

מערכות לומדות תרגיל 5

שיר שבח 322407701

31 במאי 2022

1 חלק תיאורטי

1.1 Regularization

1. בשאלה הבאה תראו כי למרות שאומדן Ridge מוטה הוא יכול להשיג MSE נמוך בהשוואה לאומדן LS. יהי $X \in \mathbb{R}^{m \times d}$ מטריצה קבועה, $y \in \mathbb{R}^d$ וקטור response, והניחו כי $X^T X$ הפיך. נסמן \hat{w} פתרון LS ו- \hat{w}_λ פתרון ridge עבור פרמטר רגולריזציה $\lambda \geq 0$ (כאשר $\hat{w}_0 = \hat{w}$).

הניחו כי המודל הלינארי נכון, כלומר $y = Xw + \varepsilon$ כאשר $\varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$. היזכרו כי במקרה הזה $\mathbb{E}(\hat{w}) = w$.

א. הראו כי $\hat{w}_\lambda = A_\lambda \hat{w}$ כאשר $A_\lambda = (X^T X + \lambda I_d)^{-1} (X^T X)$.

נסמן (לפי פירוק SVD) $X = U \Sigma V^T$. הוכחנו בתרגול כי פתרון ridge הוא $\hat{w}_\lambda = V \Sigma_\lambda U^T y$ ופתרון LS הוא: $\hat{w} = X^\dagger y = V \Sigma^\dagger U^T y$.

$$\begin{aligned} A_\lambda \hat{w} &= (X^T X + \lambda I_d)^{-1} (X^T X) V \Sigma^\dagger U^T y \\ &= (V \Sigma^T U^T U \Sigma V^T + \lambda I_d)^{-1} (V \Sigma^T U^T U \Sigma V^T) V \Sigma^\dagger U^T y \\ &= (V \Sigma^T \Sigma V^T + \lambda I_d)^{-1} V \Sigma^T \Sigma \Sigma^\dagger U^T y \\ &= V (\Sigma^T \Sigma + \lambda I_d)^{-1} V^T V \Sigma^T \Sigma \Sigma^\dagger U^T y \\ &= V \underbrace{(\Sigma^T \Sigma + \lambda I_d)^{-1} \Sigma^T \Sigma \Sigma^\dagger}_{\Sigma_\lambda} U^T y \\ \otimes &= V \Sigma_\lambda \Sigma \Sigma^\dagger U^T y \\ \otimes \otimes &= V \Sigma_\lambda U^T y \\ &= \hat{w}_\lambda \end{aligned}$$

$$\begin{aligned} \otimes \quad \Sigma_\lambda &= (\Sigma^T \Sigma + \lambda I_d)^{-1} \Sigma^T \\ \Sigma \Sigma^\dagger &= I \quad \otimes \otimes \end{aligned}$$

ב. מלמעלה, הסיקו כי לכל $\lambda > 0$ ridge הוא אומדן מוטה של w . כלומר, הראו כי לכל $\lambda > 0$, $\mathbb{E}(\hat{w}_\lambda) \neq w$.

$$\mathbb{E}(\hat{w}_\lambda) = \mathbb{E}(A_\lambda \hat{w}) \stackrel{\circledast}{=} A_\lambda \mathbb{E}(\hat{w}) \stackrel{\circledast\circledast}{=} A_\lambda w$$

\circledast לינאריות התוחלת.

$\circledast\circledast$ במקרה בו המודל הלינארי נכון, מתקיים $\mathbb{E}(\hat{w}) = w$.

נרצה להוכיח כי w הוא לא ווקטור עצמי של A_λ עם ע"ע 1.

נניח בשלילה שכן, כלומר, $A_\lambda w = w$. אזי:

$$\begin{aligned} (X^T X + \lambda I_d)^{-1} (X^T X) w &= w \\ (X^T X) w &= (X^T X + \lambda I_d) w \\ X^T X w &= X^T X w + \lambda w \\ 0 &= \lambda w \end{aligned}$$

אולם $\lambda > 0$ וגם $w \neq 0$. סתירה.

לכן קיבלנו:

$$A_\lambda w \neq w$$

ג. הראו כי: $\text{Var}(\hat{w}_\lambda) = \sigma^2 A_\lambda (X^T X)^{-1} A_\lambda^T$, עבור σ^2 השונות של רעש משוער. רמז: היזכרו כי עבור מטריצה קבועה B ווקטור מקרי z מתקיים כי $\text{Var}(Bz) = B \cdot \text{Var}(z) \cdot B^T$ וכי $\text{Var}(\hat{w}) = \sigma^2 (X^T X)^{-1}$.

$$\begin{aligned} \text{Var}(\hat{w}_\lambda) &= \text{Var}(A_\lambda \hat{w}) = A_\lambda \cdot \text{Var}(\hat{w}) \cdot A_\lambda^T \\ &= A_\lambda \cdot \sigma^2 (X^T X)^{-1} \cdot A_\lambda^T \\ &= \sigma^2 A_\lambda \cdot (X^T X)^{-1} \cdot A_\lambda^T \end{aligned}$$

ד. הפיקו ביטויים מפורשים עבור ה- $bias$ (מרובעת) וה- $variance$ של \hat{w}_λ כפונקציה של λ , כלומר כתבו פירוק $bias - variance$ עבור ה- MSE (mean square error) של \hat{w}_λ .

