

Intro to machine learning  
Major home work #3  
presenting:  
Yitzchak Grinboim – 318837317  
Sasson Shmuel Lamey - 325172351

Section 1

Q1:

$$\begin{aligned} \frac{\partial \ell(\mathbf{w}, b; x_i, y_i)}{\partial b} &= \left\{ \frac{\partial}{\partial b} \left( \frac{1}{2} (w^T \cdot x_i + b - y_i)^2 \right) \right\} \stackrel{\text{chain rule}}{=} \left\{ \begin{aligned} &(w^T \cdot x_i + b - y_i) \cdot \frac{\partial}{\partial b} (w^T \cdot x_i + b - y_i) = \\ &\delta \cdot \frac{\partial}{\partial b} (\sqrt{(w^T \cdot x_i + b - y_i)^2}) \end{aligned} \right\} = \\ &\left\{ \delta \cdot \left( \frac{1}{2\sqrt{(w^T \cdot x_i + b - y_i)^2}} \right) \cdot \frac{\partial}{\partial b} ((w^T \cdot x_i + b - y_i)^2) \right\} = \left\{ \delta \cdot \left( \frac{1}{2\sqrt{(w^T \cdot x_i + b - y_i)^2}} \right) \cdot 2(w^T \cdot x_i + b - y_i) \cdot \frac{\partial}{\partial b} (w^T \cdot x_i + b - y_i) \right\} = \\ &\left\{ \delta \cdot \left( \frac{2(w^T \cdot x_i + b - y_i)}{2\sqrt{(w^T \cdot x_i + b - y_i)^2}} \right) \right\} = \left\{ \delta \cdot \left( \frac{w^T \cdot x_i + b - y_i}{|w^T \cdot x_i + b - y_i|} \right) \right\} = \left\{ \delta \cdot \text{sign}(w^T \cdot x_i + b - y_i) \right\} \end{aligned}$$

Where the upper row was for the case when  $|w^T \cdot x_i - y_i + b| \leq \delta$  and the lower otherwise.

Q2:

Some auxiliary calculations:

$$\begin{aligned} &= \begin{bmatrix} x_{1_1} & \cdots & x_{m_1} \\ \vdots & \ddots & \vdots \\ x_{1_d} & \cdots & x_{m_d} \end{bmatrix} \cdot \begin{bmatrix} \left( (w^T \cdot x_1 + b - y_1) \cdot 1_{|w^T \cdot x_1 + b - y_1| \leq \delta} \right) + \left( (w^T \cdot x_1 + b - y_1) \cdot 1_{|w^T \cdot x_1 + b - y_1| > \delta} \right) \\ \vdots \\ \left( (w^T \cdot x_m + b - y_m) \cdot 1_{|w^T \cdot x_m + b - y_m| \leq \delta} \right) + \left( (w^T \cdot x_m + b - y_m) \cdot 1_{|w^T \cdot x_m + b - y_m| > \delta} \right) \end{bmatrix} \\ &= \frac{1}{m} \sum_{i=1}^m (w^T \cdot x_i + b - y_i) \cdot x_i \stackrel{\text{linearity of } x_i^T}{=} \begin{bmatrix} x_{1_1} & \cdots & x_{m_1} \\ \vdots & \ddots & \vdots \\ x_{1_d} & \cdots & x_{m_d} \end{bmatrix} \cdot \begin{bmatrix} w^T \cdot x_1 + b - y_1 \\ \vdots \\ w^T \cdot x_m + b - y_m \end{bmatrix} = \frac{1}{m} \begin{bmatrix} x_{1_1} & \cdots & x_{m_1} \\ \vdots & \ddots & \vdots \\ x_{1_d} & \cdots & x_{m_d} \end{bmatrix} \left( \begin{bmatrix} w^T \cdot x_1 \\ \vdots \\ w^T \cdot x_m \end{bmatrix} + \begin{bmatrix} b \\ \vdots \\ b \end{bmatrix} - \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \right) \\ &= \frac{1}{m} \begin{bmatrix} x_{1_1} & \cdots & x_{m_1} \\ \vdots & \ddots & \vdots \\ x_{1_d} & \cdots & x_{m_d} \end{bmatrix} \left( \begin{bmatrix} x_{1_1} & \cdots & x_{1_d} \\ \vdots & \ddots & \vdots \\ x_{m_1} & \cdots & x_{m_d} \end{bmatrix} \cdot \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} + \begin{bmatrix} b \\ \vdots \\ b \end{bmatrix} - \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \right) = \frac{1}{m} \cdot X^T (X \cdot w + b \cdot \mathbf{1}_m - y) \\ &\frac{1}{m} \sum_{i=1}^m \delta \cdot \text{sign}(w^T \cdot x_i + b - y_i) \cdot x_i \stackrel{\text{a linear combination of } x_i^T}{=} \frac{\delta}{m} \begin{bmatrix} x_{1_1} & \cdots & x_{m_1} \\ \vdots & \ddots & \vdots \\ x_{1_d} & \cdots & x_{m_d} \end{bmatrix} \cdot \begin{bmatrix} \text{sign}(w^T \cdot x_1 + b - y_1) \\ \vdots \\ \text{sign}(w^T \cdot x_m + b - y_m) \end{bmatrix} \\ &= \frac{\delta}{m} \begin{bmatrix} x_{1_1} & \cdots & x_{m_1} \\ \vdots & \ddots & \vdots \\ x_{1_d} & \cdots & x_{m_d} \end{bmatrix} \cdot \begin{bmatrix} 1_{w^T \cdot x_1 + b - y_1 \geq 0} \\ \vdots \\ 1_{w^T \cdot x_m + b - y_m \geq 0} \end{bmatrix} = \frac{\delta}{m} \begin{bmatrix} x_{1_1} & \cdots & x_{m_1} \\ \vdots & \ddots & \vdots \\ x_{1_d} & \cdots & x_{m_d} \end{bmatrix} \cdot \text{sign} \left( \begin{bmatrix} w^T \cdot x_1 + b - y_1 \\ \vdots \\ w^T \cdot x_m + b - y_m \end{bmatrix} \right) \\ &= \frac{\delta}{m} \begin{bmatrix} x_{1_1} & \cdots & x_{m_1} \\ \vdots & \ddots & \vdots \\ x_{1_d} & \cdots & x_{m_d} \end{bmatrix} \cdot \text{sign} \left( \begin{bmatrix} x_{1_1} & \cdots & x_{1_d} \\ \vdots & \ddots & \vdots \\ x_{m_1} & \cdots & x_{m_d} \end{bmatrix} \cdot \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} + \begin{bmatrix} b \\ \vdots \\ b \end{bmatrix} - \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \right) = \frac{\delta}{m} \cdot X^T \cdot \text{sign}(X \cdot w + b \cdot \mathbf{1}_m - y) \end{aligned}$$

The 3 upper rows are when naively assuming that for any  $i$ :  $|w^T \cdot x_i - y_i + b| \leq \delta$ , and the lower 2 are when assuming naively that for any  $i$   $|w^T \cdot x_i - y_i + b| > \delta$ . But for some  $i$  it holds and for some it does not.

really it is like this:

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m \left( \left( (w^T \cdot x_i + b - y_i) \cdot 1_{|w^T \cdot x_i + b - y_i| \leq \delta} \right) + \left( (w^T \cdot x_i + b - y_i) \cdot 1_{|w^T \cdot x_i + b - y_i| > \delta} \right) \right) \cdot x_i$$

this is a linear combination of  $x_i^T$ 's and the coefficients depend on  $|w^T \cdot x_i + b - y_i|$ .

$$= \begin{bmatrix} x_{1_1} & \cdots & x_{m_1} \\ \vdots & \ddots & \vdots \\ x_{1_d} & \cdots & x_{m_d} \end{bmatrix} \cdot \begin{bmatrix} \left( (w^T \cdot x_1 + b - y_1) \cdot 1_{|w^T \cdot x_1 + b - y_1| \leq \delta} \right) + \left( (w^T \cdot x_1 + b - y_1) \cdot 1_{|w^T \cdot x_1 + b - y_1| > \delta} \right) \\ \vdots \\ \left( (w^T \cdot x_m + b - y_m) \cdot 1_{|w^T \cdot x_m + b - y_m| \leq \delta} \right) + \left( (w^T \cdot x_m + b - y_m) \cdot 1_{|w^T \cdot x_m + b - y_m| > \delta} \right) \end{bmatrix}$$

