

מבוא למערכות לומדות- תרגיל בית 3- דו"ח עבודה

מגישות:



Section 1: Linear regression implementation

(Q1)

$$L_\delta(w, b, x_i, y_i) = \begin{cases} \frac{1}{2}(w^T x_i + b - y_i)^2, & |w^T x_i + b - y_i| \leq \delta \\ \delta \left(|w^T x_i + b - y_i| - \frac{1}{2}\delta \right), & \text{else} \end{cases}$$

הפונקציה היא פונקציית מקרים.

נתייחס לנגזרת שלה לפי b לפי המקרים:

- עבור המקרה $|w^T x_i + b - y_i| < \delta$, הטעות היא ריבועית ביחס ל b והפונקציה היא גזירה ביחס ל b . אין נקודות אי גזירות בתחום זה ולכן הגדרת הנגזרת והsubderivative מתלכדות.
- עבור המקרה $|w^T x_i + b - y_i| > \delta$, הטעות היא מוגדרת באמצעות הזהה והכפלה בסקלר של פונקציית הערך המוחלט- $|w^T x_i + b - y_i|$. פונקציה זו היא גזירה ביחס ל b בכל מקום מלבד באופן פוטנציאלי בנקודה שבה $w^T x_i + b - y_i = 0$. עם זאת, מאחר שאנו בתחום שבו $|w^T x_i + b - y_i| > \delta$, הפונקציה גזירה באזורים אלו גם כן, ולכן הגדרת הנגזרת והsubderivative מתלכדות.
- עבור המקרים בהם $|w^T x_i + b - y_i| = \delta$, כלומר כאשר
- $b_1 = \delta + y_i - w^T x_i$, $b_2 = -\delta + y_i - w^T x_i$ הפונקציה עדיין גזירה ביחס ל b בנקודות אלו כיוון שהנגזרות מימין מתלכדות עם הנגזרות משמאל.

סה"כ מאחר והפונקציה גזירה ביחס ל b על כל התחום נוכל לחשב את הנגזרות החלקיות בנפרד:

$$\frac{\partial L_\delta(w, b; x_i, y_i)}{\partial b} = \begin{cases} \frac{\partial \left(\frac{1}{2}(w^T x_i + b - y_i)^2 \right)}{\partial b}, & |w^T x_i + b - y_i| \leq \delta \\ \frac{\partial \left(\delta \left(|w^T x_i + b - y_i| - \frac{1}{2}\delta \right) \right)}{\partial b}, & \text{otherwise} \end{cases}$$

נחשב כל חלק בנפרד:

$$1. \quad |w^T x_i + b - y_i| \leq \delta$$

$$\begin{aligned} \frac{\partial \left(\frac{1}{2}(w^T x_i + b - y_i)^2 \right)}{\partial b} &= \frac{1}{2} \cdot 2 \cdot (w^T x_i + b - y_i) \cdot \frac{\partial (w^T x_i + b - y_i)}{\partial b} = (w^T x_i + b - y_i) \cdot 1 \\ &\Rightarrow \frac{\partial \left(\frac{1}{2}(w^T x_i + b - y_i)^2 \right)}{\partial b} = (w^T x_i + b - y_i) \end{aligned}$$

2. Otherwise

$$\begin{aligned} \frac{\partial \left(\delta \left(|w^T x_i + b - y_i| - \frac{1}{2}\delta \right) \right)}{\partial b} &= \frac{\partial \left(\delta |w^T x_i + b - y_i| - \frac{1}{2}\delta^2 \right)}{\partial b} \\ &= \begin{cases} \frac{\partial \left(-\delta(w^T x_i + b - y_i) - \frac{1}{2}\delta^2 \right)}{\partial b}, & (w^T x_i + b - y_i) < 0 \\ \frac{\partial \left(\delta(w^T x_i + b - y_i) - \frac{1}{2}\delta^2 \right)}{\partial b}, & (w^T x_i + b - y_i) \geq 0 \end{cases} \end{aligned}$$

$$= \begin{cases} -\delta, & (w^T x_i + b - y_i) < 0 \\ \delta, & (w^T x_i + b - y_i) \geq 0 \end{cases} = \delta \cdot \text{sign}(w^T x_i + b - y_i)$$

נקבל :

$$\frac{\partial \ell_\delta(w, b; x_i, y_i)}{\partial b} = \begin{cases} (w^T x_i + b - y_i) & , \quad |w^T x_i + b - y_i| \leq \delta \\ \delta \cdot \text{sign}(w^T x_i + b - y_i) & , \quad \text{otherwise} \end{cases}$$

(Q2)

נניח שיש לנו m זוגות של דוגמאות ותיוגים, כאשר כל דוגמה ממימד d כלומר :

$$\{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m) \mid x_i \in \mathbb{R}^d, y_i \in \mathbb{R} \text{ for } i \in [m]\}$$

נסמן את האלמנטים של דוגמה מסוימת, x_i באופן הבא: $x_i = (x^{(1)}, x^{(2)}, \dots, x^{(d)})^T$

$$X = \begin{pmatrix} x_1^{(1)} & \dots & x_1^{(d)} \\ \vdots & \ddots & \vdots \\ x_m^{(1)} & \dots & x_m^{(d)} \end{pmatrix} \quad \text{אזי נגדיר את המטריצה}$$

עבור וקטור $z = (z^{(1)}, z^{(2)}, \dots, z^{(m)})^T \in \mathbb{R}^m$ נגדיר פונקציית אינדיקטור $\mathbb{I}_{\leq \delta}: \mathbb{R}^m \rightarrow \mathbb{R}^m$ באופן הבא

$$\mathbb{I}_{\leq \delta}(z)_i = \begin{cases} 1, & |z^{(i)}| \leq \delta \\ 0, & |z^{(i)}| > \delta \end{cases}$$

עבור וקטור $z = (z^{(1)}, z^{(2)}, \dots, z^{(m)})^T \in \mathbb{R}^m$ נגדיר פונקציית אינדיקטור $\mathbb{I}_{> \delta}: \mathbb{R}^m \rightarrow \mathbb{R}^m$ באופן הבא

$$\mathbb{I}_{> \delta}(z)_i = \begin{cases} 0, & |z^{(i)}| \leq \delta \\ 1, & |z^{(i)}| > \delta \end{cases} \quad \text{כאשר } z^{(i)} \text{ הוא האלמנט במיקום ה-} i \text{ בוקטור } z.$$

$$1_m - \mathbb{I}_{\leq \delta}(z) = \mathbb{I}_{> \delta}(z)_i \quad \text{נשים לב כי}$$

