

תרגיל 5 | IML (67329)

שם: רונאל חרדים | ת"ז: 208917641

חלק I

Theoretical

שאלה 1

(א)

נראה כי $\hat{w}_\lambda = A_\lambda \hat{w}$

מתקיים השוויון הבא $A_\lambda \hat{w} = A_\lambda \hat{w}(\lambda = 0)$ ואנו יודעים כי מהפתרון ללא רגולציה מתקיים

$$\hat{w}(\lambda = 0) \iff \operatorname{argmin}_w (\|y - Ww\|_2^2) = (X^T X)^{-1} X^T y$$

ולכן:

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

ראינו בכיתה כי זה שווה ל- $\hat{w}(\lambda)$ כנדרש.

(ב)

נראה כי $\lambda > 0 \Rightarrow \mathbb{E}[\hat{w}(\lambda)] \neq w$

$$\mathbb{E}[\hat{w}(\lambda)] \stackrel{(3.A)}{=} \mathbb{E}[A_\lambda \hat{w}] = \mathbb{E}\left[(X^T X + \lambda I_d)^{-1} (X^T X) \cdot \hat{w}\right]$$

מכיוון שהמטריצה X קבועה אנו יודעים כי $\mathbb{E}[A_\lambda] = A_\lambda$ ולכן

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

כעת, אם $\lambda = 0$ אזי $A_\lambda = I$, אך מכיוון ש $\lambda > 0$ מתקיים:

$$(X^T X + \lambda I_d)^{-1} (X^T X) \neq I_d, \text{ therefore } (X^T X + \lambda I_d)^{-1} (X^T X) w \neq w$$

כנדרש.

(ג)

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

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

כנדרש.

(ד)

(ה)

באופן דומה למה שראינו בכיתה, נניח כי y^* ההיפוזזה האמיתית, \bar{y} הנצפה, ו \hat{y} האומדן. תחת ההגדרות הללו מתקיים:

$$E[\|\hat{y} - y^*\|^2] = \text{Var}[\hat{y}] + \text{bias}^2[\hat{y}].$$

מכיוון ש

$$\text{Var}[\hat{y}] = E[\|\hat{y} - \bar{y}\|^2] \wedge \text{bias}[\hat{y}] = \|\bar{y} - y^*\|^2.$$

ומתקיים כי $y^* = \hat{w}X$ ללא תנאי רגולריזציה. לכן $y^* = E[\hat{w}]$, בנוסף האסטימטור שלנו הוא $\hat{y} = \hat{w}(\lambda)$ ולכן גם $y = E[\hat{y}(\lambda)]$.
 נחשב:

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

נגדיר $(X^T X + \lambda I_d) = X' \text{ ולכן } A_\lambda = X'^{-1} (X^T X)$ וגם

$$\text{Var}(\lambda) = \sigma^2 \text{Tr}(X'^{-1} (X^T X) (X^T X)^{-1} (X'^{-1} (X^T X))^T) = \sigma^2 \text{Tr}(X'^{-1} (X^T X) X'^{-1})$$

כעת, נגדור את השונות לפי X' :

$$\left. \frac{d \text{Var}(\lambda)}{dX'} \right|_{\lambda=0} = \sigma^2 \text{Tr} \left(\frac{d}{dX'} (X'^{-1} (X^T X) X'^{-1}) \right)_{\lambda=0} = -2\sigma^2 (X^T X)^{-1} (X^T X)^{-1}$$

ולפי λ :

$$\left. \frac{d \text{Var}(\lambda)}{d\lambda} \right|_{\lambda=0} = \text{Tr} \left(\frac{d \text{Var}(\lambda)}{dX'} \frac{dX'}{d\lambda} \right) = -2\sigma^2 \text{Tr}((X^T X)^{-1} (X^T X)^{-1})$$

הערך הוא אי שלילי.