So we can conclude for any entry in the result vector:

$$\nabla_w \mathcal{L}_{\mathcal{H}}(w, b) = \frac{1}{m} X^T \cdot \begin{bmatrix} Z_1 \\ \vdots \\ Z_m \end{bmatrix}$$

$$\text{when } Z_i = \begin{cases} [X \cdot w + b \cdot 1_m - y]_i, & |w^T \cdot x_i + b - y_i| \leq \delta \\ \delta \cdot [\text{sign}(X \cdot w + b \cdot 1_m - y)]_i, & |w^T \cdot x_i + b - y_i| > \delta \end{cases}$$

Q3:

We want to find a good delta which means we want it to be a threshold which will cause only for outliers (which are estimated to be a small percentage of the data) to be calculated using the  $\delta \cdot \text{sign}(w^T \cdot x_i + b - y_i)$ . so we will have a belief p, which is believed to be the percentage of outliers. And we will calculate a regular regression, then save the residuals and look at the 100p% of the biggest residuals. Then we can get the smallest of those big residuals, and select it as our threshold, meaning we will take the square root of that:

```
def good_delta(X, Y, believed_p = 0.05):
    n = len(Y) # Number of data points
    residuals = []

    w, b = ordinary_least_squares(X, Y)

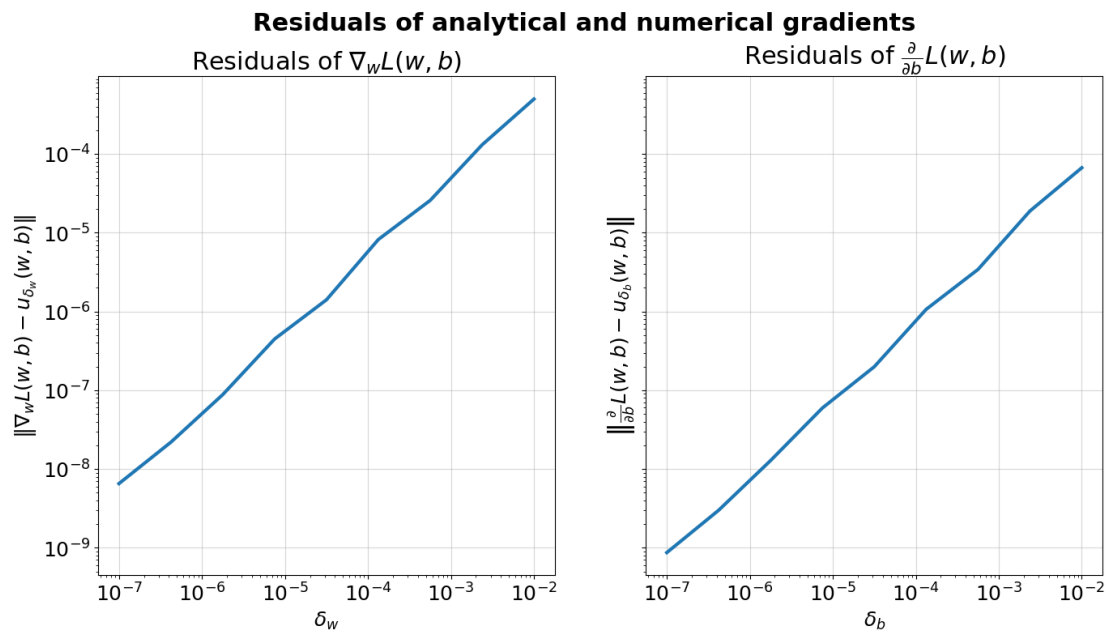
    for i in range(n):
        residual = (w.T * X[i] + b - Y[i])**2
        residuals.append(residual)

    residuals.sort()

    percentile_index = floor((1-believed_p) * n)
    delta = residuals[percentile_index]