בנוסף, עבור וקטור $z = (z^{(1)}, z^{(2)}, \dots, z^{(m)})^T \in \mathbb{R}^m$ נגדיר פונקציית אינדיקטור $\text{sign}: \mathbb{R}^m \rightarrow \mathbb{R}^m$ באופן הבא

$$\text{sign}(z)_i = \begin{cases} -1, & |z^{(i)}| < 0 \\ 1, & |z^{(i)}| > 0 \\ 0, & |z^{(i)}| = 0 \end{cases} \quad \text{כאשר } z^{(i)} \text{ הוא האלמנט במיקום ה-} i \text{ בוקטור } z.$$

אזי נוכל לבטא את $\nabla_w \mathcal{L}_H(w, b)$ באופן הבא:

$$\begin{aligned} \nabla_w \mathcal{L}_H(w, b) &= \nabla_w \left(\frac{1}{m} \sum_{i=1}^m \ell_\delta(w, b; x_i, y_i) \right) = \frac{1}{m} \left(\sum_{i=1}^m \nabla \ell_\delta(w, b; x_i, y_i) \right) \\ &= \frac{1}{m} \left(\sum_{i \text{ s.t. } |w^T x_i + b - y_i| \leq \delta} (w^T x_i + b - y_i) x_i + \sum_{i \text{ s.t. } |w^T x_i + b - y_i| > \delta} \delta \cdot \text{sign}(w^T x_i + b - y_i) x_i \right) \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{m} X^T ((Xw + b1_m - y) \odot \mathbb{I}_{\leq \delta} (Xw + b1_m - y)) + \\
&\frac{1}{m} \delta \cdot X^T (\text{sign}(Xw + b1_m - y) \odot \mathbb{I}_{> \delta} (Xw + b1_m - y)) = \\
&= \frac{1}{m} X^T ((Xw + b1_m - y) \odot \mathbb{I}_{\leq \delta} (Xw + b1_m - y)) + \\
&\frac{1}{m} \delta \cdot X^T (\text{sign}(Xw + b1_m - y) \odot (1_m - \mathbb{I}_{\leq \delta} (Xw + b1_m - y)))
\end{aligned}$$

* כאשר הסימן \odot הוא element-wise.

(נבדוק שהנוסחה חוקית מבחינת גדלים:

- $X_{m \times d} * w_{d \times 1} = v_{m \times 1}$ היא מכפלה חוקית ונותנת לנו וקטור עמודה בעל m מימדים, לכן גם הסכום חוקי $(X_{m \times d} * w_{d \times 1} + b1_m - y) \in R^m$ כיוון שכל הוקטורים שייכים ל R^m .
- הקלט ל $\mathbb{I}_{\leq \delta}$ הוא $z \in R^m$ ואכן $(X_{m \times d} * w_{d \times 1} + b1_m - y) \in R^m$ והפלט הוא גם כן וקטור $u \in R^m$
- מכפלת הדמדר (איבר - איבר) בין $(X_{m \times d} * w_{d \times 1} + b1_m - y) \in R^m$ ובין $\mathbb{I}_{\leq \delta} (Xw + b1_m - y) \in R^m$ ולידית כי שניהם מאותו מימד והפלט הוא גם כן R^m
- המטריצה $X^T \in R^{d \times m}$ ולכן חוקי להכפיל בין המטריצה לווקטור $(Xw + b1_m - y) \odot \mathbb{I}_{\leq \delta} (Xw + b1_m - y) \in R^m$ ונקבל פלט של וקטור $v \in R^d$ בגודל המצופה לוקטור גרדיאנט כיוון שאנו גוזרים לפי $w \in R^d$.
- סה"כ הביטוי $\frac{1}{m} X^T ((Xw + b1_m - y) \odot \mathbb{I}_{\leq \delta} (Xw + b1_m - y))$ מחזיר לנו $v \in R^d$ כפי שאנחנו רוצים.
- הביטוי $(1_m - \mathbb{I}_{\leq \delta} (Xw + b1_m - y))$ הוא ולידי כיוון שראינו כי $\mathbb{I}_{\leq \delta} (Xw + b1_m - y) \in R^m$ ולכן חישוב ההפרש בין 2 בוקטורים הוא לידי ומתקבל וקטור ששייך ל R^m .
- לפי איך שהגדרנו את sign לוקטורים ממימד m, אזי הביטוי $\text{sign}(Xw + b1_m - y)$ הוא ולידי ומחזיר וקטור ששייך ל R^m .
- $(\text{sign}(Xw + b1_m - y) \odot (1_m - \mathbb{I}_{\leq \delta} (Xw + b1_m - y)))$ ולידי כיוון שאנו מבצעים מכפלה איבר-איבר בין שני וקטורים ממימד m.
- כמו קודם, גם המכפלה $X^T (\text{sign}(Xw + b1_m - y) \odot (1_m - \mathbb{I}_{\leq \delta} (Xw + b1_m - y)))$ חוקית כיוון שכופלים $X^T \in R^{d \times m}$ בוקטור מימד m. נקבל פלט של וקטור $v \in R^d$ בגודל המצופה לוקטור גרדיאנט כיוון שאנו גוזרים לפי $w \in R^d$.
- סה"כ נקבל מהסכום של 2 הביטויים, 2 וקטורים ממימד d, לכן פעולת החיבור הזו היא אכן ולידית ואנחנו מקבלים וקטור גרדיאנט ממימד d כמצופה.)

כעת נפתח נוסחה עבור

$$\begin{aligned}
\frac{\partial L_H(w, b)}{\partial b} &= \frac{\partial}{\partial b} \left(\frac{1}{m} \sum_{i=1}^m \ell_\delta(w, b; x_i, y_i) \right) = \frac{1}{m} \left(\sum_{i=1}^m \frac{\partial}{\partial b} (w^T x_i + b - y_i) \right) \\
&= \frac{1}{m} \left(\sum_{i \text{ s.t. } |w^T x_i + b - y_i| \leq \delta} (w^T x_i + b - y_i) + \sum_{i \text{ s.t. } |w^T x_i + b - y_i| > \delta} \delta \cdot \text{sign}(w^T x_i + b - y_i) \right) = \\
&\frac{1}{m} (Xw + b1_m - y)^T \mathbb{I}_{\leq \delta} (Xw + b1_m - y) + \\
&\frac{1}{m} \delta \cdot (\text{sign}(Xw + b1_m - y))^T (1_m - \mathbb{I}_{\leq \delta} (Xw + b1_m - y))
\end{aligned}$$

(נבדוק שהנוסחה חוקית מבחינת גדלים:

- $X_{m \times d} * w_{d \times 1} = v_{m \times 1}$ היא מכפלה חוקית ונותנת לנו וקטור עמודה בעל m מימדים, לכן גם הסכום $(X_{m \times d} * w_{d \times 1} + b_{1 \times m} - y)$ כיוון שכל הוקטורים שייכים ל R^m .
- הקלט l הוא $z \in R^m$ ואכן $(X_{m \times d} * w_{d \times 1} + b_{1 \times m} - y) \in R^m$ והפלט הוא גם כן וקטור $u \in R^m$
- $(Xw + b_{1 \times m} - y)^T \mathbb{I}_{\leq \delta} (Xw + b_{1 \times m} - y)$ ולידי מאחר ומדובר במכפלה פנימית בין 2 וקטורים מימד m ובסופו של דבר הפלט מהמכפלה הזו הוא מספר כפי שהיינו מצפים.
- גם הביטוי $(Xw + b_{1 \times m} - y)^T \mathbb{I}_{\leq \delta} (Xw + b_{1 \times m} - y)$ הוא מספר, כמצופה.
- הביטוי $(1_m - \mathbb{I}_{\leq \delta} (Xw + b_{1 \times m} - y))$ הוא ולידי כיוון שראינו כי $(Xw + b_{1 \times m} - y) \in R^m$ ולכן חישוב ההפרש בין 2 בוקטורים הוא לידי ומתקבל וקטור ששייך ל R^m .
- לפי איך שהגדרנו את sign לוקטורים ממימד m , אזי הביטוי $\text{sign}(Xw + b_{1 \times m} - y)$ הוא ולידי ומחזיר וקטור ששייך ל R^m .
- $((\text{sign}(Xw + b_{1 \times m} - y))^T (1_m - \mathbb{I}_{\leq \delta} (Xw + b_{1 \times m} - y)))$ ולידי מאחר ומדובר במכפלה פנימית בין 2 וקטורים מימד m ובסופו של דבר הפלט מהמכפלה הזו הוא מספר כפי שהיינו מצפים.
- סה"כ הביטוי $(\text{sign}(Xw + b_{1 \times m} - y))^T (1_m - \mathbb{I}_{\leq \delta} (Xw + b_{1 \times m} - y))$ מחזיר מספר ממשי מאחר ומדובר במכפלה פנימית.
- לכן הסכום בין שני הביטויים של הנגזרות לפי b חוקי ומחזיר לנו מספר ממשי, כמצופה.

(Q3)

נשים לב כי הפרמטר דלתא בפונקציית הhuber loss הוא גודל המרחק (בערך המוחלט) בין החיזוי שלנו לתיג האמיתי כלומר residuals. הפרמטר דלתא קובע איך פונקציית loss תתנהג- כאשר residuals קטן מספיק, הפונקציה מתנהגת כמו שגיאת MSE רגילה. עבור ערכים גדולים יותר, פונקציית loss מתנהגת כמו שגיאת הערך המוחלט. מכאן שדלתא הוא סף כלשהו שמגדיר לנו האם בסבירות גבוהה מדידה מסויימת היא חריגה (outlier) ובהתאם נרצה לתת לה משקל נמוך יותר בחישוב loss.

לצורך הגדרת דלתא נרצה לבצע הערכה אילו נקודות עשויות להיחשב "outliers". בהינתן שהמידע היחיד שיש לנו על הבעיה הן הדגימות בדאטה סט, נרצה לקבל הערכה מהדאטהסט על מה גודל residuals האפשריים שיכולים להיות לנו. לצורך כך נרצה לאמן רגרסור לינארי (עם שגיאה ריבועית רגילה) ולהבין אילו ערכי residuals הם יחסית חריגים ביחס לדגימות האחרות.

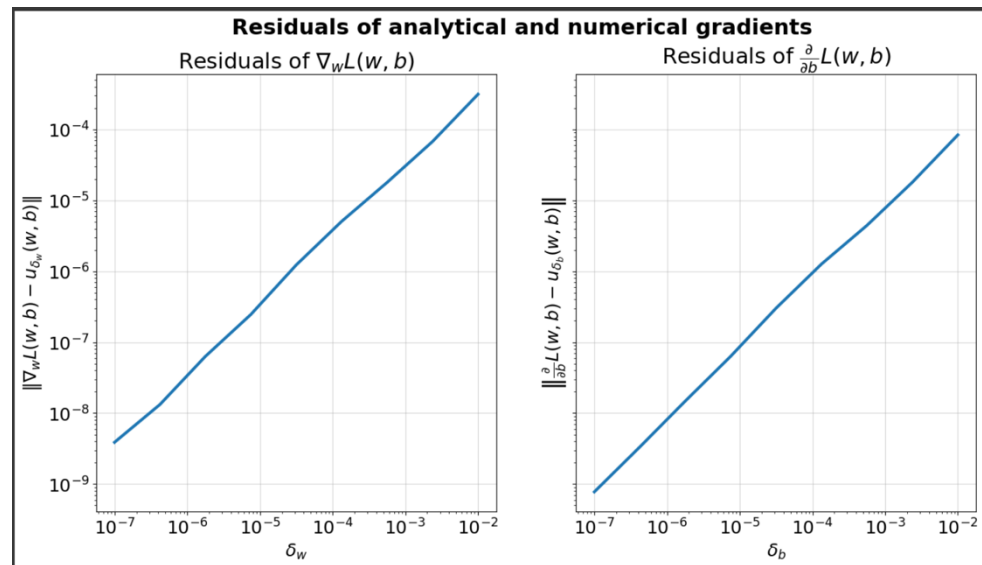
אנחנו בחרנו להשתמש בממד הטווח הבין רבעוני (IQR) שראינו בתרגול 1 כדי לקבוע מהו outlier. ע"פ המקובל ערכים קיצוניים הן תצפיות שנפלות מתחת ל- $Q1 - 1.5(IQR)$ או מעל $Q3 + 1.5(IQR)$. מכיוון שבמקרה זה אנו מסתכלים על הערך המוחלט של residuals, מספיק להסתכל רק על $Q3 + 1.5(IQR)$ כדי לקבוע את דלתא.