עבור ה $bias$:

$$\begin{aligned} bias(\lambda) &= \|\bar{y} - y^*\| = \|E[\hat{w}(\lambda)] - E[\hat{w}]\| = \|E[A_\lambda w] - E[\hat{w}]\| = \|A_\lambda E[w] - w\| \\ &= \|A_\lambda w - w\| = \|(A_\lambda - I) w\| \\ bias^2(\lambda) &= \|(A_\lambda - I) w\|^2 = w^T (A_\lambda - I)^T (A_\lambda - I) w \\ &\text{in terms of } X' \end{aligned}$$

עבור X' :

$$bias^2(\lambda) = w^T (X'^{-1} (X^T X) - I)^T (X'^{-1} (X^T X) - I) w$$

נגזור :

$$\begin{aligned} \left. \frac{d bias^2}{d\lambda} \right|_{\lambda=0} &= \left. \frac{d}{d\lambda} \left(\sum_i \left(\sum_j X'_{ij} w_j \right)^2 \right) \right|_{\lambda=0} \\ &= 2 \sum_i \left(\sum_j X'_{ij} w_j \right)_{\lambda=0} \cdot \left. \frac{d}{d\lambda} \left(\sum_j X_{ij} w_j \right) \right|_{\lambda=0} = 0 \end{aligned}$$

השגיאה נתונה ע"י:

$$MSE = bias^2 + var = w^T (A_\lambda - I)^T (A_\lambda - I) w + \sigma^2 \text{Tr} \left(A_\lambda (X^T X)^{-1} A_\lambda^T \right)$$

הנגזרת היא סכום הנגזרות וקטנה מ 0 . מסכיון שהנגזרת של ה $bias^2 = 0$ והנגזרת של השונות שלילית.

(ו)

אנו יודעים כי מודש לינארי בלי רגולריזציה מתקיים כאשר נשווה $\lambda = 0$. כעת ממה שראינו בסעיף קודם שעבור עברכים מסויימים מתקיים $MSE(\lambda) < 0 \Rightarrow \lambda > 0$ ניתן להסיק כי λ כזה מספק $MSE(\lambda) < MSE(0)$ כלומר הרגולריזציה מקטינה את הטעות.

שאלה 3

בהינתן קרנל $k(x, x')$ נגדיר את הקרנל המנורמל להיות עבור כל x - $\tilde{k}(x, x) = 1$ נגדיר:

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

ראישה נשים לב כי הקרנד המנורמל ווליד, והוא מייצג מ"פ $\psi/\|\psi\|$:

$$\begin{aligned}\tilde{k}(x, x') &= \frac{\langle \psi(x), \psi(x') \rangle}{\sqrt{\langle \psi(x), \psi(x) \rangle \langle \psi(x'), \psi(x') \rangle}} = \frac{\langle \psi(x), \psi(x') \rangle}{\sqrt{\|\psi(x)\|^2 \|\psi(x')\|^2}} = \frac{\langle \psi(x), \psi(x') \rangle}{\|\psi(x)\| \cdot \|\psi(x')\|} \\ &= \left| \frac{\psi(x)}{\|\psi(x)\|}, \frac{\psi(x')}{\|\psi(x')\|} \right|\end{aligned}$$

ומתקיים

$$\tilde{k}(x, x) = \frac{\langle \psi(x), \psi(x) \rangle}{\|\psi(x)\| \cdot \|\psi(x)\|} = \frac{\|\psi(x)\|^2}{\|\psi(x)\|^2} = 1$$

כנדרש.

שאלה 4

נביא דוגמה כנדרש:

עבור $d = 2, S = S_1 \cup S_2$ עבור $S_i = \{(x, y) : r_i \leq \sqrt{x^2 + y^2} \leq r_{i+1}\}$ המייצג סט של נקודות בתוך טבעת עם רדיוס פנימי r_i , ועם רדיוס חיצוני r_{i+1} כך ש $r_1 < r_2 < r_3$. נגדיר את פונקציית המיפוי באופן הבא: עבור נקודות x, y

$$(x, y) \in S, \psi(x, y) = (x, y, \sqrt{x^2 + y^2})$$

הפונקציה ψ ממפה את הנקודות לעצמן, עם אבולוציה $z = \sqrt{x^2 + y^2}$. מכיון ש $r_1 < r_2 < r_3$ הדאטה יהיה מופרד לינארית ב R^3 כנדרש.

