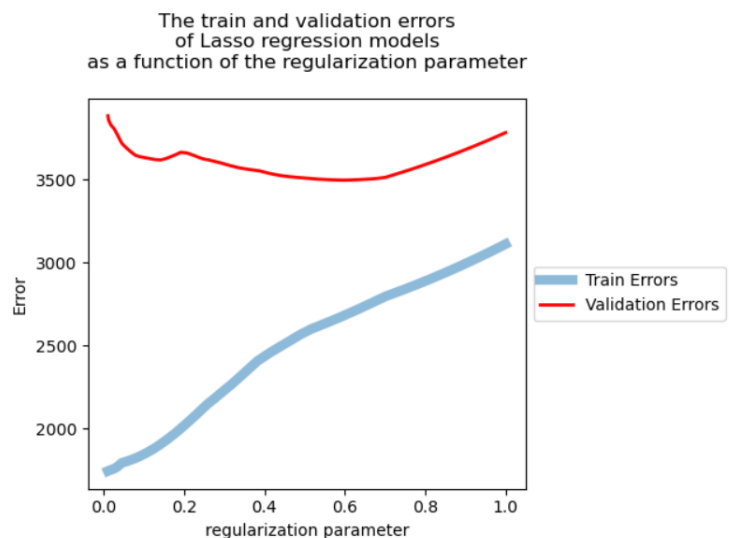
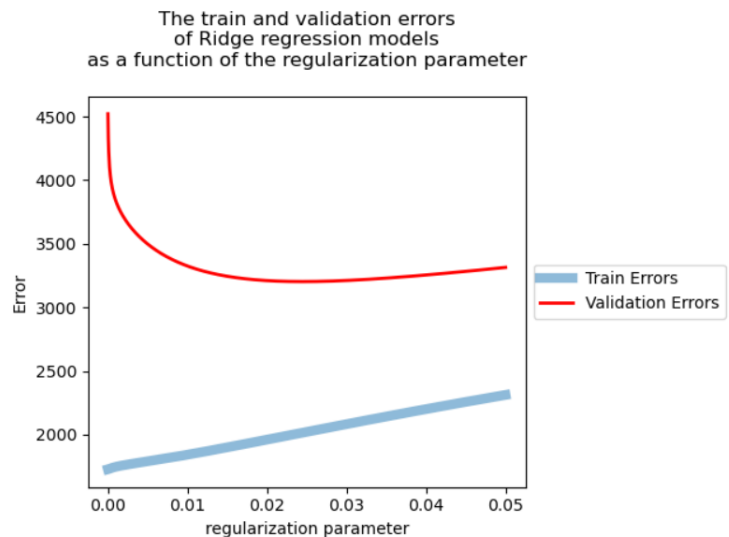
תשובה:

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 הרגיולריזציה 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. צ.לשרטט את השגיאות על פני קבוצות האימון והואלידציה (ממוצע השגיאות על פני k -החלקים) כפונקציה של λ , על פני הטווח שבחרתי.

פיתרון:



ה. צ.למצוא את ה- λ הטובה ביותר עבור *Ridge* וכן"ל עבור *Lasso*.

תשובה:

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

ו. צ.לחשב 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 הרגילריזציה עזרה לשפר את התוצאות על פני קבוצת המבחן, בהשוואה לתוצאות שהתקבלו ללא רגילריזציה על אותו הדאטא.

תשובה:

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

2 ex5.py

```
1 import numpy as np
2 from sklearn.preprocessing import PolynomialFeatures
3 from sklearn.pipeline import make_pipeline
4 from matplotlib import pyplot as plt
5 from sklearn import datasets
6 from sklearn.linear_model import Ridge, Lasso
7 from sklearn.linear_model import LinearRegression
8
9
10 LOW = -3.2
11 HIGH = 2.2
12 SIZE = 1500
13 TRAIN_FACTOR = 1000
14 MAX_DEGREE = 15
15 FOLDS_AMOUNT = 5
16 DEGS = range(1, MAX_DEGREE + 1)
17 FIRST_ITERATION_SIGMA = 1
18 SECOND_ITERATION_SIGMA = 5
19 RIDGE_PARAMS = np.linspace(0, 0.05, 1000)
20 LASSO_PARAMS = np.linspace(0.01, 1, 1000)
21
22
23 def default_labels_func():
24     """
25     The labels function for the samples of question four
26     :param sigma: The standard deviation of which to create noise from
27     :return: The function
28     """
29     return lambda x: (x + 3) * (x + 2) * (x + 1) * (x - 1) * (x - 2)
30
31
32
33 def generate_samples(low, high, size):
34     """
35     Generates samples uniformly from an interval
36     :param low: The lower bound of the interval
37     :param high: The upper bound of the interval
38     :param size: The amount of samples to generate
39     :return: The samples
40     """
41     return np.random.uniform(low, high, size)
42
43
44 def generate_labels(f, all_samples):
45     """
46     Generates the labels of all the samples that were already generated
47     :param f: The label function
48     :param all_samples: The samples to get their labels
49     :return: The labels
50     """
51     return np.apply_along_axis(f, 0, all_samples)
52
53
54 def add_noise(labels, sigma, mean=0):
55     return labels + np.random.normal(mean, sigma, len(labels))
56
57
58 def generate_data(sigma):
59     """
```



```

60     Generates the data for question four
61     :param sigma: The standard deviation for the noise of the labels
62     :return: Train samples, train labels, test samples, test labels
63     """
64     X = generate_samples(LOW, HIGH, SIZE)
65     y = generate_labels(default_labels_func(), X)
66     y = add_noise(y, sigma)
67     D, T = split_data(X, y, TRAIN_FACTOR)
68     all_train_samples = D[0]
69     all_train_samples = all_train_samples.reshape(-1, 1)
70     all_train_labels = D[1]
71     test_samples = T[0]
72     test_samples = test_samples.reshape(-1, 1)
73     test_labels = T[1]
74     return all_train_samples, all_train_labels, test_samples, test_labels
75
76
77 def split_data(all_samples, all_labels, train_factor):
78     """
79     Splits the data into train set and test set
80     :param all_samples: The samples to split
81     :param all_labels: The labels to split
82     :param train_factor: A number which specifies how many samples should
83     the train set contain
84     :return: The train set and test set
85     """
86     train_set = all_samples[:train_factor], all_labels[:train_factor]
87     test_set = all_samples[train_factor:], all_labels[train_factor:]
88     return train_set, test_set
89
90
91 def get_k_folds(k, train_samples, train_labels):
92     """
93     Splits the train samples into folds
94     :param k: The amount of folds to split the data into
95     :param train_samples: The samples to split
96     :param train_labels: The labels to split
97     :return: A list of the folds of the samples, where each fold should play
98     the role of the validation set, and also returns a list of the samples
99     where the i'th item consists of all the samples but the i'th fold,
100    which should play the role of a train set that matches the i'th fold as
101    the validation set"""
102     total_amount = len(train_samples)
103     amount_per_fold = total_amount // k
104     validation_excluded = []
105     folds = []
106
107     for i in range(1, k + 1):
108         a = train_samples[: (i - 1) * amount_per_fold]
109         b = train_samples[i * amount_per_fold:]
110         curr_samples = np.concatenate((a, b))
111         a = train_labels[: (i - 1) * amount_per_fold]
112         b = train_labels[i * amount_per_fold:]
113         curr_labels = np.concatenate((a, b))
114         validation_excluded.append((curr_samples, curr_labels))
115
116     counter = 0
117     for i in range(k):
118         fold_samples = train_samples[counter:counter + amount_per_fold]
119         fold_labels = train_labels[counter:counter + amount_per_fold]
120         folds.append((fold_samples, fold_labels))
121         counter += amount_per_fold
122
123     return folds, validation_excluded
124
125
126 def two_folds_get_single_model(model, train_samples, train_labels,
127                                validation_samples, validation_labels):

