

תרגיל 5 IML - מאי ביבי

שאלה 1:

א. נוכיח שעבור $h^* \in \text{ERM}_{\mathcal{H}_k}(S_{\text{all}})$ מתקיים בהסתברות של לפחות $1 - \delta$:

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

בתרגול הוכחנו שלכל היפותזה h ולכל δ מתקיים

$$\mathbb{P} \left(|L_{S_{\text{all}}}(h) - L_{\mathcal{D}}(h)| \leq \sqrt{\frac{\ln(2/\delta)}{2m}} \right) \geq 1 - \delta$$

בפרט עבור כל $h_i \in \mathcal{H}_k$ עבור $\frac{\delta}{|\mathcal{H}_k|}$ מתקיים

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

מחסם האיחוד, לכל $h_i \in \mathcal{H}_k$

$$\begin{aligned} \mathbb{P} \left(|L_{S_{\text{all}}}(h_i) - L_{\mathcal{D}}(h_i)| \geq \sqrt{\frac{2 \ln \left(\frac{2|\mathcal{H}_k|}{\delta} \right)}{2m}} \right) &\leq \mathbb{P} \left(\bigcup_{h_i \in \mathcal{H}_k} |L_{S_{\text{all}}}(h_i) - L_{\mathcal{D}}(h_i)| \geq \sqrt{\frac{\ln \left(\frac{2|\mathcal{H}_k|}{\delta} \right)}{2m}} \right) \\ &\leq \sum_{h_i \in |\mathcal{H}_k|} \mathbb{P} \left(|L_{S_{\text{all}}}(h_i) - L_{\mathcal{D}}(h_i)| \geq \sqrt{\frac{\ln \left(\frac{2|\mathcal{H}_k|}{\delta} \right)}{2m}} \right) \leq \sum_{h_i \in |\mathcal{H}_k|} \frac{\delta}{|\mathcal{H}_k|} = |\mathcal{H}_k| \frac{\delta}{|\mathcal{H}_k|} = \delta \end{aligned}$$

סה"כ קיבלנו שלכל $h_i \in \mathcal{H}_k$ מתקיים עם הסתברות של לפחות $1 - \delta$:

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

לכן עבור $h_i \in \mathcal{H}_k$ מתקיים עם הסתברות של לפחות $1 - \delta$:

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

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

ובפרט עבור $h_i = \min_{h \in \mathcal{H}_k} L(h)$ נקבל

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

כדרוש.

ב. מהסעיף הקודם

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

ולכל $i \in [k]$ מתקיים

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

לכן, עם הסתברות של לפחות $1 - \delta$:

$$\begin{aligned} L_{\mathcal{D}}(h^*) &\leq L_{\mathcal{D}}(h_j) + \sqrt{\frac{2 \ln(4k/\delta)}{\alpha m}} \leq \min_{h \in \mathcal{H}_i} L_{\mathcal{D}}(h) + \sqrt{\frac{2 \ln(4|\mathcal{H}_i|/\delta)}{(1-\alpha)m}} + \sqrt{\frac{2 \ln(4k/\delta)}{\alpha m}} \\ &= \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{2 \ln(4|\mathcal{H}_i|/\delta)}{(1-\alpha)m}} + \sqrt{\frac{2 \ln(4k/\delta)}{\alpha m}} \end{aligned}$$

שאלה 2:

כיוון ש- $X^T X = I_d$, הפיכה, ולכן הפתרון ל-LS הוא

$$\hat{w}^{LS} = (X^T X)^{-1} X^T y = (I_n)^{-1} X^T y = X^T y$$

א. בהרצאה ראינו שניתן להגיע למשוואה הבאה בדרך לפתרון ridge:

$$X^T y = (X^T X + \lambda I) \hat{w}^{ridge} = (I + \lambda I) \hat{w}^{ridge} \hat{w}^{ridge}$$

$$\implies \hat{w}^{ridge} = \frac{1}{1 + \lambda} X^T y = \frac{\hat{w}^{LS}}{1 + \lambda}$$

ב. יהי $w \in \mathbb{R}^d$. מתקיים

$$\|y - Xw\|^2 = \|X^T y - X^T X w\|^2 = \|\hat{w}^{LS} - w\|^2 = \sum_{i=1}^d (\hat{w}_i^{LS} - w_i)^2$$