רמז: היזכרו כי עבור המקרה הרב-משתני ה- MSE מוגדר להיות:

$$MSE(\hat{y}) = \mathbb{E}(\|\hat{y} - y\|^2) = \mathbb{E}\left((\hat{y} - y)^T (\hat{y} - y)\right)$$

כאשר y הוא הערך האמיתי ו- \hat{y} זה האומדן.

ראינו שמתקיים כי עבור MSE יש פירוק ל- $bias - variance$:

$$\begin{aligned} \mathbb{E}(\|\hat{y} - y\|^2) &= \mathbb{E}(\|\hat{y} - \mathbb{E}(\hat{y})\|^2) + \|\mathbb{E}(\hat{y}) - y\|^2 \\ &= \text{Var}(\hat{y}) + bias^2(\hat{y}) \end{aligned}$$

ולכן:

$$\mathbb{E}(\|\hat{w}_\lambda - \hat{w}\|^2) = \underbrace{\mathbb{E}(\|\hat{w}_\lambda - \mathbb{E}(\hat{w}_\lambda)\|^2)}_{Var(\hat{w}_\lambda)} + \underbrace{\|\mathbb{E}(\hat{w}_\lambda) - \hat{w}\|^2}_{bias^2(\hat{w}_\lambda)}$$

שונות:

$$\begin{aligned} Var(\hat{w}_\lambda) &= \mathbb{E}(\|\hat{w}_\lambda - \mathbb{E}(\hat{w}_\lambda)\|^2) \\ &= \mathbb{E}(\|A_\lambda \hat{w} - \mathbb{E}(A_\lambda \hat{w})\|^2) \\ &= \mathbb{E}(\|A_\lambda \hat{w} - A_\lambda w\|^2) \\ &= \mathbb{E}(\|A_\lambda(\hat{w} - w)\|^2) \\ &= \mathbb{E}(((\hat{w} - w) A_\lambda)^T (\hat{w} - w) A_\lambda) \\ &= \mathbb{E}(A_\lambda^T (\hat{w} - w)^T (\hat{w} - w) A_\lambda) \\ &= A_\lambda^T \mathbb{E}((\hat{w} - w)^2) A_\lambda \\ &= A_\lambda^T Var(\hat{w}) A_\lambda \\ &= \sigma^2 A_\lambda (X^T X)^{-1} A_\lambda^T \end{aligned}$$

קיבלנו בדיוק לפי הסעיף הקודם.
הטיה בריבוע:

$$\begin{aligned} bias^2(\hat{w}_\lambda) &= \|\mathbb{E}(\hat{w}_\lambda) - \mathbb{E}(\hat{w})\|^2 \\ &= \|A_\lambda w - w\|^2 \\ &= \|(A_\lambda - I)w\|^2 \\ &= ((A_\lambda - I)w)^T (A_\lambda - I)w \\ &= w^T (A_\lambda - I)^T (A_\lambda - I)w \end{aligned}$$

ה. הראו ע"י גזירה כי:

$$\frac{\partial}{\partial \lambda} \mathbf{MSE}(\hat{w}_\lambda) |_{\lambda=0} = \frac{\partial}{\partial \lambda} bias^2(\hat{w}_\lambda) |_{\lambda=0} + \frac{\partial}{\partial \lambda} Var(\hat{w}_\lambda) |_{\lambda=0} < 0$$

כלומר, חשבו את הנגזרת של הפונקציה לעיל ביחס ל λ בנקודה $\lambda = 0$.

$$\begin{aligned}\frac{\partial}{\partial \lambda} \text{MSE}(\hat{w}_\lambda) &= \frac{\partial}{\partial \lambda} (\text{Var}(\hat{w}_\lambda) + \text{bias}^2(\hat{w}_\lambda)) \\ &\stackrel{\circledast}{=} \frac{\partial}{\partial \lambda} (\text{Var}(\hat{w}_\lambda)) + \frac{\partial}{\partial \lambda} (\text{bias}^2(\hat{w}_\lambda))\end{aligned}$$

\circledast לינאריות הנגזרת.
נגזור את ה bias^2 :

$$\begin{aligned}\frac{\partial}{\partial \lambda} (\text{bias}^2(\hat{w}_\lambda)) &= \frac{\partial}{\partial \lambda} (w^T (A_\lambda - I)^T (A_\lambda - I) w) \\ &= \frac{\partial}{\partial \lambda} (\| (A_\lambda - I) w \|^2) \\ &= \frac{\partial}{\partial \lambda} \left(\sum_{i=1}^n ((A_\lambda - I)_i w)^2 \right) \\ &= \frac{\partial}{\partial \lambda} \left(\sum_{i=1}^n \left(\sum_{j=1}^n (A_\lambda - I)_{ij} w_j \right)^2 \right) \\ &= 2 \sum_{i=1}^n \left(\sum_{j=1}^n (A_\lambda - I)_{ij} w_j \right) \cdot \left(\sum_{j=1}^n (A_\lambda - I)_{ij} w_j \right)'\end{aligned}$$

נשים לב כי מתקיים: $A_\lambda - I = 0$ עבור $\lambda = 0$. לכן כל הביטוי הנ"ל שווה ל-0. אז $\frac{\partial}{\partial \lambda} (\text{bias}^2(\hat{w}_\lambda)) = 0$.