חלק II

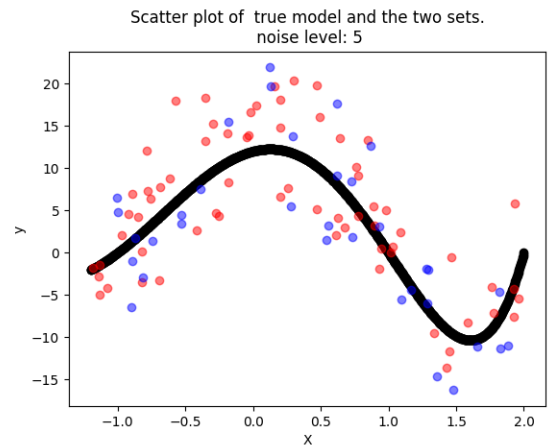
:Practical

שאלה 1

גרף הדגימות עם רעש 5:

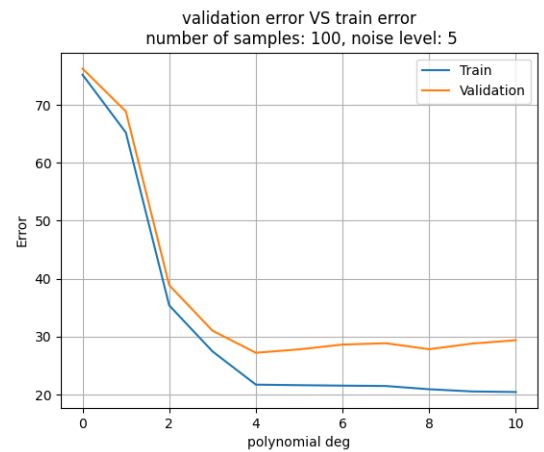
$train = red$

$test = blue$



שאלה 2

הגרף עבור $poly\ fit$ עם $k = 0...10$:



הסבר: אנו רואים כי בטעות על סט האימון יורדת ככל ש k עולה, כי האלגוריתם התאמן עליו כבר ומתאים את עצמו לסט טוב יותר ככל שדרגת הפולינום עולה. בעוד שעבור סט המבחן הוא חוזה טוב יותר עבור $k = 4$. כי ה var גבוה על סט האימון ככל ש k עולה.