נציין שכמובן שהבחירה שלנו בממד IQR לצורך קביעת outliers אינה יחידה וישנן אפשרויות נוספות

פסודו קוד:

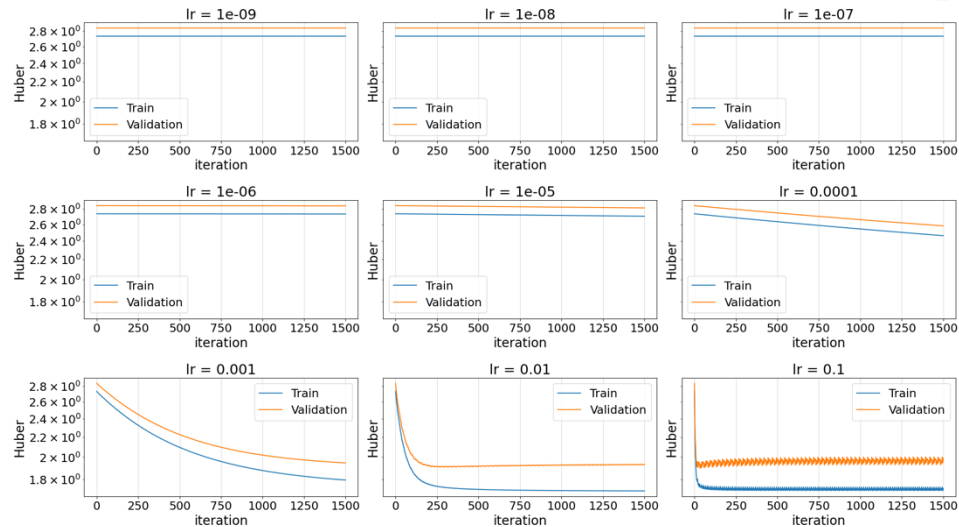
1. אמן את הדאטה עם רגרסור לינארי ופונקציית loss של MSE
2. חשב את residuals בערך מוחלט
3. א. חשב את האחוזון ה-25 של הדאטה ($Q1$)
ב. חשב את האחוזון של ה-75 של הדאטה ($Q3$)
ג. חשב את הטווח הבין רבעוני (IQR): $IQR = Q3 - Q1$
4. הגדר $\delta = Q3 + 1.5 * IQR$
5. החזר את δ

(Q4)



(Q5)

Training and Validation Losses Across Different Learning Rates for Huber Regression with `huber_delta=3.67`



השתמשנו ב `huber_delta=3.67` (חשוב רק על החלק של ה `train-validation` מתוך `train-validation`) ע"פ האלגוריתם שייצרנו בסעיף 3.

ניתן להצדיק שרואים בגרפים כך:

* עבור קצבי למידה נמוכים ($1e-09$ - $1e-5$), ניתן לראות שבקושי יש ירידה ב `loss` הן ב `train` והן ב `validation` מאחר וכל הנראה גודל הצעד ב `gradient` קטן מדי ולכן ההתכנסות איטית מאוד ולא קורית ב 1500 צעדים.

* כאשר קצב הלמידה עולה ל 0.0001 או רואים כי כן יש ירידה קטנה ב `loss`, הן ב `validation` והן ב `train`, אך עדיין הירידה הזו קטנה ביחס ל `loss` שמגיעים אליו בקצבי למידה גבוהים יותר. עבור קצב למידה הזה, אנו אכן מתקרבים לכיוון המינימום אך כמות הצעדים לא מספיקה עדין להתכנסות.

* עבור קצבי הלמידה 0.01, 0.001 או רואים כבר שחלה ירידה משמעותית ב-loss במהלך האימון והם מגיעים לערכי loss יחסית דומים בסוף האימון עבור שני הקצבים. עם זאת, הירידה עבור הקצב 0.001 היא מתונה יותר, בהשוואה ל-0.01, מאחר שגודל הצעד קטן יותר ולכן אנחנו מתקדמים לעבר המינימום בצורה איטית יותר.

* עבור קצב הלמידה 0.1, ההתכנסות היא מהירה מאד אבל יש הרבה יותר תנודות בערך loss בתהליך האימון וגם בקבוצת ולידיה. התנהגות זו מצביעה על כך שקצב למידה זה גדול מדי ולכן האלגוריתם מפספס מעט את המינימום בכל איטרציה.

* עוד נשים לב, שערך loss של קבוצת האימון ושל קבוצת הוולידציה דומה במהלך כל תהליך האימון של המסווג, ושווקטור w כלשהו שאנו מוצאים במהלך תהליך האימון שמקטין את loss על הtrain, מקטין את loss על הvalidation. התנהגות זו היא התנהגות טובה, שמראה לנו שהמסווג שלנו אכן לומד במהלך תהליך האימון ויודע לבצע הכללה על דאטה שלא ראה (validation). לא נראה שהגענו עדיין לנק' overfit במהלך תהליך האימון כיוון שלא ראינו אינדיקציה לכך שהtrain ממשיך לרדת אבל loss של validation כבר מתחיל לעלות. כמו כן, ברור שעבור קצב הלמידה 0.01 הייתה התכנסות מהירה יחסית למינימום כלשהו לאחר מס' נמוך של צעדים, ואימון נוסף בקושי מקטין את loss על הtrain והvalidation.

לדעתנו, קצב הלמידה האופטימלי הוא $lr=0.01$ כיוון שהוא מתכנס למינימום, ללא תנודות כמעט ובצורה מהירה יותר מאשר הקצב 0.001. בנוסף, הוא מקטין את loss גם על הtrain וגם על הvalidation והloss שקיבלנו על קבוצת הוולידציה היה הנמוך ביותר עבור קבוצה זו.

הגדלה של מספר הצעדים עבור lr זה אינה הכרחית כיון שניתן לראות שהאלגוריתם הגיע כבר לנק' מינימום כבר לאחר מס' מועט של צעדים, ואימון נוסף כבר בקושי מקטין את loss.

(Q6)

robustness של פונקציית loss huber מנוצלת בצורה טובה היותר כאשר הדאטה מכיל outliers, כאשר הדאטה מגיע מהתפלגות עם זנב ארוך (יש הסתברות נמוכה לקבל ערכים מאוד רחוקים מהערכים שבדר"כ מקבלים) או כאשר ההנחה שלנו שהרעש מתפלג נורמלית נשברת. במקרים כאלו, רגרסיה מסוג OLS (בעיית הריבועים הפחותים הרגילה) עובדת פחות טוב כיוון שהיא רגישה מאוד למדידות חריגות. מהגדרת בעיית הריבועים הפחותים באמצעות שגיאה ריבועית, ערכי מדידות חריגים שיש להם שגיאה גדולה יותר מהשגיאה הממוצעת, משפיעים הרבה יותר על גודל loss מאשר המדידות האחרות ובכך גורמים לקו הרגרסיה להתאים את עצמו הרבה יותר לרעש. המאפיינים של פונקציית loss huber מסייעים לה להתמודד עם בעיית הרעש- פונקציית loss זו מתנהגת כמו שגיאה ריבועית כאשר residuals קטנים ומתנהגת כמו שגיאת הערך המוחלט כאשר residuals גדולים מסף מסויים.

פונקציית loss huber מאפשרת לנצל את היתרונות של שתי השיטות- השגיאה האבסולוטית אינה גזירה ב-0 וקשה לאופטימיזציה, בעוד שפונק' loss huber מתנהגת סביב 0 כמו שגיאה ריבועית ולכן גזירה. עם זאת, עבור ערכי שגיאה גדולים, פונק' loss huber פחות רגישה לרעש כיוון שנותנת פחות חשיבות לרעש ואנומליות (שם השגיאה לינארית).