בעיית ה-subset selection היא

$$\arg \min_{w \in \mathbb{R}^d} \|y - Xw\|^2 + \lambda \|w\|_0 = \arg \min_{w \in \mathbb{R}^d} \sum_{i=1}^d (\hat{w}_i^{LS} - w_i)^2 + \lambda \|w\|_0 = \arg \min_{w \in \mathbb{R}^d} \sum_{i=1}^d (\hat{w}_i^{LS} - w_i)^2 + \lambda \cdot \mathbf{1}_{w_i \neq 0}$$

עבור $i \in [d]$

$$\arg \min_{w_i \in \mathbb{R}} (\hat{w}_i^{LS} - w_i)^2 + \lambda \cdot \mathbf{1}_{w_i \neq 0}$$

נשים לב שעבור $(\hat{w}_i^{LS})^2 > \lambda$ אנחנו נרצה $w_i = \hat{w}_i^{LS}$ ואז

$$(\hat{w}_i^{LS} - w_i)^2 + \lambda \cdot \mathbf{1}_{w_i \neq 0} = \lambda$$

וזה באמת המינימום במקרה זה. עבור $(\hat{w}_i^{LS})^2 \leq \lambda$, נרצה $w_i = 0$, ואז נקבל

$$(\hat{w}_i^{LS} - w_i)^2 + \lambda \cdot \mathbf{1}_{w_i \neq 0} = \hat{w}_i^{LS}$$

וזה באמת המינימום במקרה זה. ונשים לב ששני המרים שתיארו הם בדיוק:

$$\eta_{\sqrt{\lambda}}(\hat{w}^{LS})_i = \begin{cases} \hat{w}_i^{LS} & |\hat{w}_i^{LS}| \geq \sqrt{\lambda} \\ 0 & \text{else} \end{cases}$$

שאלה 3:

א. בהרצאה ראינו:

$$\hat{w} = (X^T X)^{-1} X^T y$$

$$X^T y = (X^T X + \lambda I) \hat{w}(\lambda)$$

$$\Rightarrow \hat{w}(\lambda) = (X^T X + \lambda I)^{-1} X^T y$$

מתקיים:

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

כדורש

ב. A_λ לא מקרית ולכן

$$\mathbb{E}(\hat{w}(\lambda)) = \mathbb{E}(A_\lambda \hat{w}) = A_\lambda \mathbb{E}(\hat{w}) = A_\lambda w = (X^T X + \lambda I_d)^{-1} (X^T X) w$$

$$A_\lambda w = (X^T X + \lambda I_d)^{-1} (X^T X) w \neq w \text{ עבור } \lambda > 0 \text{ מתקיים}$$

ג. אנו יודעים שמתקיים:

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

מתקיים

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

ד. בהרצאה ראינו שמתקיים

$$\mathbb{E}(\|\hat{y} - y^*\|^2) = \mathbb{E}(\|\hat{y} - \bar{y}\|^2) + \|\bar{y} - y^*\|^2 = \text{Var}(\hat{y}) + \text{Bias}(\hat{y})^2$$

כאשר y^* הם ערכי האמת, \hat{y} הוא הפתרון שבוחר האלגוריתם ו- $\bar{y} = \mathbb{E}(\hat{y})$.
אצלנו מתקיים $y^* = w$, $\hat{y} = \hat{w}(\lambda)$, $\bar{y} = \mathbb{E}(\hat{w}(\lambda))$. נקבל

$$\text{Bias}(\lambda)^2 = \|\mathbb{E}(\hat{w}(\lambda)) - w\|^2 = \|A_\lambda w - w\|^2 = \|(A_\lambda - I) w\|^2$$

$$\text{Var}(\lambda) = \text{Tr}(\text{Var}(\hat{w}(\lambda))) = \text{Tr}(\sigma^2 A_\lambda (X^T X)^{-1} A_\lambda^T) = \sigma^2 \text{Tr}(A_\lambda (X^T X)^{-1} A_\lambda^T)$$

לכן

$$\text{MSE}(\lambda) = \text{Bias}(\lambda)^2 + \text{Var}(\lambda) = \|(A_\lambda - I) w\|^2 + \sigma^2 \text{Tr}(A_\lambda (X^T X)^{-1} A_\lambda^T)$$

