



# INTROUCTION TO MACHINE LEARNING - 236756

MAJOR HW4

Itay Lasch & Yonatan Battat

## Section 1: Linear regression implementation

A1:

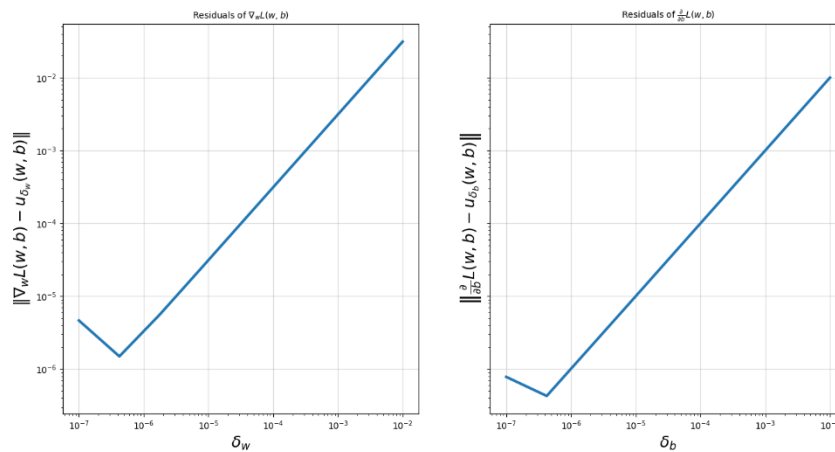
The MSE loss function for our regressor:

$$L(\underline{w}, b) = \frac{1}{m} \left\| X\underline{w} + \underline{1}_m \cdot b - \underline{y} \right\|_2^2$$

$$\begin{aligned} \frac{\partial}{\partial b} L(\underline{w}, b) &= \frac{2}{m} \cdot (\underline{1}_m)^T \cdot (X\underline{w} + \underline{1}_m \cdot b - \underline{y}) = \frac{2}{m} \cdot (\underline{1}_m)^T \cdot (X\underline{w} - \underline{y}) + \frac{2}{m} \cdot m b \\ &= \frac{2}{m} \cdot (\underline{1}_m)^T \cdot (X\underline{w} - \underline{y}) + 2b \end{aligned}$$

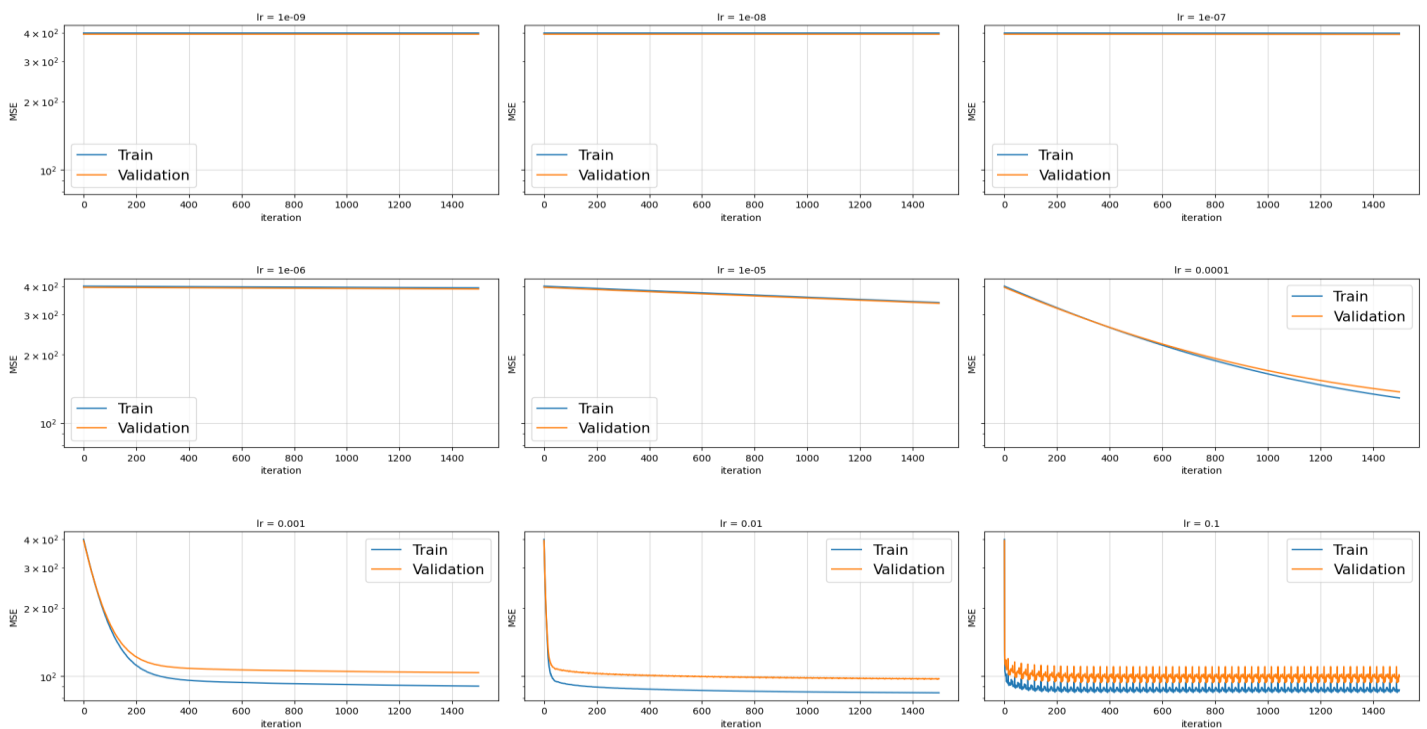
A2:

**Residuals of analytical and numerical gradients**



A3:

**Loss per iteration for different step size(lr)**



As we can see from the graphs above, for  $lr \in \{1e-9, \dots, 1e-6\}$  the loss is almost constant and there is a very small improvement per iteration. The reason is that the step size is too small and therefore  $w, b$  barely change between iterations, since our number of iterations is limited ( $= 1500$ ) we do not converge to ideal results for these  $lr$  values. Moreover for  $lr = \{1e-9, 1e-8, 1e-7\}$  we tested the loss per iteration with 10000 iterations and still saw that there is no change in the loss and no improvement per iteration which made us conclude that for these small  $lr$  values there are numerical errors that convert the steps to zeros. For  $lr$  values higher than  $1e-5$  including we can see that MSE is starting to decrease which means there is an improvement in loss per iteration. The higher  $lr$ , in particular  $lr = 1e-3, 1e-2$ , the decrease in loss is bigger and it converges faster to its minimum value.

We can see for  $lr = 0.1$  there is not convergence to a minimum value, this is happening because the step size is too big which makes MSE jump between different ideal  $w$ 's and not reach the actual minimum.

The best  $lr$  size in our opinion is  $lr = 0.01$ . As we explained it converges to minimum value very quick and although we saw that for  $lr = 0.1$  we get even a smaller train and validation loss, we prefer  $lr = 0.01$  since unlike  $lr = 0$ , it converges to a single minimum value and the difference in losses between the two  $lr$ 's is quite small.

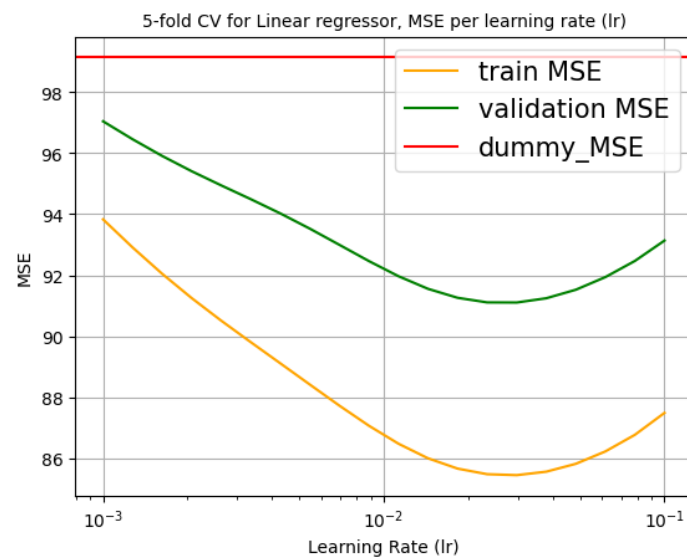
We don't think it makes sense to increase the number of iterations, we can see that as we approach the minimum value the slope of the validation MSE is getting smaller and smaller. That happens because as we get closer to the minimum value the gradient gets smaller as well, so each step is smaller than the one before him. Thus, adding more iterations won't make a significant change in the minimum validation loss (We verified this with 10000 iterations).

## Section 2: Evaluation and Baseline

A4:

Model	Section	Train MSE	Valid MSE
		Cross validated	
Dummy	2	99.13589	99.3029

A5:



The optimal  $lr$  is  $lr \cong 0.0297$ , which results with validation MSE of  $\cong 91.106$

Model	Section	Train MSE	Valid MSE
		Cross validated	
Dummy	2	99.13589	99.3029
Linear	2	85.45293	91.10669

A6:

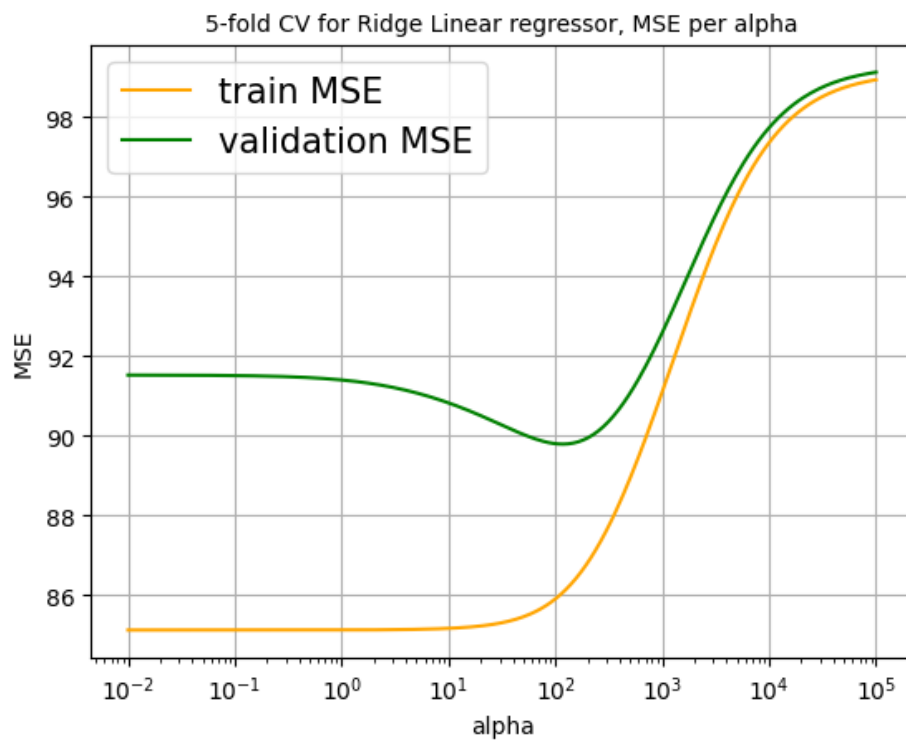
Had we chosen not to normalize the features, it wouldn't impact the Dummy model at all since it is always returns the mean value of the target variable 'contamination\_level' as a prediction to all samples. Since the target variable is not affected by normalization the MSE will stay the same.