נגזור את ה var :

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

לכן:

$$\begin{aligned}\sigma^2 \text{tr} \left((X^T X + \lambda I)^{-1} (X^T X) (X^T X)^{-1} A_\lambda^T \right) &= \sigma^2 \text{tr} \left((X^T X + \lambda I)^{-1} \left((X^T X + \lambda I)^{-1} (X^T X) \right)^T \right) \\ &= \sigma^2 \text{tr} \left((X^T X + \lambda I)^{-1} (X^T X) \left((X^T X + \lambda I)^{-1} \right)^T \right)\end{aligned}$$

$X^T X = U D U^T$ הוא: SVD הפירוק וגם: $(X^T X)^{-1} = (U D U^T)^{-1} = U D^{-1} U^T$

בנוסף, $X^T X + \lambda I$ מטריצה סימטרית לכן הפירוק SVD הוא: $X^T X = V R V^T$. וגם: $(X^T X + \lambda I)^{-1} = (V R V^T)^{-1} = V R^{-1} V^T$

$$\begin{aligned}
 &= \sigma^2 \text{tr} \left((V R^{-1} V^T) (U D U^T) (V R^{-1} V^T)^T \right) \\
 &= \sigma^2 \text{tr} \left((V R^{-1} V^T) (V R^{-1} V^T) (U D U^T) \right) \\
 &= \sigma^2 \text{tr} \left(V R^{-1} R^{-1} V^T (V R V^T - \lambda I) \right) \\
 &= \sigma^2 \text{tr} \left(V R^{-1} R^{-1} V^T V R V^T - \lambda V R^{-1} R^{-1} V^T \right) \\
 &= \sigma^2 \text{tr} \left(V R^{-1} V^T - \lambda V (R^{-1})^2 V^T \right) \\
 &= \sigma^2 \left(\text{tr} (V R^{-1} V^T) - \lambda \text{tr} \left(V (R^{-1})^2 V^T \right) \right) \\
 &= \sigma^2 \left(\text{tr} (R^{-1}) - \lambda \text{tr} \left((R^{-1})^2 \right) \right)
 \end{aligned}$$

משום שמתקיים:

$$R = \begin{bmatrix} \alpha_1 & & \\ & \ddots & \\ & & \alpha_d \end{bmatrix}, D = \begin{bmatrix} b_1 & & \\ & \ddots & \\ & & b_d \end{bmatrix}$$

אזי מתקיים עבור ווקטור עצמי u_i של $X^T X$:

$$\begin{aligned}
 X^T X u_i &= b_i u_i \\
 (X^T X + \lambda I) u_i &= b_i u_i + \lambda u_i = (b_i + \lambda) u_i
 \end{aligned}$$

לכן קיבלנו כי $\alpha_i = b_i + \lambda$ ואז:

$$R = D + \lambda I \rightarrow R^{-1} = \begin{bmatrix} \frac{1}{b_1 + \lambda} & & \\ & \ddots & \\ & & \frac{1}{b_d + \lambda} \end{bmatrix}$$

לכן $\text{tr}((R^{-1})) = \sum_{i=1}^d \left(\frac{1}{b_i + \lambda} \right)$

ולכן הביטוי שהגענו אליו קודם שווה ל:

$$= \sigma^2 \left(\text{tr} \left(\sum_{i=1}^d \left(\frac{1}{b_i + \lambda} \right) \right) - \lambda \text{tr} \left(\sum_{i=1}^d \left(\frac{1}{b_i + \lambda} \right)^2 \right) \right)$$

נגזור את הביטוי הזה שפיתחנו:

$$\begin{aligned}
\frac{\partial \text{var}(\hat{w}_\lambda)}{\partial \lambda} &= \frac{\partial}{\partial \lambda} \left(\sigma^2 \left(\text{tr} \left(\sum_{i=1}^d \left(\frac{1}{b_i + \lambda} \right) \right) - \lambda \text{tr} \left(\sum_{i=1}^d \left(\frac{1}{b_i + \lambda} \right)^2 \right) \right) \right) \\
&= -\sigma^2 \left(\sum_{i=1}^d \left(\frac{1}{b_i + \lambda} \right)^2 + \lambda \cdot \sum_{i=1}^d -2 \frac{1}{(b_i + \lambda)^3} \right) \\
&= -\sigma^2 \sum_{i=1}^d \left(\frac{1}{(b_i + \lambda)^2} - \frac{2\lambda}{(b_i + \lambda)^3} \right)
\end{aligned}$$

נציב $\lambda = 0$ ונבדוק את ערך הנגזרת:

$$-\sigma^2 \cdot \sum_{i=1}^d \frac{1}{b_i^2} < 0$$

לכן קיבלנו

$$\frac{\partial}{\partial \lambda} \text{MSE}(\hat{w}_\lambda) |_{\lambda=0} < 0$$

1. הסיקו מכך, כי אם המודל הלינארי נכון רגולריזציה קטנה של *ridge* עוזרת להפחית את MSE .