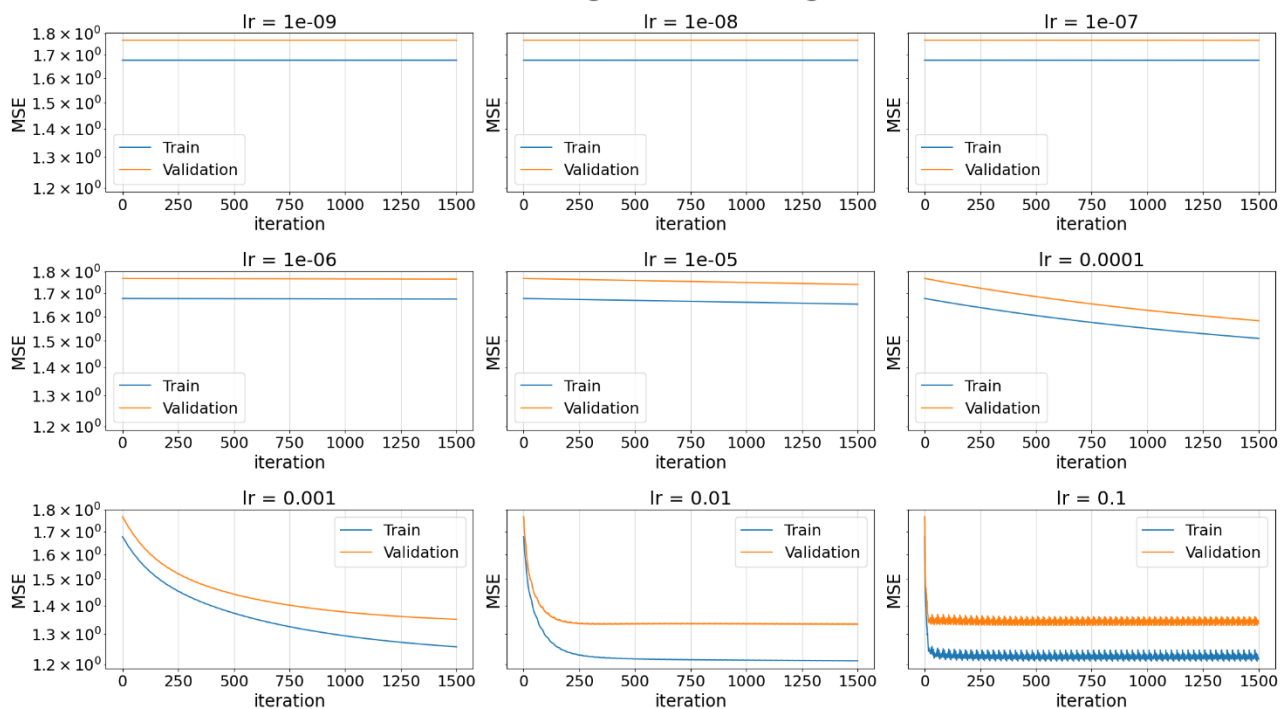
    return delta
```

Q4:



Q5:

## robust regression learning rates



As can be seen in the graphs the best learning rate that gives the best convergence is 0.01. This can be simply explained by saying that for a small learning rate 1500 steps might not be enough steps to reach convergence (as can be seen again, in the graphs) and for a big learning rate, it might be too big in a way which will make it 'hop' around the minimum point (or even worse, but in the graphs we see the 'hop' (hop tralala)) and not converge. Enlarging the number of iterations doesn't seem to make a lot of sense because it looks like it already converged, and has nowhere to go anymore. Furthermore, it might start to 'hop' if too many iterations had been made.

Q6:

The Huber loss can be best exploited when there are outliers in the data, because it offers a middle ground between the sensitivity of MSE to outliers and the alternative, not to regard outliers, which can sometimes be a mistake because the model can't really know which are outliers. Unlike ordinary least squares, which can be heavily effected by outliers, Huber loss minimizes their impact by giving a different loss value for large residuals values. This makes it particularly useful in scenarios with noisy data or just data with little outliers, resulting in giving us more reliable regression models.