Section 2: Evaluation and Baseline

(Q7) עבור הדאמי יש לחשב שגיאה ריבועית על האימון

Model	Section	Train Huber Loss (cross validated) Mean (std)	Train MSE (cross validated) Mean (std)	Valid MSE (cross validated) Mean (std)
Dummy	2	2.29(+/- 0.1)	5.27 (+/- 0.26)	5.29 (+/- 1.051)

(Q8)

הערכים שבחרנו לבדוק עבור learning rate הם 31 ערכים ברווחי log בטווח $10^{-3} - 10^0$ (כך יצא שההפרש בין חזקה לחזקה הוא 0.1)

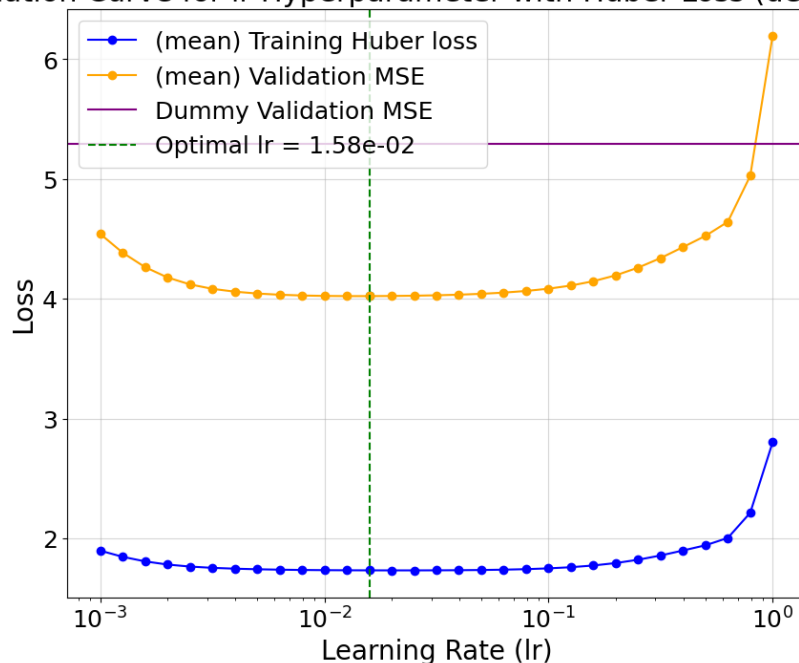
בחרנו בטווח זה כיוון שראינו בשאלה 5 שהאלגוריתם שלנו מתכנס עבור learning rates בין 10^{-3} , 10^{-2} ושבאזור 10^{-1} יש קצת רעש אך עדיין הייתה התכנסות. לכן החלטנו להרחיב את טווח החיפוש על טווח learning rates הזה כדי לאפסם את learning rate כך שנקבל שגיאת ולידציה הכי קטנה.

ע"פ ההוראות, אימנו עבור כל learning rate מודל רגרסיה המאומן על huber_loss (השתמשנו ב cross validation עם חלוקה 5 קבוצות כך שלמעשה לכל learning rate אומנו 5 מסוגים).

עבור כל learning rate חישבנו את ממוצע שגיאת האימון וממוצע שגיאת המבחן (ממוצע על חמשת הסגמנטים מה cross validation).

להלן הגרף שקיבלנו:

Validation Curve for lr Hyperparameter with Huber Loss (delta=3.65)



כפי שניתן לראות, הן שגיאת האימון והן שגיאת המבחן מתנהגות באופן הבא: ככל שהlr עולה השגיאה הולכת ויורדת ולאחר מכן עבור lr גדולים מדי השגיאה מתחילה לעלות בחזרה. ההתנהגות הזו צפויה, כי כפי שראינו קודם, עבור lr גדולים מדי, אלגוריתם הsgd כבר מתחיל להתבדר ולכן לא מקטין את השגיאה. עבור lr קטנים, ההתכנסות של האלגוריתם איטית יותר ולכן גודל loss שמגיעים אליו במספר קבוע של צעדים הוא עדין גדול לעומת ה שאפשר להגיע עם lr גדול יותר.

נשים לב כי אנו מודדים את השגיאה באימון והשגיאה במבחן באמצעות מטריקות שונות, וההפרש בין שגיאת האימון לשגיאת המבחן נובע בין השאר גם מההבדל בחישוב (עבור huber loss, לא מעלים בריבוע שגיאות גדולות מסף דלתא).

ה learning rate האופטימלי שמצאנו באימון הוא:

```
Optimal learning rate: 0.02
Training Huber loss at optimal lr: 1.73(+/- 0.10)
Validation MSE at optimal lr: 4.02(+/- 1.00)
```

Model	Section	Train Huber Loss (cross validated)	Train MSE (cross validated)	Valid MSE (cross validated)
Dummy	2	2.29(+/- 0.1)	5.27 (+/- 0.26)	5.29 (+/- 1.051)
Linear	2	1.73(+/- 0.1)		4.02 (+/- 1)

(Q9)

עבור המודל dummy שחזרה תמיד ערך אחד והוא הע הממוצע, לנרמול אין שום אפקט. זאת מכיוון שמודל זה בכלל לא מתייחס מה הם הפיצירים של דוגמה ספציפית בעת חיזוי, אלא הוא תמיד חוזר את אותו ערך הע הממוצע (ערך יחיד).

עבור המודל הלינארי עם huber-loss –

ברמה התאורטית-

שיטות הנרמול אליהן נחשפנו בקורס - min-max ו standardization הן טרנספורמציות לינאריות:

$$\text{min max scaling : } x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}} = cx + b \text{ where } c = \frac{1}{x_{\max} - x_{\min}}, b = \frac{-x_{\min}}{x_{\max} - x_{\min}}$$

$$\text{standardization: } x' = \frac{x - \mu}{\sigma} = mx + n \text{ where } m = \frac{1}{\sigma}, n = \frac{-\mu}{\sigma}$$

מאחר ושיטות הנרמול שלנו הן טרנספורמציות לינאריות, ניתן להוכיח כי לכל מסווג w המאופיין ע"י המשוואה

$$y = w^T \bar{x} + b \text{ ניתן למצוא מסווג } w' \text{ שיפעל על הפיצירים המנורמלים } x' \text{ ויחזיר את אותה פרדיקציה כלומר } y' = w'^T \bar{x}' + b'.$$

(ההוכחה מצורפת למטה).

מההוכחה אנו מקבלים שקיים קשר חד-חד ערכי ועל בין הפרמטרים w, b עבור x לפני נרמול לפרמטרים w', b' עבור x לאחר נרמול. הסט של ההיפר-מישורים האפשריים לפני נרמול ואחרי נרמול הוא למעשה אותו סט. מכאן, שהערך של ה training loss (הביצועים), ובמיוחד הערך האופטימלי של ה training loss יישאר אותו הדבר ברמה התאורטית במהלך תהליך אימון (כלומר תמיד נוכל להגיע לאותו ערך של מינימום, לא לאותו arg min).