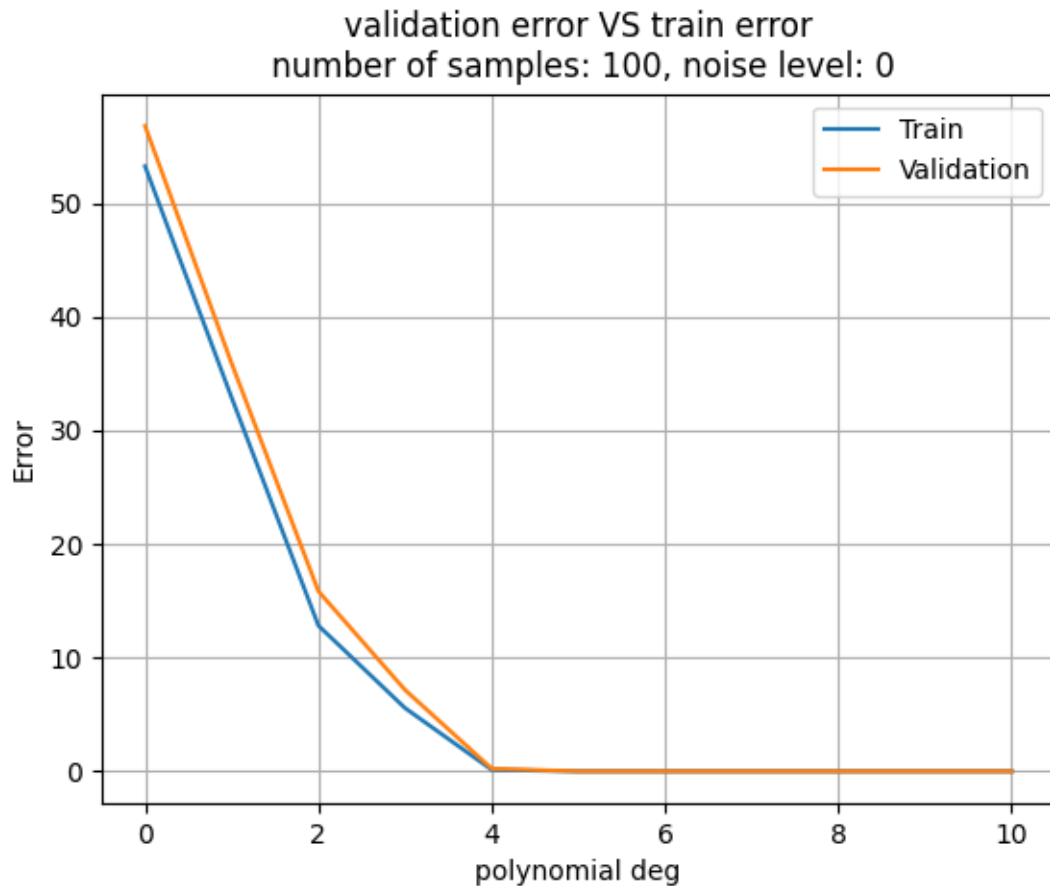
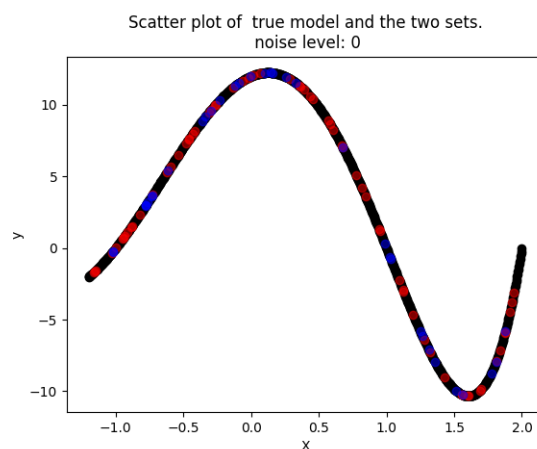
שאלה 3

הערך k הטוב ביותר היה $k = 4$.
הטעות שהתקבלה על סט הבחינה הייתה $test\ error = 27.16$.
 $validation\ error = 25.19$, לכן לא התרקיים שוויון בניהם, אך הם קרובים יחסית.

שאלה 4

$train = red$

$test = blue$

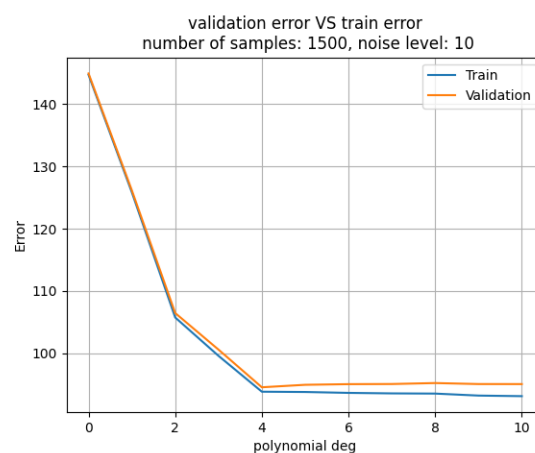
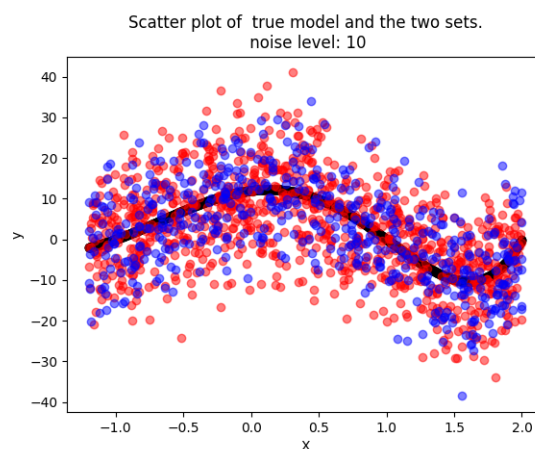