לפי הסעיף הקודם ראינו שבסביבת $\lambda = 0$, כאשר מגדילים את λ MSE יורד - משום שהנגזרת שלילית, וזה אומר כי עבור רגולריזציה קטנה - שמתקרבת ל- $\lambda = 0$, טעות ה- MSE פוחתת.

PCA 1.2

2. יהי $X : \Omega \rightarrow \mathbb{R}^d$ משתנה מקרי עם תוחלת אפס וקובריאנס $\Sigma \in \mathbb{R}^{d \times d}$. הראו כי לכל $v \in \mathbb{R}^d$, כאשר $\|v\|_2 = 1$, השונות של $\langle v, X \rangle$ זה לא גדול יותר מהשונות המתקבלת על ידי הטמעת PCA של X על תת מרחב חד מימדי (הניחו כי PCA משתמש ב- Σ בפועל).

נפתח את הביטוי $\text{Var}(v^T X)$ ונרצה להגיע לכך ש $\text{Var}(v^T X) \leq \text{Var}(u_1^T X)$.

$$\begin{aligned}
\text{Var}(v^T X) &= \mathbb{E} \left((v^T X - \mathbb{E}(v^T X))^2 \right) \\
&= \mathbb{E} \left(\left(v^T X - \mathbb{E} \left(\sum_{i=1}^d v_i X_i \right) \right)^2 \right) \\
&= \mathbb{E} \left(\left(v^T X - \sum_{i=1}^d v_i \mathbb{E}(X_i) \right)^2 \right) \\
&\stackrel{\circledast}{=} \mathbb{E} \left((v^T X)^2 \right)
\end{aligned}$$

$$\textcircled{*} \text{ נשים לב כי אם } \mathbb{E}(X) = 0 \text{ אזי } \mathbb{E} \begin{pmatrix} X_1 \\ \vdots \\ X_d \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} \text{ לכן גם לכל } i \in [d] : \mathbb{E}(X_i) = 0 \text{ ואז } \sum_{i=1}^d v_i \mathbb{E}(X_i) = 0.$$

נמשיך לפתח את הביטוי:

$$\begin{aligned} \text{Var}(v^T X) &= \mathbb{E} \left((v^T X)^2 \right) = \mathbb{E} \left(v^T X (v^T X)^T \right) \\ &= \mathbb{E} (v^T X X^T v) = v^T \mathbb{E} (X X^T) v \end{aligned}$$

נשים לב כי XX^T זוהי מכפלה חיצונית. ובה מקבלים:

$$XX^T = \begin{bmatrix} X_1 \\ \vdots \\ X_d \end{bmatrix} \begin{bmatrix} X_1 & \cdots & X_d \end{bmatrix} = \begin{bmatrix} X_1 X_1 & \cdots & X_1 X_d \\ \vdots & \ddots & \vdots \\ X_d X_1 & \cdots & X_d X_d \end{bmatrix}$$

ולכן

$$\mathbb{E}(XX^T) = \mathbb{E} \left(\begin{bmatrix} X_1 X_1 & \cdots & X_1 X_d \\ \vdots & \ddots & \vdots \\ X_d X_1 & \cdots & X_d X_d \end{bmatrix} \right) = \begin{bmatrix} \mathbb{E}(X_1 X_1) & \cdots & \mathbb{E}(X_1 X_d) \\ \vdots & \ddots & \vdots \\ \mathbb{E}(X_d X_1) & \cdots & \mathbb{E}(X_d X_d) \end{bmatrix}$$

ונשים לב כי מתקיים שלכל $i, j \in [d]$:

$$\mathbb{E}(X_i X_j) = \text{cov}(X_i, X_j)$$

משום ש:

$$\begin{aligned} \text{cov}(X_i, X_j) &= \mathbb{E}((X_i - \mathbb{E}(X_i))(X_j - \mathbb{E}(X_j))) \\ [\forall k \in [d] \quad \mathbb{E}(X_k) &= 0] = \mathbb{E}(X_i X_j) \end{aligned}$$

לכן קיבלנו כי $\mathbb{E}(XX^T) = \Sigma$.