As for the linear model, we wouldn't see a difference as well. As we saw in the lectures LS problem is convex and has a global minimum, and normalizing preserves convexity since it performs convexity preserving operations such as scaling and summing. Since our problem is convex with or without normalizing the features, with the right learning step and enough iterations we are guaranteed to reach that global minimum. So, both models' training performances will not change.

\* Each model will probably need different number of iterations and different learning rate which will yield different  $w$ .

### Section 3: Ridge linear regression

A7:



The optimal  $\alpha$  is  $\alpha \cong 118.46$ , which results with validation MSE of  $\cong 89.79$

A8:

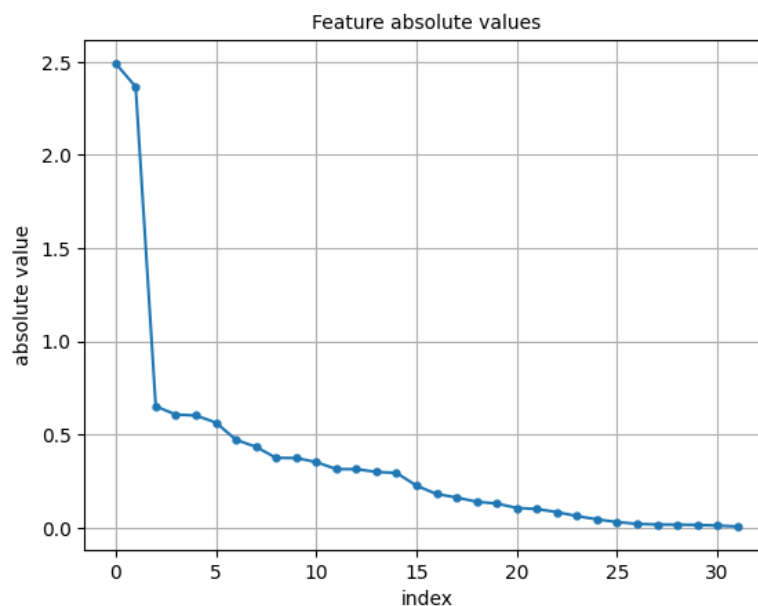
Model	Section	Train MSE	Valid MSE
		Cross validated	
Dummy	2	99.13589	99.3029
Linear	2	85.45293	91.10669
Ridge Linear	3	86.07993	89.7909

A9:

The 5 features with the largest coefficients (in absolute value) are:

1. PCR\_01 – 2.49206403
2. sugar\_levels – 2.36611658
3. household\_income – 0.6518341
4. sport\_activity – 0.60772657
5. happiness\_score – 0.60249684

A10:



A11:

The magnitude of each coefficient is interesting because as we learned each feature gets a coefficient based on how much correlated it is with our target variable, i.e. the more expressive and meaningful a feature is for our prediction the bigger its corresponding coefficient would be.

This can help us for example with feature selection when we want to decrease the number of dimensions in our problem, in order lower the memory and time complexity without losing a lot accuracy wise.

A12:

If we had not normalized the training performance of our linear ridge regressor would suffer and might miss the optimal solution. There are a few problems that led us to this answer.

First, as we learned in the tutorial:  $w^{Ridge} = \operatorname{argmin}_w \frac{1}{m} (\sum_{i=1}^m (y_i - w^T x_i)^2 + \lambda \|w\|_2^2)$

$= \operatorname{argmin}_w \frac{1}{m} (\sum_{i=1}^m (y_i - w^T x_i)^2), s.t. \|w\|_2^2 \leq c.$

In words, it means that we put constraints on the norm of  $w$  which can be seen as constraints on the coefficients of the features of sample  $x$ .

Let's look at a dataset that has different features with different scales. Now we want to look at two different features, one on a large scale and the other on a small scale. There is a large difference in scales between these two features, in order to compensate for that and let both of them have an impact on the prediction we are likely to give the small scale feature a bigger coefficient, but due to the constraint on the coefficients this solution might won't be in our possible solutions space and therefore we won't be able to get that result.