בתהליך האימון שלנו אנו מחפשים מינימום עבור פונ' מטרה קמורה שהיא huber-loss. loss הוא בסופו של דבר פונקציה של w, b . כפי שראינו, הנרמול לא אמור להשפיע על ערך loss שאנו יכולים להשיג כיוון שכל ערך loss שניתן להשיג עם w, b שמוצאים לפני נרמול אפשר למצוא $w'b'$ שיתנו לנו את אותו ערך loss ולהפך.

עם זאת, כפי שראינו בהרצאה, ברמה הפרקטית, הביצועים יכולים להיפגע כאשר לא מנרמלים את הדאטה. עבור דאטהסט מנורמל/לא מנורמל נקבל גרדיאנט שונה שכן במקרה של huberloss הגרדיאנט תלוי גם בערך של הפיצירים ולא קבוע. עבור פיצירים שהסקלה שלהם גדולה יחסית, נצטרך learning rate נמוך יותר (כיוון שהנגזרת החלקית לפי המשקולת המתאימה לפיציר תהיה גדולה יותר) מאשר פיצירים שהסקלה שלהם יחסית קטנה. כתוצאה מכך, נצטרך לבחור η קטן מספיק לפיציר עם הסקאלה הכי גדולה. יכול להיווצר מצב שנבחר η נמוך מדי עבור הפיצירים עם הסקלה הקטנה וההתכנסות תהיה הרבה יותר איטית, כלומר נצטרך יותר איטרציות כדי להתכנס.

הוכחה:

נניח שיש לנו דאטהסט $X \in R^{m \times d}, y \in R^m, w \in R^d, b \in R$ ואנו מבצעים טרנספורמציה לינארית לכל אחד מהפיצירים כך ש $x'_i = a_i x_i + c_i, a_i \neq 0$. כלומר עבור כל $x_{org} \in X = (x_1, x_2, \dots, x_d)$ אנו מקבלים

$$x_{new} = (x'_1, x'_2, \dots, x'_d) = (a_1 x_1 + c_1, a_2 x_2 + c_2, \dots, a_d x_d + c_d)$$

נניח ש $w = (w_1, w_2, \dots, w_d)$ אזי נגדיר, $w' = (w'_1, w'_2, \dots, w'_d) = \left(\frac{w_1}{a_1}, \frac{w_2}{a_2}, \dots, \frac{w_d}{a_d}\right)$

ונקבל: $b' = b - \sum_{i=1}^d \frac{w_i a_i}{c_i}$

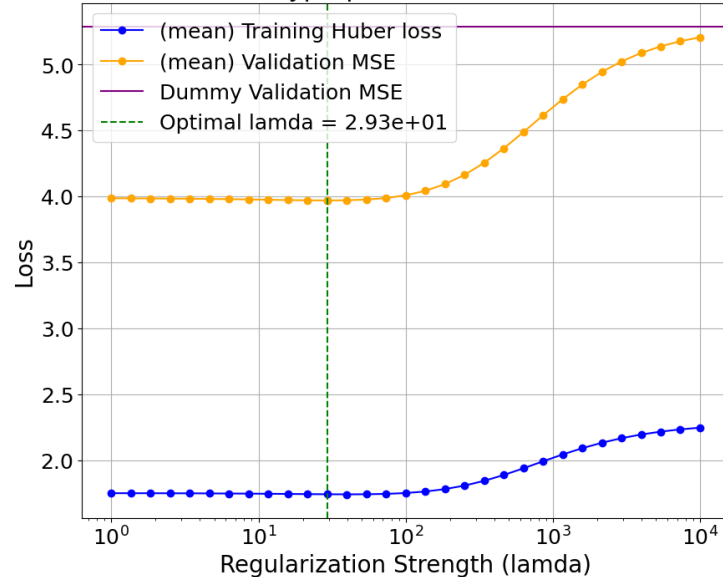
$$\begin{aligned} w'^T x_{new} + b' &= \sum_{i=1}^d w'_i x'_i + b' = \sum_{i=1}^d \frac{w_i}{a_i} * (a_i x_i + c_i) + b - \sum_{i=1}^d \frac{w_i a_i}{c_i} = \sum_{i=1}^d w_i x_i + \sum_{i=1}^d \frac{w_i a_i}{c_i} + b - \sum_{i=1}^d \frac{w_i a_i}{c_i} \\ &= \sum_{i=1}^d w_i x_i + b = w^T x_{old} + b \end{aligned}$$

Section 3: Ridge linear regression

(Q10)

הערכים שבחרנו לבדוק עבור regularization parameter הם 31 ערכים ברווחי log בטווח $10^0 - 10^4$ (לאחר שניסינו על טווח גדול יותר בהתחלה). להלן התוצאות:

Validation Curve for lamda Hyperparameter with Huber Loss (epsilon=3.65)



נשים לב כי עבור ערכי למדא קטנים, אין כמעט השפעה של הרגולריזציה על הביצועים של המודל, כאשר עבור ערכי למדא גדולים רואים כי loss גדל והמודל מתחיל להכנס כבר לunder fit, הן בחן train והן בחן test הביצועים מתחילים לרדת. זו התנהגות צפויה עבור הרגרסיה, שכן ראינו בתרגול כי ככל שאנו מגדילים את הפרמטר למדא, אנחנו מציבים מגבלה על גודל הנורמה האפשרי של argmin כך שהווקטור האופטימלי הולך ומתקרב ל0 כל שלמדא גדל כי נותנים יותר משקל על הנורמה ולכן יכולת ההכללה מתחילה לרדת בשלב הזה ומקבלים מודלים פשוטים מדי.
הלמדא האופטימלי שמצאנו באימון הוא:

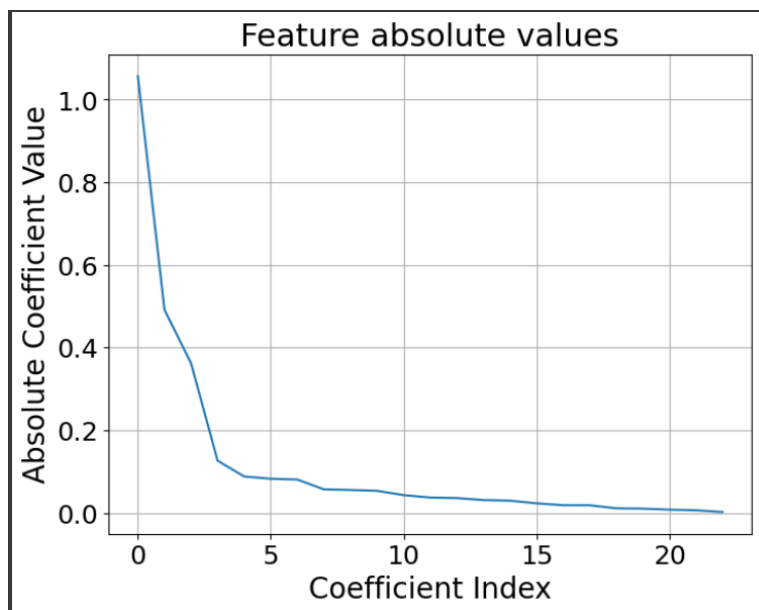
Optimal regularization strength (lamda): 2.93e+01
Training Huber loss at optimal lamda: 1.74(+/- 0.10)
Validation MSE at optimal lamda: 3.97(+/- 0.95)