```

```

128     """
129     Gets a single model for the two folds version
130     :param model: The model to train
131     :param train_samples: Train samples
132     :param train_labels: Train labels
133     :param validation_samples: Validation samples
134     :param validation_labels: Validation labels
135     :return: The trained model, its train error and its validation error
136     """
137     model.fit(train_samples, train_labels)
138     y_hat_train = model.predict(train_samples)
139     train_error = np.mean((train_labels - y_hat_train) ** 2)
140     y_hat_validation = model.predict(validation_samples)
141     validation_error = np.mean((validation_labels - y_hat_validation) ** 2)
142
143     return model, train_error, validation_error
144
145
146 def poly_generator(deg):
147     """
148     A generator for a polynomial regression learner
149     :param deg: The degree of the polynomial to generate
150     :return: The generated learner
151     """
152     return make_pipeline(PolynomialFeatures(deg), LinearRegression())
153
154
155 def two_folds_get_all_models(model_generator, params, train_samples,
156                             train_labels, validation_samples,
157                             validation_labels):
158     """
159     Gets all the models of the two folds case
160     :param model_generator: The model generator
161     :param params: The parameters from which to generate the learners
162     :param train_samples: Train samples
163     :param train_labels: Train labels
164     :param validation_samples: Validation samples
165     :param validation_labels: Validation labels
166     :return: The models
167     """
168     models = []
169     for param in params:
170         model = model_generator(param)
171         model = two_folds_get_single_model(
172             model, train_samples, train_labels, validation_samples,
173             validation_labels)
174         models.append(model)
175     return models
176
177
178 def multiple_folds_get_single_model(model, folds, validation_excluded):
179     """
180     Gets a single model from the multiple folds case
181     :param model: The model to train
182     :param folds: The validation sets
183     :param validation_excluded: The train sets
184     :return: The trained model, its mean train error and its mean validation
185     error
186     """
187     folds_amount = len(folds)
188     train_losses = []
189     validation_losses = []
190     for i in range(folds_amount):
191         train_samples = validation_excluded[i][0]
192         train_labels = validation_excluded[i][1]
193         model.fit(train_samples, train_labels)
194
195         y_hat = model.predict(train_samples)

```

```

196         y = train_labels
197         train_losses.append(np.mean((y - y_hat) ** 2))
198
199         fold_samples = folds[i][0]
200         y_hat = model.predict(fold_samples)
201         y = folds[i][1]
202         validation_losses.append(np.mean((y - y_hat) ** 2))
203
204     train_error = np.mean(train_losses)
205     validation_error = np.mean(validation_losses)
206
207     return model, train_error, validation_error
208
209
210 def multiple_folds_get_all_models(model_generator, params, folds,
211                                 validation_excluded):
212     """
213     Gets all the models for the multiple folds case
214     :param model_generator: The model generator
215     :param params: The parameters from which to generate the models
216     :param folds: The validation sets
217     :param validation_excluded: The train sets
218     :return: The trained models as triplets, each one consists of the model,
219             its mean train error and its mean validation error
220     """
221     models = []
222     for param in params:
223         model = model_generator(param)
224         poly = multiple_folds_get_single_model(model, folds,
225                                                validation_excluded)
226         models.append(poly)
227     return models
228
229
230 def get_best_model(models):
231     """
232     Gets the best model (triplet)
233     :param models: The models to choose from
234     :return: The best model (triplet)
235     """
236     return min(models, key=lambda poly: poly[2])
237
238
239 def get_errors(models):
240     """
241     Gets a list of the mean train errors and a list of the mean validation
242     errors of the given models
243     :param models: The models to get their errors
244     :return: The list of train errors and the list of validation errors
245     """
246     train_errors = []
247     validations_errors = []
248     for model in models:
249         train_errors.append(model[1])
250         validations_errors.append(model[2])
251     return train_errors, validations_errors
252
253
254 def get_multiple_folds_models(model_generator, params,
255                               all_train_samples,
256                               all_train_labels):
257     """
258     Gets the models for the multiple folds case
259     :param model_generator: The model generator
260     :param params: The parameters from which to generate the models
261     :param all_train_samples: Train samples
262     :param all_train_labels: Train labels
263     :return: The models