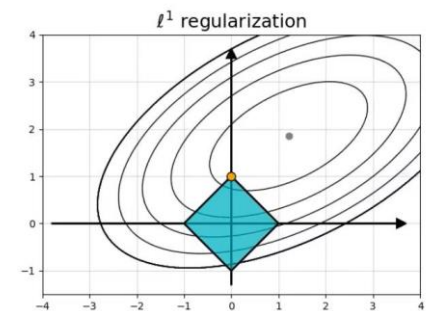
Another problem is that Ridge regularization uses  $l_2$  norm which is sensitive to different scales of the features. If we are deciding to not normalize our data, the features with bigger scales will have a larger impact on the regularization part of the ERM loss and therefore would get smaller coefficients.

As we explained in A11 the magnitude coefficient is directly proportional to how much the corresponding feature is correlated\uncorrelated with the target variable. Unregularized data could lead our regressor to yield a suboptimal classifier that his coefficients are decided based on the scales of the features instead of how much they actually impact our predictions.

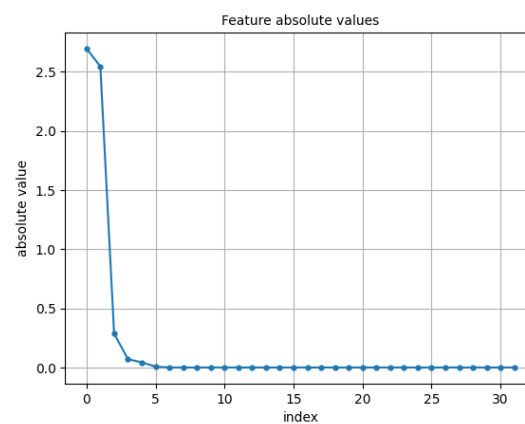
A13:

If we had used Lasso regressor instead of Ridge I would expect most of the coefficients to zero out. As we saw in tutorial 9, Lasso regressor uses  $l_1$  norm, and so our possible plain space is different than the one of Ridge regressor. To understand more visually we can look at the plot that shows the ERM loss as function of  $w$  on a dataset with 2 features.

As we can see there are a few options to choose a proper  $\underline{w}$  that will yield us the minimum ERM loss, but since we use Ridge/Lasso regression we also emphasize on the norm of  $\underline{w}$ . Therefore, in Lasso regression, which uses the  $l_1$  norm, SGD will choose the  $\underline{w}^*$  in our possible option space that has the smallest norm of them all. In order to achieve that the SGD will yield a  $\underline{w}^*$  that most of its features are zero out i.e. their coefficients are zero.



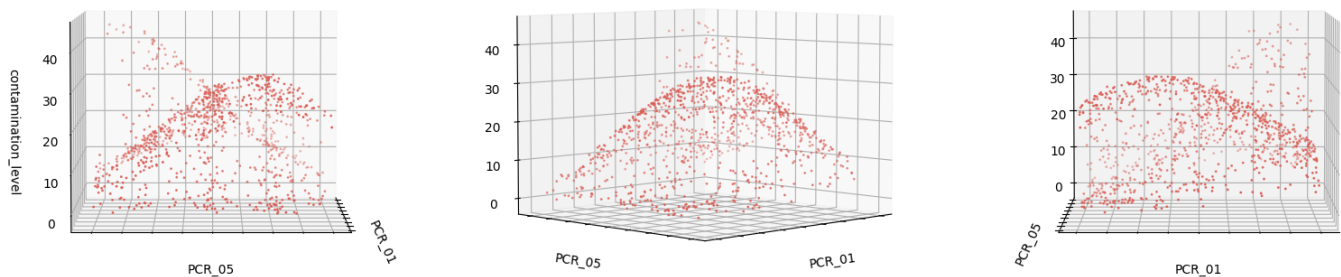
Just to clarify here is a coefficient graph if we were actually using Lasso regression:



## Section 4: Polynomial fitting (visualization)

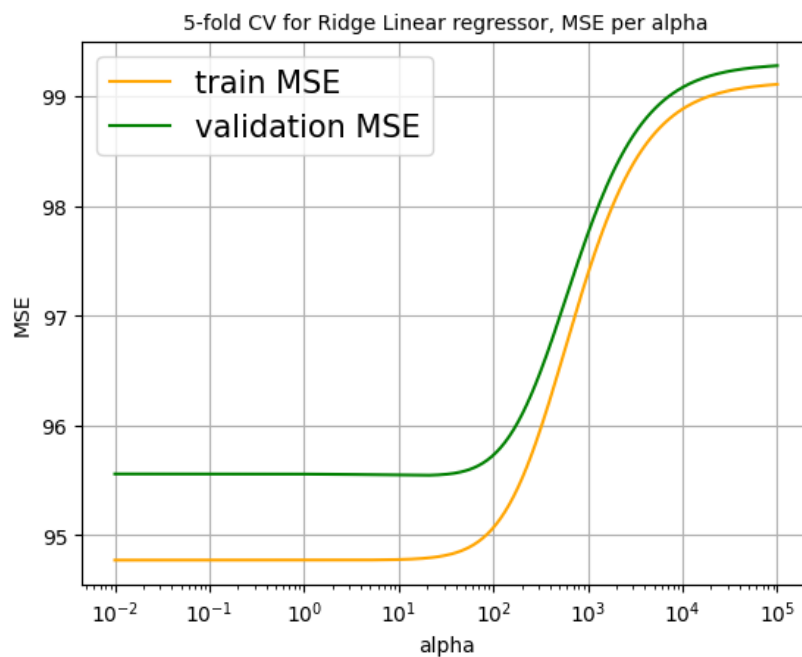
A14:

contamination level as function of PCR\_01 & PCR\_05



As we can see from the graphs above, it seems that the connection between 'contamination\_level' and the features 'PCR\_01', 'PCR\_05' cannot be expressed using a straight line and therefore linear regression won't work on the data as it currently is. It seems that the dots are scattered in some paraboloid shape so linear regression with polonomial feature mapping might work for us.

A15:

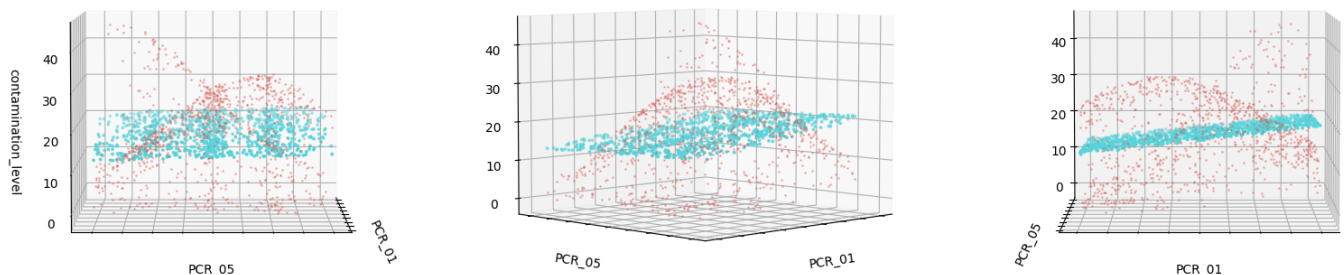


The optimal  $\alpha$  is  $\alpha \approx 16.119$  which results with validation MSE of 95.5436.



A16:

contamination level as function of PCR\_01 & PCR\_05 vs. predictions



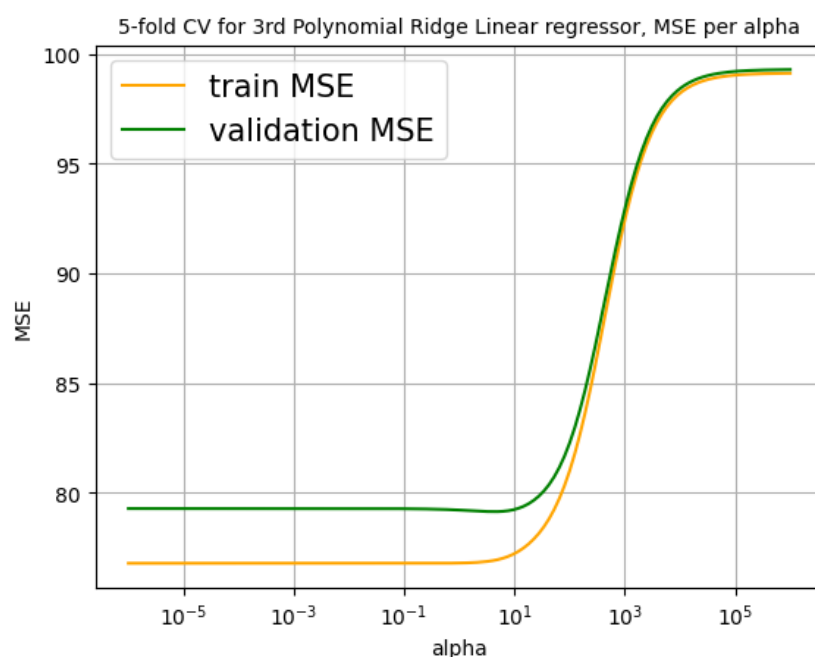
A17:

It is important to normalize the features once more after the polynomial feature mapping for the same reason explained in (A12). We start with the 2 features 'PCR\_01', 'PCR\_05' let's mark them as  $x_1, x_2$  appropriately. Now after the 3<sup>rd</sup> degree polynomial feature mapping the features are:  $x_1, x_2, x_1^2, x_2^2, x_1 \cdot x_2, x_1^3, x_2^3, x_1^2 \cdot x_2, x_1 \cdot x_2^2$ .

Since we normalized the features with MinMax normalization to be in between (-1,1) prior to the feature mapping, all the new features will be much smaller than the original ones.

To prevent the problem that occur from different scales in the features, as explained in (A12), it is important to normalize the features to the same scale once more.