נגזור ונציב $\lambda = 0$:

$$\left. \frac{d}{d\lambda} \text{Bias}(\lambda)^2 \right|_{\lambda=0} = \left. \frac{d}{d\lambda} (\|(A_\lambda - I) w\|^2) \right|_{\lambda=0} = \left. \frac{d}{d\lambda} \left(\sum_i \sum_j ((A_\lambda - I)_{i,j} w_j)^2 \right) \right|_{\lambda=0}$$

$$\begin{aligned}
&= 2 \sum_i \sum_j (A_\lambda - I)_{i,j} w_j \Big|_{\lambda=0} \cdot \frac{d}{d\lambda} \sum_j (A_\lambda - I)_{i,j} w_j \Big|_{\lambda=0} = 2 \sum_i \sum_j 0_{i,j} w_j \Big|_{\lambda=0} \cdot \frac{d}{d\lambda} \sum_j (A_\lambda - I)_{i,j} w_j \Big|_{\lambda=0} = 0 \\
&\frac{d}{d\lambda} \text{Var}(\lambda) \Big|_{\lambda=0} = \frac{d}{d\lambda} \sigma^2 \text{Tr} \left(A_\lambda (X^T X)^{-1} A_\lambda^T \right) \Big|_{\lambda=0} = \frac{d}{d\lambda} \sigma^2 \sum_i \left(A_\lambda (X^T X)^{-1} A_\lambda^T \right)_{i,i} \Big|_{\lambda=0} \\
&= \frac{d}{d\lambda} \sigma^2 \sum_i \left((X^T X + \lambda I_d)^{-1} (X^T X) (X^T X)^{-1} \left((X^T X + \lambda I_d)^{-1} (X^T X) \right)^T \right)_{i,i} \Big|_{\lambda=0} \\
&= \frac{d}{d\lambda} \sigma^2 \sum_i \left((X^T X + \lambda I_d)^{-1} (X^T X)^T (X^T X + \lambda I_d)^{-1T} \right)_{i,i} \Big|_{\lambda=0} \\
&= \frac{d}{d\lambda} \sigma^2 \sum_i \left((X^T X + \lambda I_d)^{-1} (X^T X) (X^T X + \lambda I_d)^{-1} \right)_{i,i} \Big|_{\lambda=0} \dots < 0 \\
&\text{לא כל כך ידעתי איך להמשיך..}
\end{aligned}$$

$$\frac{d}{d\lambda} \text{MSE}(\lambda) \Big|_{\lambda=0} = \frac{d}{d\lambda} \Big|_{\lambda=0} \left(\text{Bias}(\lambda)^2 + \frac{d}{d\lambda} \text{Var}(\lambda) \right) \Big|_{\lambda=0} < 0$$

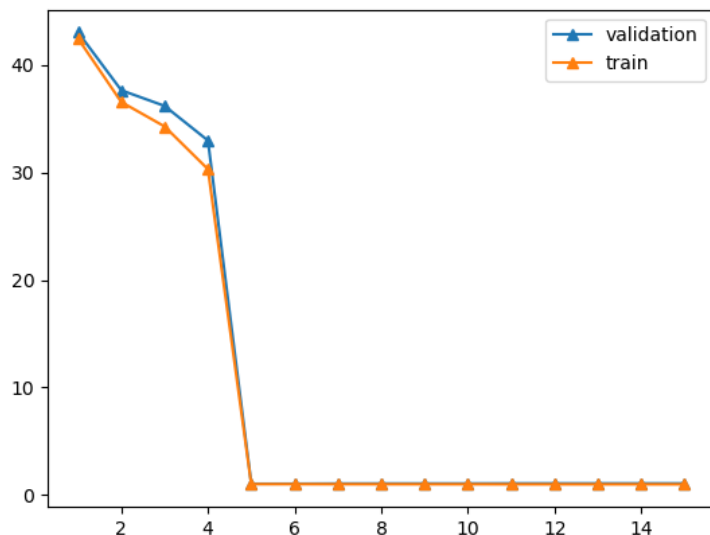
ה. הראינו שהנגזרת של $\text{MSE}(\lambda)$ בנקודה 0 היא שלילית, כלומר הפונקציה בירידה ולכן קיים $\lambda > 0$ עבורו $\text{MSE}(\lambda) < \text{MSE}(0)$.