$$\text{Var}(v^T X) = v^T \Sigma v$$

נכתוב את v כצירוף לינארי של ווקטורי הבסיס העצמיים $(u_1, \dots, u_d) : v = \sum_{i=1}^d \alpha_i u_i$ עבור $(\alpha_1, \dots, \alpha_d)$ סקלרים לא כולם אפס. נשים לב כי מתקיים:

$$\begin{aligned}
1 = \|v\|^2 &= \left\| \sum_{i=1}^d \alpha_i u_i \right\|^2 = \left\langle \sum_{i=1}^d \alpha_i u_i, \sum_{i=1}^d \alpha_i u_i \right\rangle \\
&= \sum_{i=1}^d \sum_{j=1}^d \alpha_i \alpha_j \langle u_i, u_j \rangle \\
[\text{orthogonal basis}] &= \sum_{i=1}^d \alpha_i^2 \|u_i\|^2 = \sum_{i=1}^d \alpha_i^2
\end{aligned}$$

נמשיך לפתח. נרשום $\vec{\alpha} = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_d \end{pmatrix}$ ונקבל:

$$\begin{aligned}
Var(v^T X) &= v^T \Sigma v = \left(\sum_{i=1}^d \alpha_i u_i \right)^T \Sigma \left(\sum_{i=1}^d \alpha_i u_i \right) \\
&= (\vec{\alpha}^T U)^T \Sigma \vec{\alpha} U \\
&= (U \vec{\alpha})^T (U D U^T) U \vec{\alpha} \\
&= \vec{\alpha}^T U^T U D \vec{\alpha} \\
&= \vec{\alpha}^T D \vec{\alpha}
\end{aligned}$$

מתקיים

$$\begin{aligned}
\alpha^T D \alpha &= \begin{bmatrix} \alpha_1 & \cdots & \alpha_d \end{bmatrix} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_d \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_d \end{bmatrix} = \begin{bmatrix} \alpha_1 & \cdots & \alpha_d \end{bmatrix} \begin{bmatrix} \alpha_1 \lambda_1 \\ \vdots \\ \alpha_d \lambda_d \end{bmatrix} \\
&= \alpha_1^2 \lambda_1 + \alpha_2^2 \lambda_2 + \dots + \alpha_d^2 \lambda_d \leq \lambda_1 \cdot \sum_{i=1}^d \alpha_i^2 = \lambda_1
\end{aligned}$$

לכן אם נציב $\alpha_1 = a, \alpha_2 = 0, \dots, \alpha_d = 0$ אזי

$$v = \sum_{i=1}^d \alpha_i u_i = u_1$$

לכן

$$var(u_1^T X) = \alpha_1^2 \lambda_1 = \lambda_1$$

וקיבלנו כי

$$\text{var}(v^T X) \leq \text{var}(u_1^T X)$$

מה שהיה להוכיח.

Kernels

3. יהי $k(x, x')$ קרנל PSD חוקי. ספקו קרנל PSD חוקי $\tilde{k}(x, x')$, הנבנה מ- k , כאשר הוא מנורמל. כלומר, לכל x מתקיים $\tilde{k}(x, x) = 1$. הוכיחו את תשובתכם.

משום שהקרנל PSD k חוקי, אזי גם לכל f פונקציה מתקיים כי גם

$$\tilde{k}(x, y) = f(x) k(x, y) f(y)$$

חוקית.

לכן אם נגדיר $f(x) = \frac{1}{\sqrt{k(x, x)}}$ אזי נקבל:

$$\tilde{k}(x, x') = \frac{1}{\sqrt{k(x, x)}} k(x, x') \frac{1}{\sqrt{k(x', x')}} = 1$$

וגם מתקיים כי:

$$\tilde{k}(x, x) = \frac{1}{\sqrt{k(x, x)}} k(x, x) \frac{1}{\sqrt{k(x, x)}} = \frac{k(x, x)}{k(x, x)} = 1$$

כנדרש.

4. נתון data set $S = \{(x_i, y_i)\}_{i=1}^m$ כאשר $x_i \in \mathbb{R}^d$ ו- $y_i \in \{\pm 1\}$, וממפה פיצ'רים $\psi : \mathbb{R}^d \rightarrow \mathcal{F}$ כאשר \mathcal{F} זה מרחב פיצ'רים כלשהו. תנו דוגמה של S data set ומפת פיצ'רים כך ש- S לא ניתן לחלוקה לינארית ב- \mathbb{R}^d (עבור $d \geq 2$) אבל הדאטא שעבר שינוי $S_\psi = \{(\psi(x_i), y_i)\}_{i=1}^m$ ניתן לחלוקה לינארית ב- \mathcal{F} .

$S = \{((-2, -2), 1), ((1, 1), -1), ((2, 2), -1)\}$
נשים לב כי הדאטא הזה אינו ניתן לחלוקה לינארית.

נגדיר $\psi\left(\begin{pmatrix} x \\ y \end{pmatrix}\right) = (x^2, y^2, x^2 + y^2)$ ואז מתקיים:

$$\psi\left(\begin{pmatrix} 2 \\ 2 \end{pmatrix}\right) = (4, 4, 8)$$

$$\psi\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}\right) = (1, 1, 2)$$

$$\psi\left(\begin{pmatrix} -2 \\ -2 \end{pmatrix}\right) = (4, 4, 8)$$

והנקודות האלו ניתנות לחלוקה לינארית במימד 3.