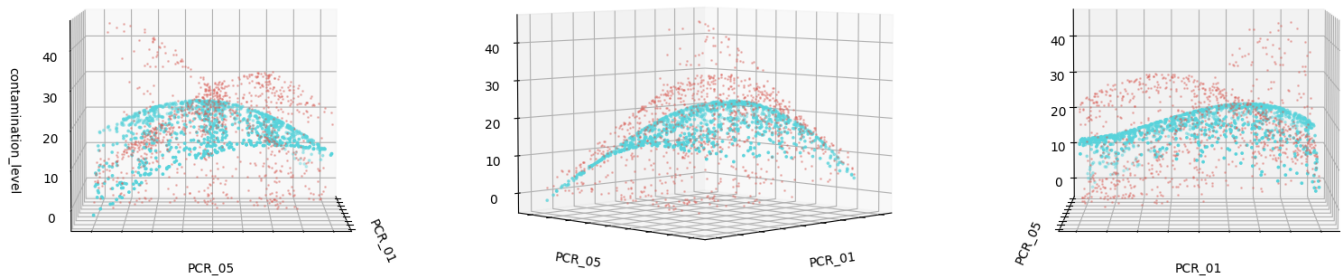
A18:



The optimal  $\alpha$  is  $\alpha \cong 4.023$  which results with validation MSE of 79.1425.

A19:

## Predictions vs. actual labels while using Polynomial mapping



A20:

There are a few different aspects to discuss when comparing the two different models:

First, as we can see after using the 3<sup>rd</sup> degree polynomial feature mapping our validation MSE was lower than the model on the original data (79.14 compared to 95.54), which means our predictions were closer to the actual label after the feature mapping.

Also we can see from the scattered plots in (A16) and (A19) that on the original data set our predictions created a straight plane that crossed the actual target variable values. On the other hand, in the polynomial model, we can see that even though our predictions were not perfect, most of the time they were pretty close to the actual values and depicted the same shapes as them.

The reason for those differences is that the feature mapping added complexity to the model and enabled it to fit the data more accurately.

## Section 5: Fitting Gradient Boosted Machines (GBM)

### Task before Q21:

The preprocessing we did is the one we did in the beginning of the exercise and is identical to the preprocessing we did in previous exercises and for the same reasons. We normalized the features using StandardScaler and MinMaxScaler in the following manner:

MinMaxScaler: PCR\_01, PCR\_02, PCR\_03, PCR\_09, sport\_activity

StandardScaler: PCR\_04, PCR\_05, PCR\_10, sugar\_levels

### A21:

Squared error – As we learned in the lecture this loss function tries to fit the residuals as much as possible at each iteration. This is relevant for us since we use bounded decision trees as weak learners and therefore at each iteration, we won't be able to fit the data perfectly and we would have residuals between our predictions and the actual target variable. The squared error loss function penalizes large deviations from the target variable while neglecting small residuals.

By emphasizing the residuals at each iteration, we are minimizing the MSE and with enough iterations can get good results.

Absolute error – Using this loss function, we are trying to minimize the mean of absolute differences between our predictions and the actual target variable. As opposed to the squared error loss function the absolute error is more robust to outliers because it doesn't magnify the results of large residuals.

It's important to note that this loss function is not differentiable at zero, but we can overcome that using methods like SGD.

Overall it seems that this loss function will perform well on our data and can be relevant for our cause.

Hubert – The Hubert loss function combines the two above loss functions by defining a threshold of the residuals. This way we can enjoy both worlds, for small residuals in prediction, under some  $\delta$  we would want to use the squared error function to neglect those residuals and for residuals above this  $\delta$  we would use the absolute error function to target them as our main objective to minimize with robustness to outliers.

We decided above that both MSE and MAE loss function are relevant

Quantile – The quantile loss function behaves a bit differently than the rest of the loss functions above. It does not try to minimize the MSE or MAE, but instead it tries to minimize a loss function that captures the difference between the predicted quantiles and the actual quantiles. Therefore, for each sample,  $h(x_i)$  is supposed to return a value in  $[0,1]$  which will represent the quantile of  $x_i$ , i.e. what is the percentage of samples that their contamination level is less than  $x_i$ 's contamination level. For this to work we need to hold a quantile for each  $x_i$  as well as it's contamination level variable. This loss function works with values between  $[0,1]$  and we want to work with values across all  $\mathbb{R}$ .

Since this loss function requires unnecessary calculations of quantiles and doesn't work with values in our preferred domain, we find it less relevant for our scenario and would probably not work with it.

#### A22:

Let's see the similarities and differences between Random Forests and Gradient Boosted Trees:

Similarities:

- Both algorithms use assembling to put together information from different small depth-bounded decision trees and make predictions with it.
- Both algorithms are nonlinear and can capture non-linear relationships between the features and the target variable due to the fact they are both using decision trees.

Differences:

- Training set: Random Forest creates different subsets of our training dataset using bootstrapping while in Gradient Boosted Trees we are using the entire dataset.
- Trees creation: Random Forest creates individual independent small decision trees based on each subset. Gradient Boosted Trees creates each tree serially where each tree is dependent on the trees that were created before him.
- Final prediction: Random Forests decides the final prediction value based on the prediction values decided by each independent tree (majority in classification or mean in regression) whereas in Gradient Boosted Trees the values predicted by each tree are weighted and the total sum of the weighted values is the prediction of the algorithm.
- Usage of features: Random Forests randomly choose at each iteration  $d'$  features and creating the decision tree based on them, whereas Gradient Boosted Trees takes into consideration all the features in each iteration.
- Variance and bias: Random Forests increases variance between the trees by random selection of features which in turn results in less variance in the predictions as we learned in the lecture, whereas Gradient Boosted Trees reduces bias by lowering the error at each iteration.