שאלה 4:

רעש עם $\sigma = 1$

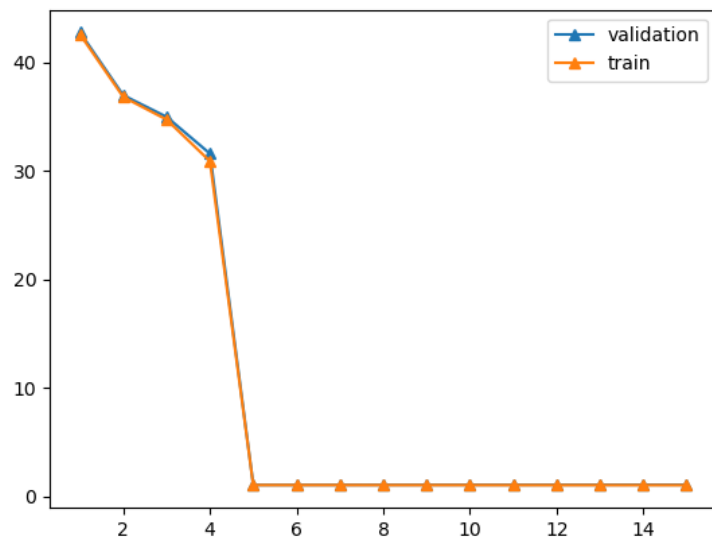
K-fold ולדיציה עם $k = 2$:

loss of KFold validation&training K=2 sigma=1



K-fold ולדיציה עם $k = 5$:

loss of KFold validation&training K=5 sigma=1

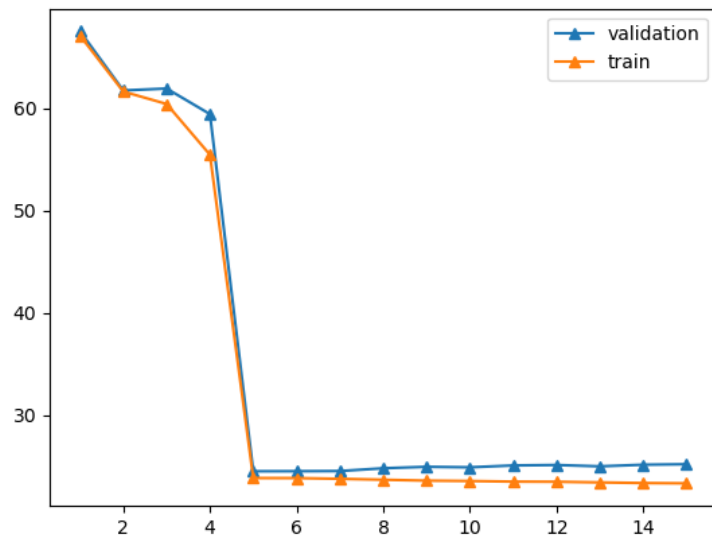


הדרגה עם השגיאה הכי נמוכה היא $d = 5$, וזה הגיוני כי הפולינום שלנו מדרגה 5.
 שגיאת הולדיציה עבור $d = 5$ היא 1.0435
 שגיאת test עבור $d = 5$ היא 22.1962

רעש עם $\sigma = 5$:

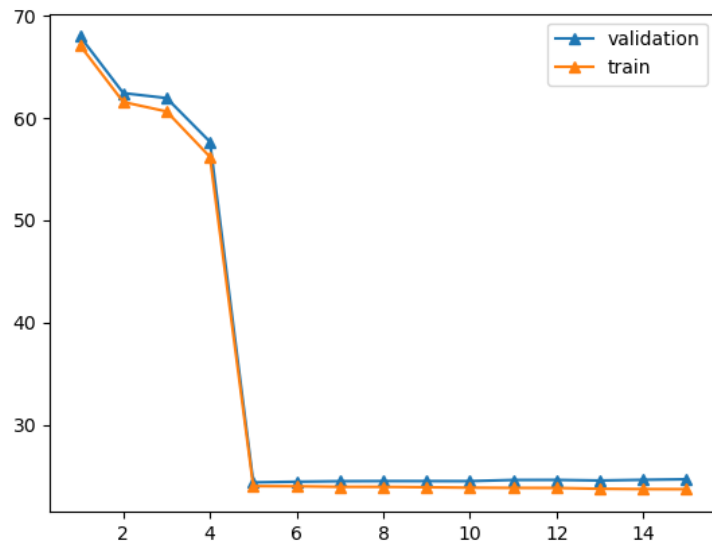
K-fold ולדיציה עם $k = 2$:

loss of KFold validation&training K=2 sigma=5



K-fold ולדיציה עם $k = 5$:

loss of KFold validation&training K=5 sigma=5



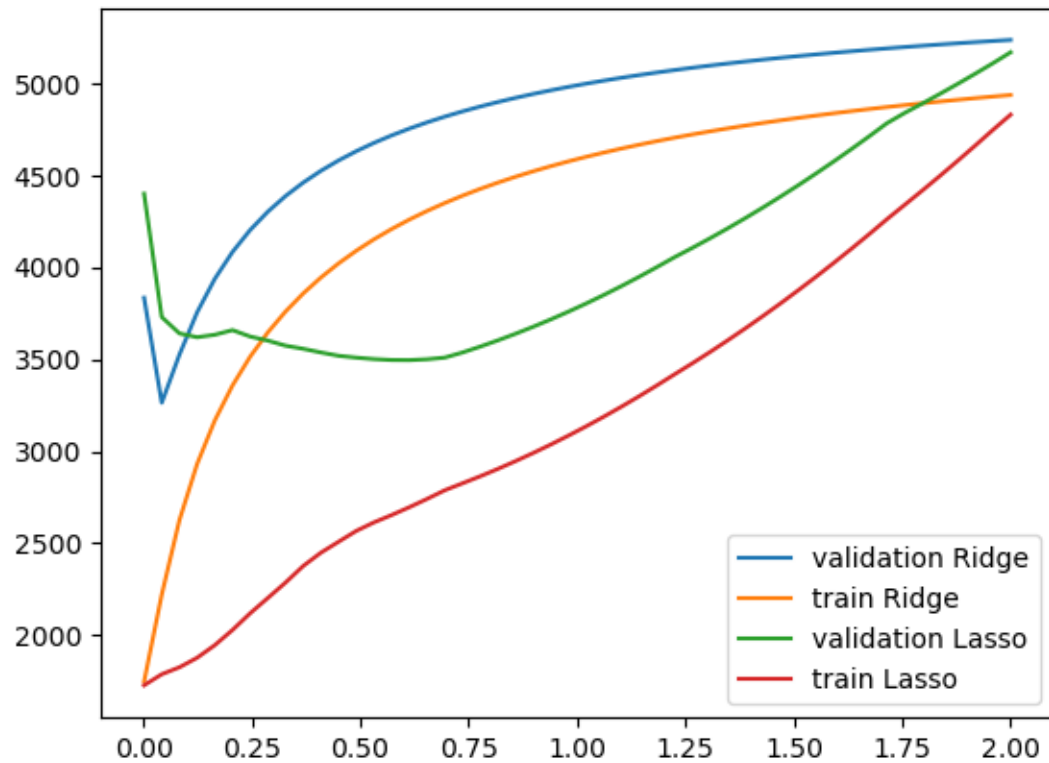
הדרגה עם השגיאה הכי נמוכה היא $d = 5$, וזה הגיוני כי הפולינום שלנו מדרגה 5.
 שגיאת הולדיציה עבור $d = 5$ היא 24.5348
 שגיאת ה-test עבור $d = 5$ היא 107.4876

כצפוי לדאטה עם יותר רעש, קיבלתי שגיאות גדולות יותר. אך עדיין קיבלנו $d = 5$

שאלה 5:

שיחקתי עם הערכים כדי לראות מתי השגיאה נמוכה יותר. בדקתי ערכי λ קטנים מאוד, כלומר מעט רגולריזציה, וגם ערכי $\lambda > 1$, הרבה רגולריזציה. לכן בדקתי ערכים בין 0 ל-2.

loss of KFold validation&training



```
best lambda for ridge: 0.04179591836734694
best lambda for lasso: 0.612938775510204
```

```
Ridge test error: 3191.397109707721
Lasso test error: 3652.376475971041
LinearRegression test error: 3612.249688324898
```

ניתן לראות כי ה- λ הכי טובה עבור ridge היא קטנה מאוד, כלומר אנחנו מאוד מתקרבים לרגרסיה לינארית רגילה. אך השגיאה של ridge קטנה יותר משל הרגרסיה הרגילה, כלומר ridge כן משפר את השגיאה.