5. עבור כל אחת מהפונקציות הבאות, הוכיחו אם היא קרנל PSD חוקי או הביאו דוגמה נגדית.

א. $k(x, y) = \exp(\|x - y\|^2)$

$$x = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, y = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

$$k(e_1, e_2) = \exp\left(\left\|\begin{pmatrix} 1 \\ 0 \end{pmatrix} - \begin{pmatrix} 0 \\ 1 \end{pmatrix}\right\|^2\right) = e^2$$

$$G = \begin{bmatrix} 1 & e^2 \\ e^2 & 1 \end{bmatrix}$$

נחשב ערכים עצמיים:

$$\begin{aligned} \det\left(\begin{bmatrix} 1 & e^2 \\ e^2 & 1 \end{bmatrix} - \lambda I\right) &= \det\left(\begin{bmatrix} 1-\lambda & e^2 \\ e^2 & 1-\lambda \end{bmatrix}\right) \\ &= (1-\lambda)^2 - (e^2)^2 \end{aligned}$$

$$\begin{aligned} (1-\lambda)^2 - e^4 &= 0 \\ e^4 &= (1-\lambda)^2 \\ e^2 &= 1-\lambda \\ \lambda_1 &= 1-e^2 \end{aligned}$$

או

$$\begin{aligned} -e^2 &= 1-\lambda \\ \lambda_2 &= 1+e^2 \end{aligned}$$

קיבלנו כי ערך עצמי אחד הוא שלילי - $\lambda_1 = 1 - e^2$, לכן זאת לא מטריצת PSD .

ב. $k(x, y) = k_1(x, y) - k_2(x, y)$ עבור k_1 ו k_2 קרנלים חוקיים.

נגדיר $k_1(x, y) = 5$ ו $k_2(x, y) = 10$. הוכחנו בתרגול כי קרנל קבוע (שנותן לכל x, y אותו מספר) הוא קרנל חוקי. אולם $k(x, y) = k_1(x, y) - k_2(x, y) = -5$ ואז מתקיים שהמטריצת G של k 2×2 היא:

$$G = \begin{bmatrix} -5 & -5 \\ -5 & -5 \end{bmatrix}$$

ואז עבור $x = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ מתקיים:

$$xGx^T = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{bmatrix} -5 & -5 \\ -5 & -5 \end{bmatrix} \begin{pmatrix} 1 & 0 \end{pmatrix} = -5 < 0$$

לכן לפי התנאים של מטריצת PSD , G לא מקיימת את אחד התנאים. לכן k לא PSD חוקי.

ג. $k(x, y) = k_a(x_a, y_a) + k_b(x_b, y_b)$ עבור k_1 ו k_2 קרנלים חוקיים, $x = \begin{bmatrix} \frac{x_a}{x_b} \end{bmatrix}, y = \begin{bmatrix} \frac{y_a}{y_b} \end{bmatrix}$

מתקיים:

$$\begin{aligned} G &= \begin{bmatrix} k_a(x_a, x_a) + k_b(x_b, x_b) & k_a(x_a, y_a) + k_b(x_b, y_b) \\ k_a(y_a, x_a) + k_b(y_b, x_b) & k_a(y_a, y_a) + k_b(y_b, y_b) \end{bmatrix} \\ &= \begin{bmatrix} k_a(x_a, x_a) & k_a(x_a, y_a) \\ k_a(y_a, x_a) & k_a(y_a, y_a) \end{bmatrix} + \begin{bmatrix} k_b(x_b, x_b) & k_b(x_b, y_b) \\ k_b(y_b, x_b) & k_b(y_b, y_b) \end{bmatrix} \\ &= G_a + G_b \end{aligned}$$

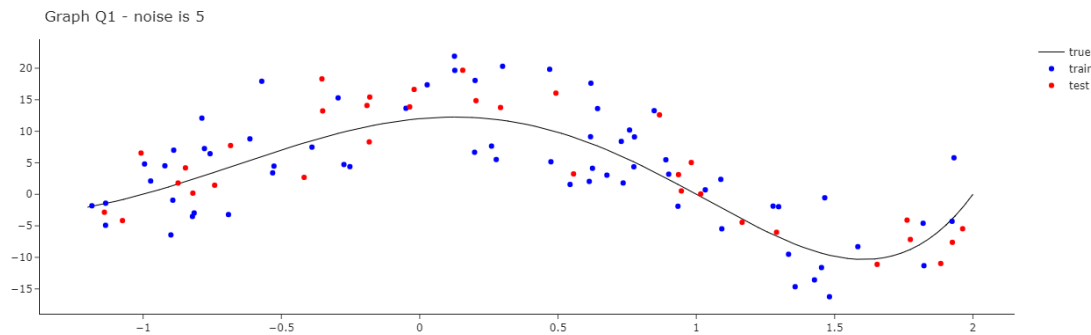
לכן:

$$x^T G x = x^T (G_a + G_b) x = x^T G_a x + x^T G_b x \geq 0$$

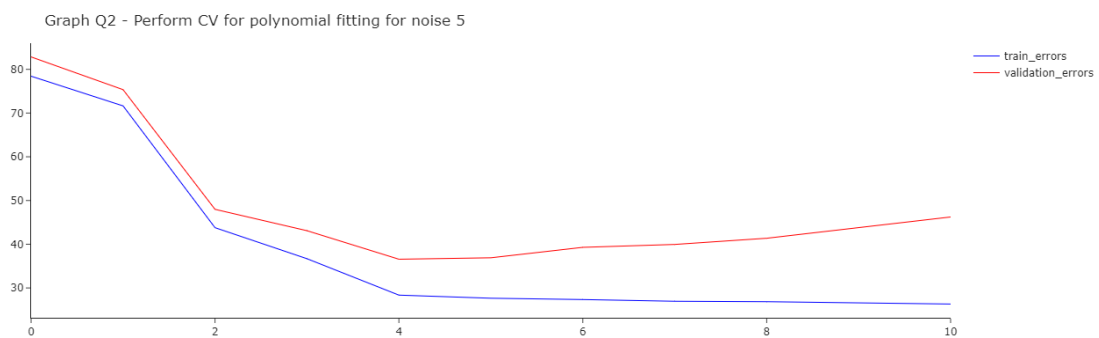
לכן G מטריצת PSD ולכן k חוקי.

חלק פרקטי

1.



2.

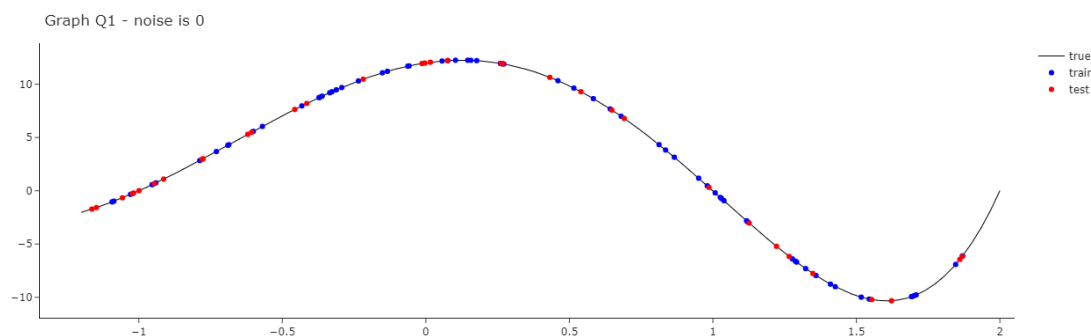


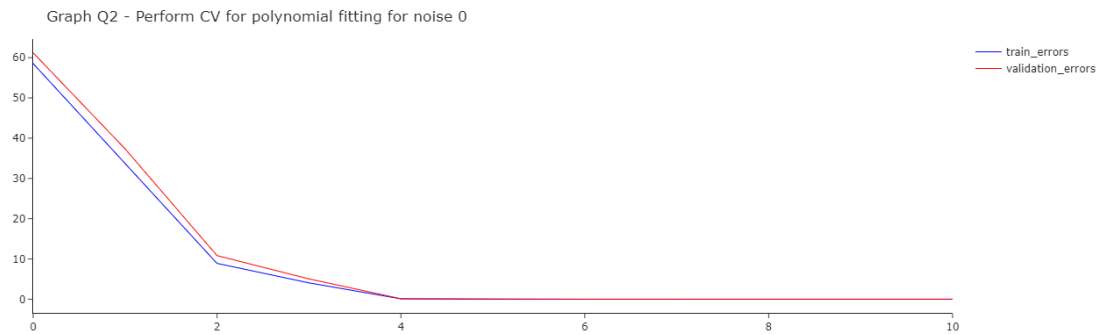
נשים לב כי הגרף של validation error לא מבצע overfit לכן באיזשהו שלב הטעות עולה. ול- train error יש overfitting לכן הטעות יורדת ויורדת (אולם השונות עולה).

3.

הא בו validation error הכי נמוך זה $k=4$. בו ה- test error הוא -16.967. ה- validation error הוא 36.572 ולכן שני ה- error -ים לא דומים.