הסבר: עבור $k = 4$ מתקיים מפגש בין הטעויות של הטסט סט ו $validation$ משום שאין רעז אזי הטעות קרובה והאגוריתם חוזה טוב בשני המקרים. עבור $train set$ הטעות שווה ל 2.15 - שונה אך לא בהרבה.

שאלה 5

$train = red$

$test = blue$

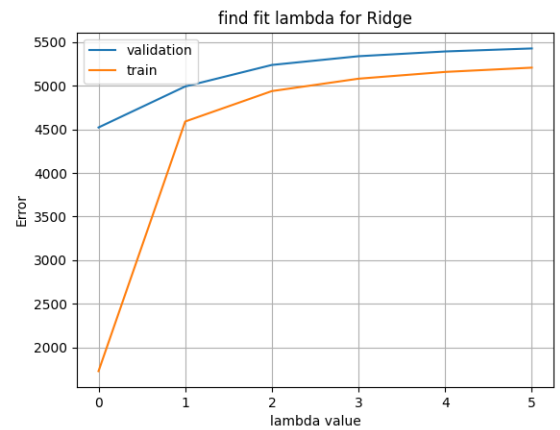


הסבר: נשים לב כי ככל שמספר הדגימות עולה טווח הטעות בין סט האימון ל $validation$ מצטצם והם מתנהגים כאילו אין רעש כלל. ככל שיש לאלגוריתם יותר דגימות הוא מאומן טוב יותר ולכן האומדן שלו קרוב יותר למציאות. זה נתמך גם ע"י חוק המספרים הגדולים - התכנסות לתוחלת.

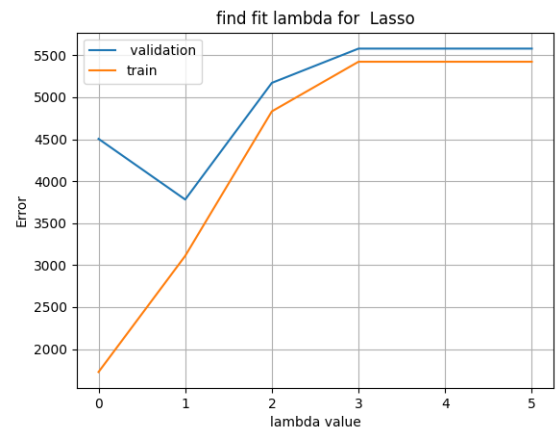
שאלה 7

הטווח הטוב ביותר עבור האלגוריתמים הוא: 6 – 0

$Ridge$ plot



:*Lasso plot*



נסביר את ההבדלים בין הגרפים ובין הסטים:

אנו רואים כי יש הבדלים בין האלגוריתמים, משום שהם ממשקלים את טעויות האלגוריתם עם נורמה שונה. אך שניהם מקבלים ערך נמוך כאשר למדה קטנה.

שאלה 8

ההפסד המינימלי מתקבל עבור פרמטרי רגולריזציה הבאים:

1 :*Lasso*

0 :*Ridge*

נחשב את הטעות לטסט סט:

4040.64 :*Lasso*

3612.24 :*Ridge*

3612.24 :*Linear Regression*