(Q11)

Model	Section	Train Huber Loss (cross validated)	Train MSE (cross validated)	Valid MSE (cross validated)
Dummy	2	2.29(+/- 0.1)	5.27 (+/- 0.26)	5.29 (+/- 1.051)
Linear	2	1.73(+/- 0.1)		4.02 (+/- 1)
Ridge linear	3	1.74(+/- 0.1)		3.97 (+/- 0.95)

(Q12)

לא רשום מה להוסיף לדוח בשאלה הזאת, אבל אנחנו מניחות שאת זה:



```
sorted_indices: [ 4 10 14 11 1 15 13 7 5 3 9 2 8 17 21 22 6 16 0 12 19 20 18]
```

(Q13)

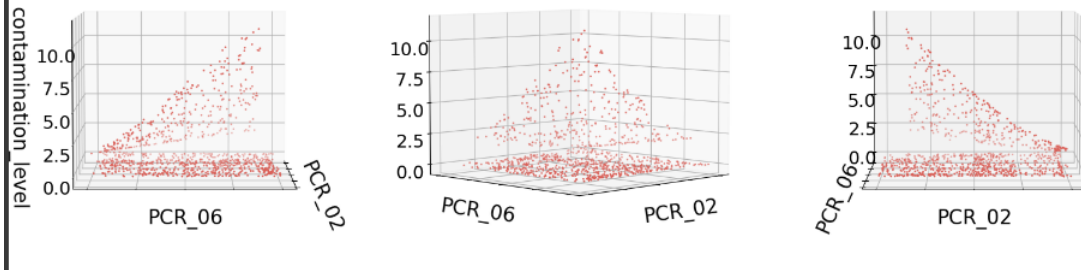
ראינו בתרגול וגם רואים בשאלה הקודמת שככל שמגדילים את הלמדא ברגולריזציה אנחנו מקבלים שהcoefficients הולכים ומתקרבים לאפס (לא מתאפסים) וזה כי אנחנו דורשים נורמה יותר ויותר קטנה על ידי הגדלת מקדם רגולריזציה על הנורמה (למדא) וזה גורם למשקלים לדעוך לכיוון אפס. ייתכן גם שהridge ייתן משקל גבוה יותר לפיצורים חשובים יותר (זה תלוי גם בנרמול הפיצורים) אך אם למדא גדול מאוד, יהיה קשה לפרש את החשיבות של כל פיצור כי כל המשקולות יהיו קטנים מאוד.

לכן אנחנו חושבות שהridge פחות מתאים למיקסום interpretability. לעומת זאת, בשימוש ברגולריזציית lasso נוכל למקסם יותר את interpretability כיוון שחלק מהמקדמים של הפיצורים מתאפסים. אנחנו עושים בעצם סוג feature selection די בקלות (בשיעורי בית 1 עבדנו מאוד קשה בשביל לעשות זאת). לתכונות עם מקדמים גדולים תהיה השפעה גדולה יותר על החיזוי מה שאומר שזה נותן להם חשיבות גדולה יותר במודל הרגרסייה וזה נותן לנו רמז אילו תכונות טובות לניבוי, לכן נרצה שתכונות שלא טובות לניבוי יתאפסו (עם זאת חשוב לנרמל את הדאטה כדי שזה יתקיים).

Section 4: Feature Mappings (visualization)

(Q14)

Contamination Level vs PCR_02 and PCR_06



אנו מחפשים להבין האם ניתן להשתמש במודל לינארי כדי לחזות את המשתנה הרציף contamination level, לכן נרצה להבין האם קיים איזושהו קשר לינארי בין משתני החיזוי למשתנה הפרדיקציה (אנחנו לא מחפשים הפרדה לינארית כי לא מדובר בסיווג).

ניתן לראות שיש "2 קבוצות" שונות של נקודות – הקבוצה האחת הן נקודות שמקבלות contamination level אפס או קרוב מאוד לאפס בלי תלות בערכים PCR_02 ו PCR_06 (כלומר מתנהגות כמו המישור הקבוע $z=0$) וקבוצה שניה יוצרות ושאר קירוב של היפר מישור עולה. לכן לא נראה שיש רק מישור אחד שיתאים בצורה טובה לחיזוי הדאטה. אם נאמן מודל רגרסיה על דאטה זה, נקבל מודל שינסה להחזיר את הערך הממוצע של y לכל x בקירוב (כי מדובר בבעיית LS עם שגיאה ריבועית) ולכן המישור שנקבל יהיה מישור ש"יפריד" בין 2 הקבוצות, כלומר יהיה ממוקם בפוזיציה שנמצאת בין המישור ל $z=0$ למישור המשופע, והוא לא יהיה קירוב טוב מאד לאף אחת מהקבוצות. (נשים לב שרוב הנקודות נמצאות ב-contamination level אפס ולכן המישור שנקבל יהיה יותר קרוב אליהם ויקבל שיפוע קטן כלפי הנקודות העולות). לאור מה שכתבנו, אנו חושבות שלמדל את הבעיה עם מודל רגרסיה לינארית לא יחזה מספיק את משתנה המטרה, ולכן אולי כדאי להשתמש בפונקציות קרנל כלשהן שיכולות למפות את הבעיה למימד בו היא מתאימה לחיזוי לינארי.