```

```

264     """
265     S, V = get_k_folds(FOLDS_AMOUNT, all_train_samples, all_train_labels)
266     return multiple_folds_get_all_models(model_generator, params, S, V)
267
268
269 class Data:
270     """
271     Data object for question four
272     """
273     def __init__(self, sigma):
274         """
275         Initializes a data object for question four according to the noise
276         of the labels
277         :param sigma: The noise to give to the labels
278         """
279         self.all_train_samples, self.all_train_labels, self.test_samples, \
280         self.test_labels = generate_data(sigma)
281
282
283 def fourth_question_items_b_c(all_train_samples, all_train_labels,
284                               param_getter, param_description,
285                               regression_description):
286     """
287     Runs a general version of items b and c of question four
288     :param all_train_samples: Train samples
289     :param all_train_labels: Train labels
290     :param param_getter: Gets the parameter from the model
291     :param param_description: Description of the parameter
292     :param regression_description: Description of the regression
293     :return: None
294     """
295     S, V = get_k_folds(2, all_train_samples, all_train_labels)
296     S, V = S[0], V[0]
297     train_samples = S[0]
298     train_labels = S[1]
299     validation_samples = V[0]
300     validation_labels = V[1]
301     models = two_folds_get_all_models(poly_generator, DEGS, train_samples,
302                                       train_labels, validation_samples,
303                                       validation_labels)
304     best_model = get_best_model(models)
305     best_param = param_getter(best_model)
306     print()
307     print(' The best {} with two folds for {} is: '.
308           format(param_description, regression_description), best_param)
309
310
311 def fourth_question_items_d_e(model_generator, params, all_train_samples,
312                               all_train_labels, xlabel, title1, title2,
313                               title3, param_getter, param_description,
314                               regression_description):
315     """
316     Runs a general version of items d and e of question four
317     :param model_generator: The model generator
318     :param params: The parameters from which to generate the models
319     :param all_train_samples: Train samples
320     :param all_train_labels: Train labels
321     :param xlabel: The labels for the X axis
322     :param title1: One part of the title for the graph
323     :param title2: Second part of the title for the graph
324     :param title3: Third part of the title for the graph
325     :param param_getter: Gets the parameter from the model
326     :param param_description: Description of the parameter
327     :param regression_description: Description of the regression
328     :return: None
329     """
330     models = get_multiple_folds_models(model_generator, params,
331                                       all_train_samples,

```

```

332         all_train_labels)
333
334     best_model = get_best_model(models)
335     best_param = param_getter(best_model)
336     print()
337     print(' The best {} with cross validation for {} is: '.
338           format(param_description, regression_description), best_param)
339
340     train_errors, validation_errors = get_errors(models)
341
342     plt.plot(params, train_errors, linewidth=6, alpha=0.5, label='Train '
343              'Errors')
344     plt.plot(params, validation_errors, 'r-', linewidth=2,
345              label='Validation Errors')
346     plt.legend(loc='center left', bbox_to_anchor=(1, 0.5))
347     plt.xlabel(xlabel)
348     plt.ylabel('Error')
349     plt.title('The train and validation errors\nof {}\nas a function '
350              'of {}\n{}'.format(title1, title2, title3))
351     plt.show()
352
353
354 def fourth_question_items_f_g(model_generator, all_train_samples,
355                               all_train_labels,
356                               test_samples, test_labels, best_param):
357     """
358     Runs a general version of items f and g of question four
359     :param model_generator: The model generator
360     :param all_train_samples: Train samples
361     :param all_train_labels: Train labels
362     :param test_samples: Test samples
363     :param test_labels: Test labels
364     :param best_param: The parameter of the best model
365     :return: None
366     """
367     model = model_generator(best_param)
368     model.fit(all_train_samples, all_train_labels)
369
370     y_hat = model.predict(all_train_samples)
371     y = all_train_labels
372     train_error = np.mean((y - y_hat) ** 2)
373     y_hat = model.predict(test_samples)
374     y = test_labels
375
376     test_error = np.mean((y - y_hat) ** 2)
377     diff = np.fabs(train_error - test_error)
378
379     print(' The train error is: {}'.format(train_error))
380     print(' The test error is: {}'.format(test_error))
381     print(' The difference between the train error and the test error is: '
382           '{}'.format(diff))
383
384
385 def run_all_items(model_generator, params, xlabel, title1, title2, title3,
386                  best_param, all_train_samples, all_train_labels,
387                  test_samples, test_labels, param_getter, param_description,
388                  regression_description):
389     """
390     Runs a general version of all items from question four
391     :param model_generator: The model generator
392     :param params: The parameters from which to generate the models
393     :param xlabel: The label of the X axis
394     :param title1: One part of the title for the graph
395     :param title2: Second part of the title for the graph
396     :param title3: Third part of the title for the graph
397     :param best_param: The parameter of the best model
398     :param all_train_samples: Train samples
399     :param all_train_labels: Train labels