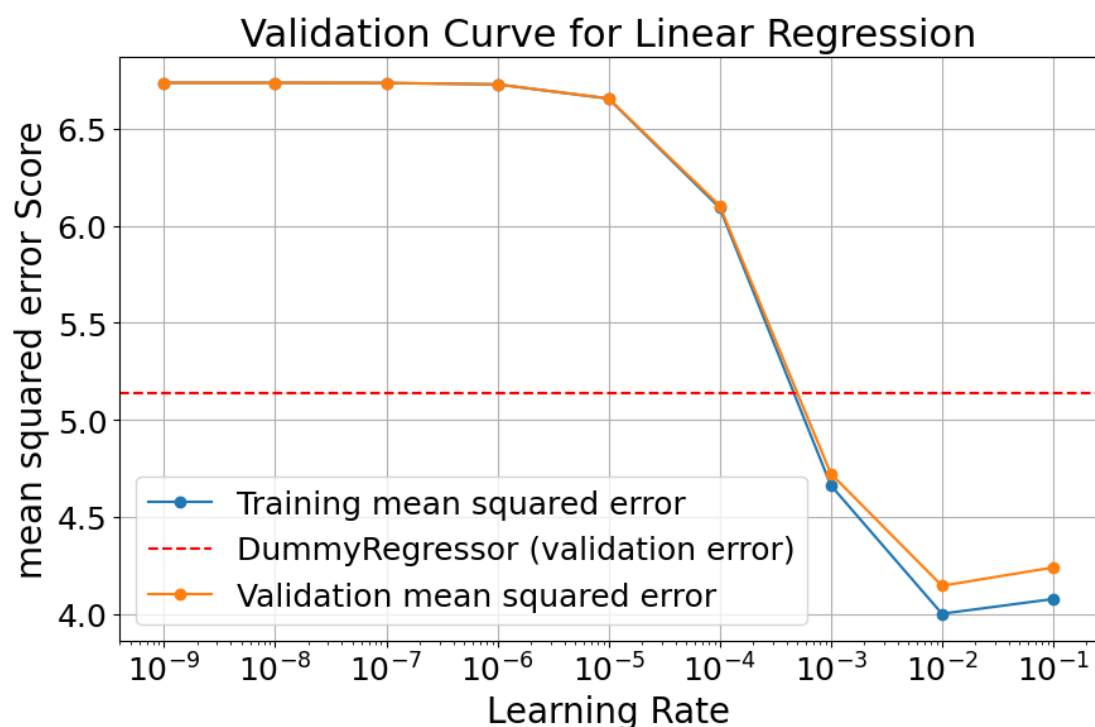
## Section 2:

Q7:

The errors are:

Model	Section	Train huber loss	Valid MSE
		Cross validated	
Dummy	2	1.62162	5.13973

Q8:



Model	section	Train huber loss	Valid mse
		Cross validated	
dummy	2	1.62162	5.13973
Linear regressor	2	1.23043	4.14708

Q9:

for the linear regression model it is likely that not normalizing beforehand would not matter very much. That is because for linear regression, every change that you make in the scaling of some feature, will be apparent accordingly in the way of predicting, when you come to predict a new  $x$ , it will have that feature in the same scale. But for models that use gradient descent for optimization it does matter when deciding on the direction from the gradient vector, because different scaling when not normalized will make it biased more towards the bigger values for no real good reason.