(Q15)

פה לא כל כך מובן מה להוסיף לדוח, החלטנו להוסיף את הקוד שמבצע את הפעולות הנדרשות.

```
delta_0_q_15= calculate_huber_delta(X_train_section4, y_train_section4)
Huber_regressor_model2 = HuberRegressor(alpha=best_lambda, epsilon=delta_0_q_15, fit_intercept=True)
Huber_regressor_model2.fit(X_train_section4, y_train_section4)

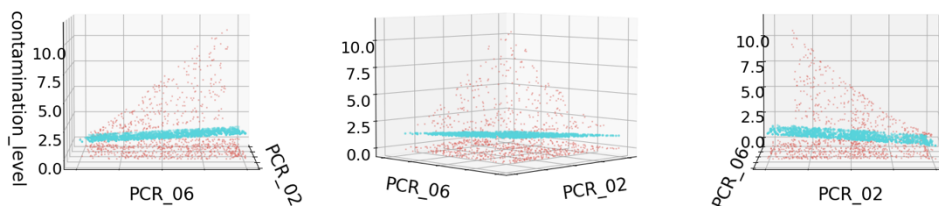
y_train_pred = Huber_regressor_model2.predict(X_train_section4)
y_test_pred = Huber_regressor_model2.predict(X_test_section4)
```

והתוצאות שקיבלנו על המודל:

```
delta_0: 3.39
Training Huber Loss: 2.2622476779753633
Training MSE: 5.066776810489589
Test MSE: 4.930536273026795
```

(Q16)

Contamination Level vs PCR_02 and PCR_06 and matched Huber Regressor



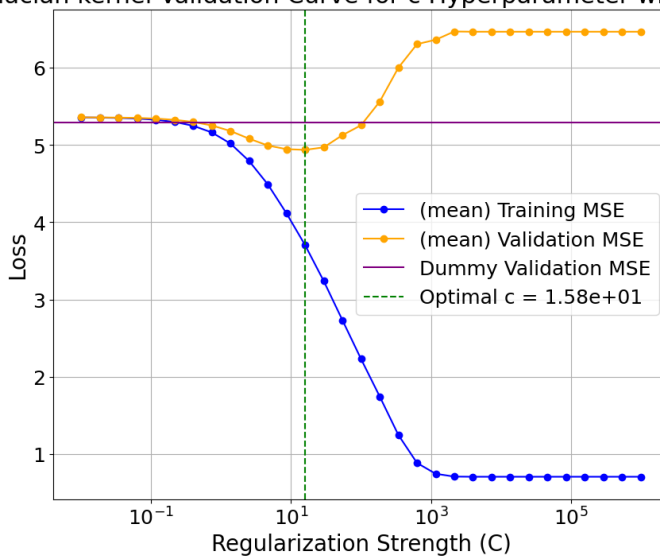
ניתן לראות שקיבלנו מודל שמתאים לציפיות שלנו מסעיף 14.

(Q17)

נדרשנו בסעיף זה לבצע hyperparameter Tuning for C על פי שני הפיצרים.

להלן הגרף שקיבלנו:

SVR laplacian kernel Validation Curve for c Hyperparameter with MSE Loss)



הפרמטר האופטימלי שקיבלנו + התוצאות שנתן (רק ע"פ 2 פיצרים):

```
Optimal regularization strength (c): 1.58e+01
Training MSE loss at optimal c: 3.70(+/- 0.26)
Validation MSE loss at optimal c: 4.94(+/- 1.02)
```

Optimal C = 15.8

התוצאות שקיבלנו הן צפויות, ככל שC קטן יותר יש יותר דגש על הקטנת הנורמה (מצב של underfit) וככל שC גדול יותר יש יותר דגש על צמצום השגיאה של מודל הSVR על קבוצת האימון (עד להגעה למצב overfit).

לאחר מכן עשינו cross validation יחיד עם הפרמטר האופטימלי שמצאנו $C = 15.8$ על כל קבוצת האימון ואלה התוצאות שקיבלנו:

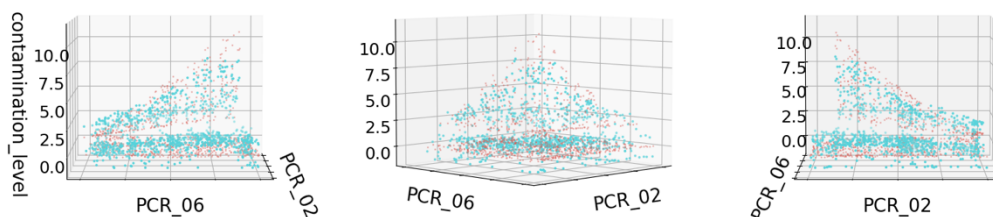
Training MSE loss at optimal c: 0.67(+/- 0.26)
Validation MSE loss at optimal c: 2.57(+/- 1.02)

אנחנו נשתמש בהם בטבלה, כפי שצוין בהוראות בפיאצה:

Model	Section	Train Huber Loss (cross validated)	Train MSE (cross validated)	Valid MSE (cross validated)
Dummy	2	2.29(+/- 0.1)	5.27 (+/- 0.26)	5.29 (+/- 1.051)
Linear	2	1.73(+/- 0.1)		4.02 (+/- 1)
Ridge linear	3	1.74(+/- 0.1)		3.97 (+/- 0.95)
SVR + Laplace	4		0.67(+/- 0.26)	2.57(+/- 1.02)

(Q18)

Contamination Level vs PCR_02 and PCR_06 and matched SVR Laplacian Regressor



(Q19)

מודל Huber regressor המתייחס למטריקה של Huber loss הוא מודל המאמן רגרסור לינארי המפגין עמידות כנגד outliers אך במקרה שלנו ניתן לראות ע"פ plot ומטריקות השגיאה כי החיזוי שלו על הנקודות יהיה פחות טוב (הסברנו בסעיף 14 שהדאטה פחות מתאים למודל לינארי).

לעומת זאת, מודל SVR + מיפוי קרנל מספק התאמה טובה יותר ל training set ולכד פרטים עדינים יותר. על אף שגם מודל svr מחפש חזאי לינארי, המיפוי של הפיצירים באמצעות קרנל הלפלסיאן מאפשר לייצר חזאי שאינו וכך ניתן לקבל ביצועים טובים יותר באיזורים עם דפוסים מורכבים כמו במקרה שלנו. מהמטריקות והplot ניתן לראות שאכן המודל מתאים עצמו בצורה מאד טובה ל train set (מה שגם יכול לגרום ל overfit) וכי שגיאת הולידציה נמוכה יותר.

אנחנו מסיקות שהשיפור בביצועים הגיע בעיקר מהמיפוי של הפיצירים באמצעות קרנל הלפלסיאן ולא שינוי המודל.

(Q20)

Model	Section	Train Loss (cross validated)	Valid MSE (cross validated)	Test MSE
Dummy	2	5.27 (+/- 0.26) (MSE)	5.29 (+/- 1.051)	5.160304602563259
Linear	2	1.73(+/- 0.1) (Huber)	4.02 (+/- 1)	3.748605086090574
Ridge linear	3	1.74(+/- 0.1) (Huber)	3.97 (+/- 0.95)	3.7388682486726066
SVR + Laplace	4	0.67(+/- 0.26) (MSE)	2.57(+/- 1.02)	2.3622709697094715

המודל הטוב ביותר שקיבלנו הוא SVR + Laplace עם שגיאה ריבועית ממוצעת הכי נמוכה על test.