#### A23:

First we ran GridSearch CV on the 4 hyper parameters 'loss', 'learning\_rate', 'subsample', 'min\_samples\_leaf' and we got that the optimal values for them are:

*loss: hubert*

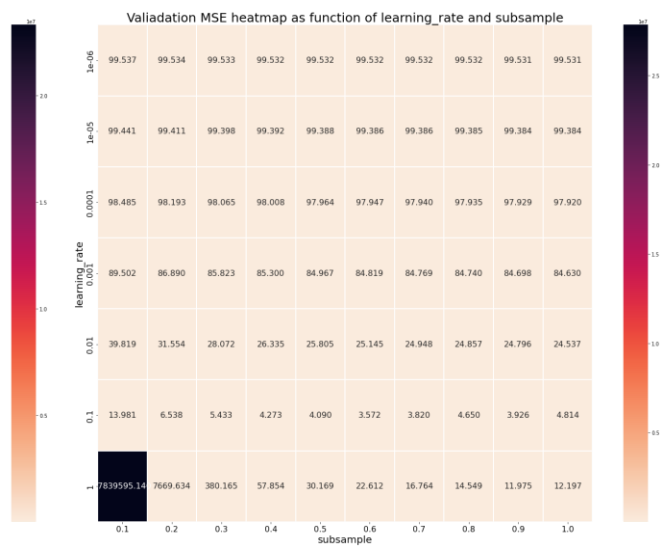
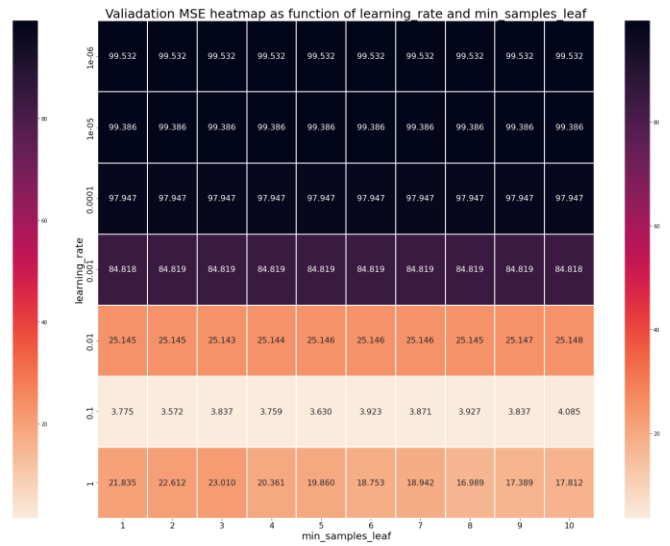
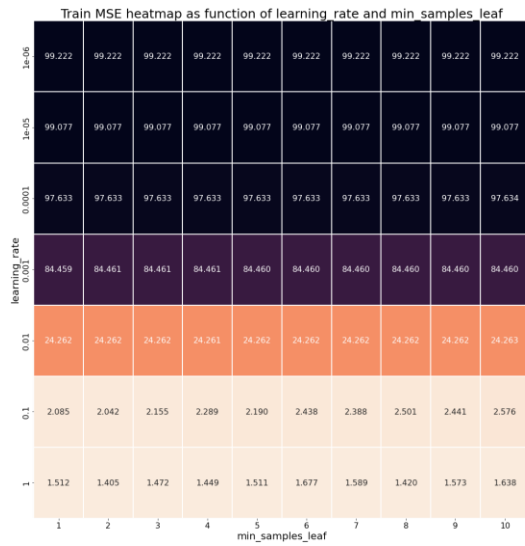
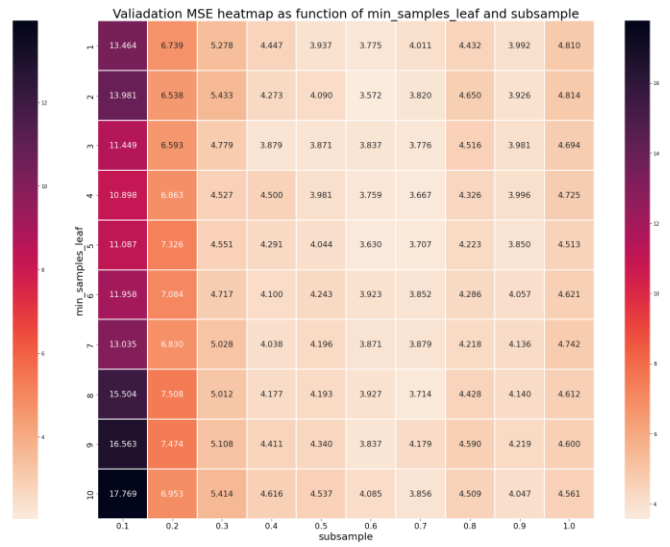
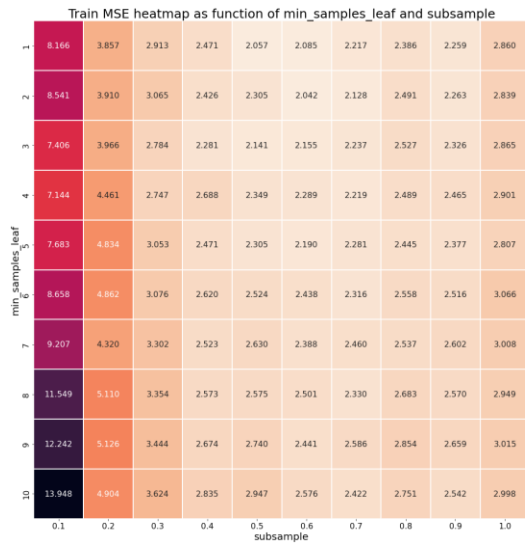
*learning\_rate: 0.1*