### Section 3:

Q10:



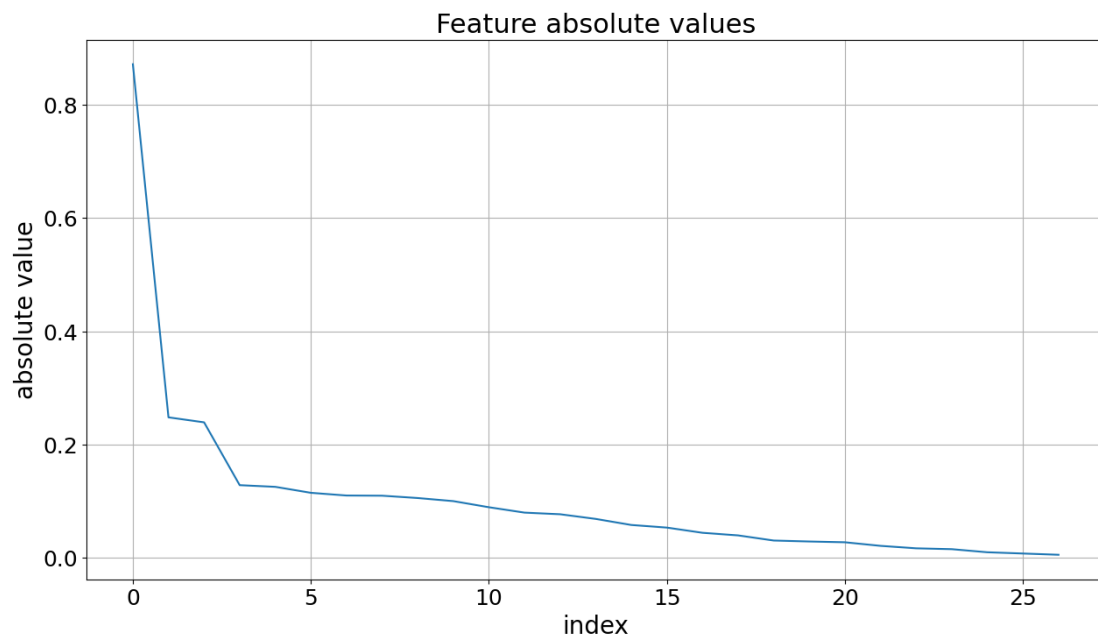
Optimal parameter we got is: 0.001

Validation accuracy is 4.079

Q11:

Model	section	Train huber loss	Valid MSE
Cross validated			
dummy	2	1.62162	5.13973
Linear regressor	2	1.23043	4.14708
Ridge Linear	3	1.22534	4.14708

Q12:



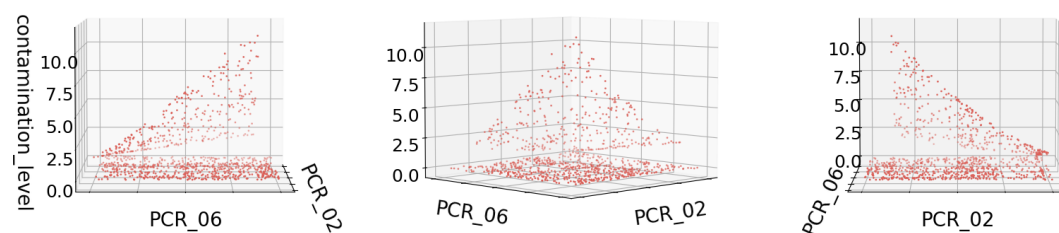
Q13:

Ridge regularization is not a good choice because it does not have the power to leave out feature by reducing their coefficients to zero. It can only lower them close to zero, but then if a feature has very large scale it can still effect the interpretability and there is no way using that to understand the roles of different features in determining behavior. Lasso regression can reduce coefficients to zero thus helps to find a subset of features that effect the model the most, and effectively improving the interp retability.

#### Section 4 Feature Mappings (visualization):

Q14:

Contamination level as a function of PCR\_02 and PCR\_06

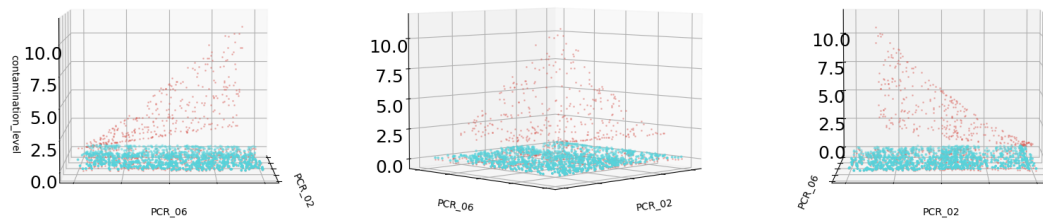


Explanation: we can see the there is some partition of the data such that in one part there