```

```

400     :param test_samples: Test samples
401     :param test_labels: Test labels
402     :return: None
403     """
404     fourth_question_items_b_c(all_train_samples, all_train_labels,
405                               param_getter, param_description,
406                               regression_description)
407     fourth_question_items_d_e(model_generator, params, all_train_samples,
408                               all_train_labels, xlabel, title1, title2,
409                               title3, param_getter, param_description,
410                               regression_description)
411     fourth_question_items_f_g(model_generator, all_train_samples,
412                               all_train_labels,
413                               test_samples, test_labels, best_param)
414
415
416 def get_deg_from_poly_model(poly):
417     """
418     Gets the degree of a given polynomial regression learner
419     :param poly: A triplet of a polynomial regression learner
420     :return: The degree of the polynomial regression learner
421     """
422     return poly[0].named_steps['polynomialfeatures'].degree
423
424
425 def run_question_four(sigma):
426     data = Data(sigma)
427
428     all_train_samples, all_train_labels, test_samples, test_labels = \
429         data.all_train_samples, data.all_train_labels, data.test_samples, \
430         data.test_labels
431
432     models = get_multiple_folds_models(poly_generator, DEGS,
433                                       all_train_samples,
434                                       all_train_labels)
435
436     best_deg = get_deg_from_poly_model(get_best_model(models))
437
438     run_all_items(poly_generator, DEGS, 'Degree',
439                 'polynomials in polynomial regression,',
440                 'as a function of the degree',
441                 'samples generated from a normal distribution\nwith mean '
442                 'with value of zero\nand standard deviation of {}'.format(
443                     sigma),
444                 best_deg, all_train_samples, all_train_labels,
445                 test_samples, test_labels, get_deg_from_poly_model,
446                 'degree', 'polynomial regression on the given data')
447
448
449 def get_diabetes_data():
450     X, y = datasets.load_diabetes(return_X_y=True)
451     train_samples = X[:50]
452     train_labels = y[:50]
453     test_samples = X[50:]
454     test_labels = y[50:]
455     return train_samples, train_labels, test_samples, test_labels
456
457
458 def ridge_generator(alpha):
459     """
460     A generator for a Ridge regression learner
461     :param alpha: The regularization parameter of the learner we want to
462     generate
463     :return: The generated learner
464     """
465     return Ridge(alpha)
466
467

```

```

468 def lasso_generator(alpha):
469     """
470     A generator for a Lasso regression Learner
471     :param alpha: The regularization parameter of the learner we want to
472     generate
473     :return: The generated learner
474     """
475     return Lasso(alpha, max_iter=10**4)
476
477
478 def get_best_reg_param(model):
479     return model[0].alpha
480
481
482 def run_question_five():
483     """
484     Run question five
485     :return: None
486     """
487     train_samples, train_labels, test_samples, test_labels = get_diabetes_data()
488
489     models = get_multiple_folds_models(ridge_generator, RIDGE_PARAMS,
490                                       train_samples,
491                                       train_labels)
492     ridge_best_param = get_best_reg_param(get_best_model(models))
493     ridge_best_model = ridge_generator(ridge_best_param)
494     ridge_best_model.fit(train_samples, train_labels)
495     ridge_y_hat = ridge_best_model.predict(test_samples)
496     test_error = np.mean((test_labels - ridge_y_hat) ** 2)
497     print()
498     print(' The test error for the best Ridge model is: {}'.format(test_error))
499
500     linear_regression_model = ridge_generator(0)
501     linear_regression_model.fit(train_samples, train_labels)
502     linear_regression_model_y_hat = linear_regression_model.predict(test_samples)
503     test_error = np.mean((test_labels - linear_regression_model_y_hat) ** 2)
504     print(' The test error for a linear regression model is: {}'.format(
505         test_error))
506
507     models = get_multiple_folds_models(lasso_generator, LASSO_PARAMS,
508                                       train_samples,
509                                       train_labels)
510     lasso_best_param = get_best_reg_param(get_best_model(models))
511     lasso_best_model = lasso_generator(lasso_best_param)
512     lasso_best_model.fit(train_samples, train_labels)
513     lasso_y_hat = lasso_best_model.predict(test_samples)
514     test_error = np.mean((test_labels - lasso_y_hat) ** 2)
515
516     print(' The test error for the best Lasso model is: {}'.format(test_error))
517     fourth_question_items_d_e(ridge_generator, RIDGE_PARAMS, train_samples,
518                               train_labels,
519                               'regularization parameter',
520                               'Ridge regression models',
521                               'the regularization parameter', '',
522                               get_best_reg_param, 'regularization parameter',
523                               'Ridge regression on the given data')
524
525     fourth_question_items_d_e(lasso_generator, LASSO_PARAMS, train_samples,
526                               train_labels,
527                               'regularization parameter',
528                               'Lasso regression models',
529                               'the regularization parameter', '',
530                               get_best_reg_param, 'regularization parameter',
531                               'Lasso regression on the given data')
532
533
534 # run_question_four(FIRST_ITERATION_SIGMA)
535 # run_question_four(SECOND_ITERATION_SIGMA)

```

```
536 # run_question_five()
```