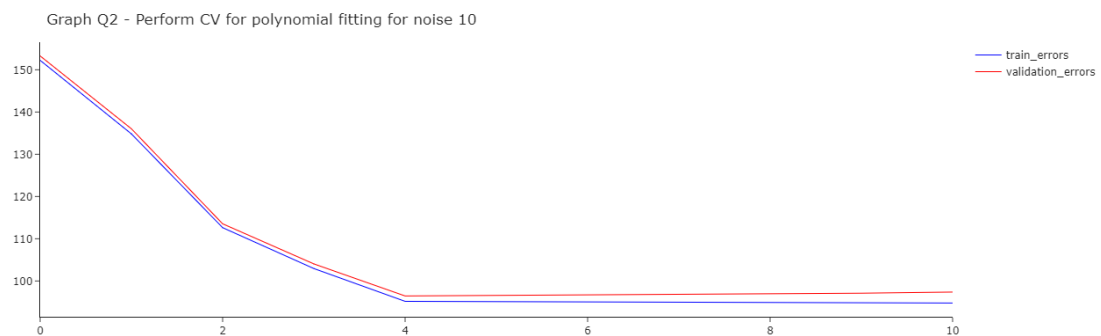
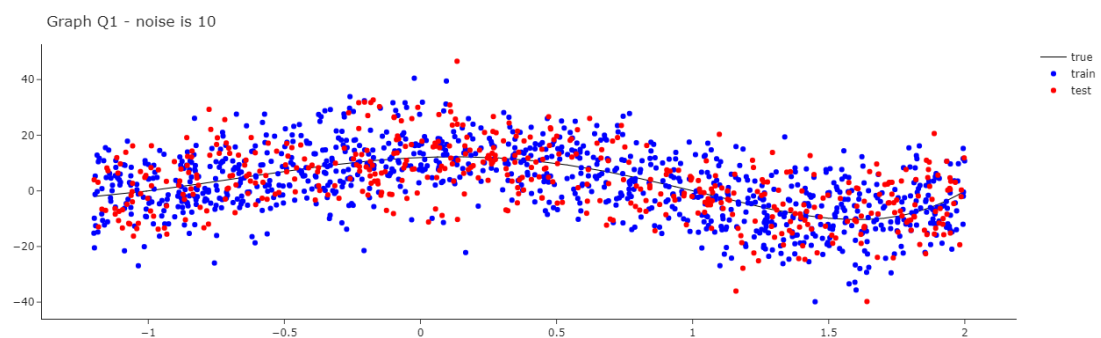
4.





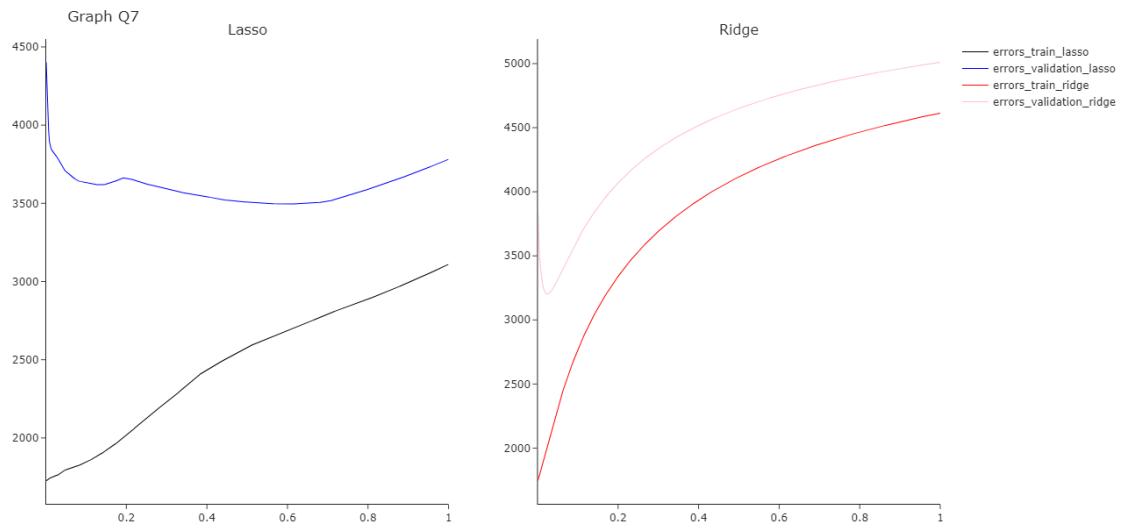
הא בו validation error הכי נמוך זה $k=5$. בו test error הוא (מעוגל) 0.0.
ה validation error הוא גם מעוגל ל 0 ולכן שני error-ים דומים.

.5



הא בו validation error הכי נמוך זה $k=4$. בו test error הוא 95.050.
ה validation error הוא 96.458 ולכן שני error-ים דומים. למרות שהעלנו את רמת הרעש, הוספנו בהרבה את מספר הדגימות. לכן הטעות בין שני error-ים דומה – בהתאם לחוק המספרים הגדולים (אולם עדיין הטעות גדולה מאוד, פשוט דומה ל validation).
.

.7



עבור טווח של לאמדות (0.001,1).

אנו שמים לב כי לasso, בגרף של הטעות של validation, יש לאמדות דומות, ולאחר מכן ככל שעולים בערך הלאמדה, הטעות עולה עד התייצבות.

8.

הלאמדה המינימלית בlasso: 0.597

הלאמדה המינימלית בridge: 0.025

הטעות בLinearRegression: 3612.249

הטעות בlasso: 3641.447

הטעות בridge: 3243.425

המודל שמקבל את הטעות הקטנה ביותר הוא ridge.