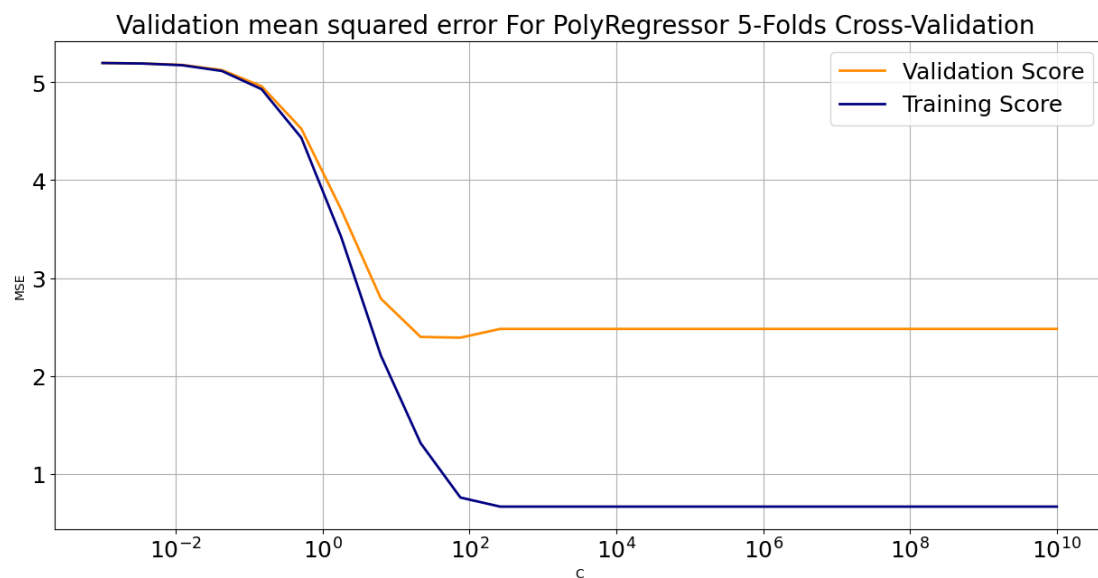
is no affect of combination of PCR2 and PCR6. But in the other part we see a pattern of relationship between them. So we might want to create some feature mapping to a have a better regression.

Q15-16:

Contamination level as a function of PCR\_02 and PCR\_06



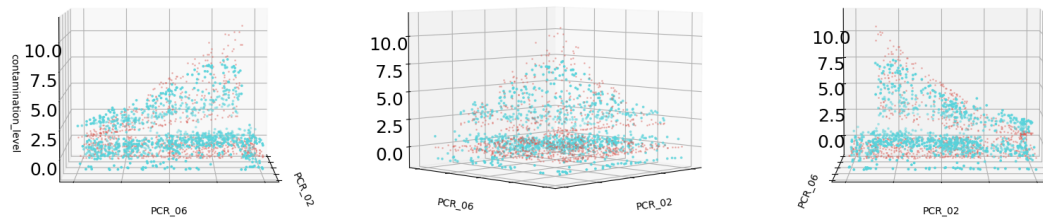
Q17:



*best C value is 74.98942093324558 ,its train accuracy is 0.7595527767361085  
and its validation accuracy is 2.391672498620812*

Q18:

## Contamination level as a function of PCR\_02 and PCR\_06



Q19:

It seems the last model was better in fitting the data when using the poly\_reg. We can see from the graph that there is better correlation between this model predictions and the actual values (than the former graph/model).

We could also see that validation accuracy was better (by almost a factor of 2).

Q20:

Model	section	Train huber loss	Valid MSE	TEST MSE
		Cross validated		RETRAINED
dummy	2	1.62162	5.13973	5.320209
Linear regressor	2	1.23043	4.14708	4.338360
Ridge Linear	3	1.22534	4.14708	4.20895
SVR + Laplace	4	0.75	2.39	0.818121