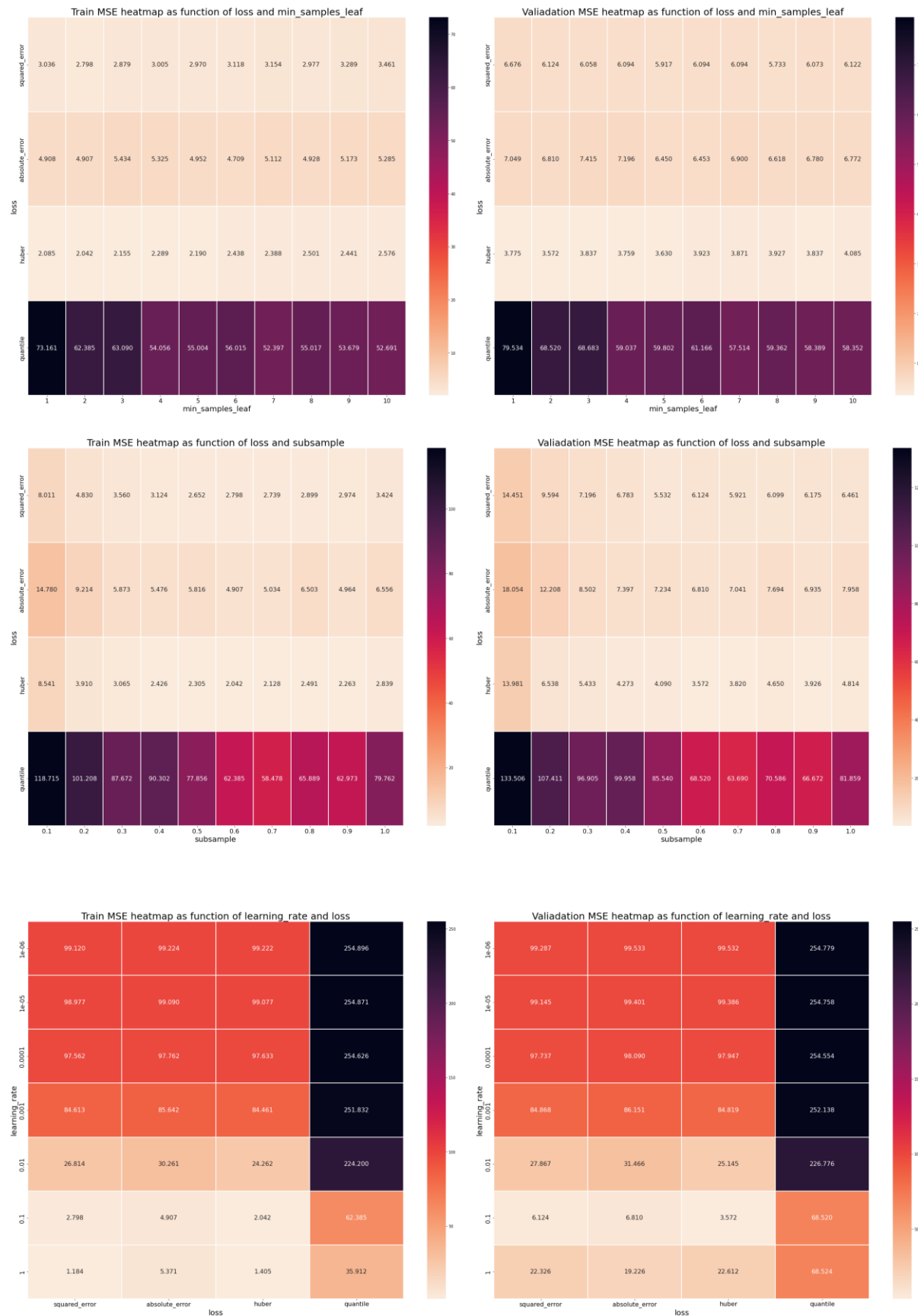
*subsample: 0.6*

*min\_sample\_leaf: 2*

With these optimal values the **validation MSE score is  $\cong 3.57$**  and the **train MSE score is  $\cong 2.04$** .

Now in order to get heatmap for each pair of parameters out of these four we set each time 2 parameters to their optimal value and get the results for the other two parameters.





As we can see from the 6 heat maps above, we get the optimal validation MSE score from the optimal parameters that we stated in the beginning of the question.

A24:

Model	Section	Train MSE	Valid MSE
		Cross validated	
Dummy	2	99.13589	99.3029
Linear	2	85.45293	91.10669
Ridge Linear	3	86.07993	89.7909
Polynomial Ridge Linear	4	76.9089	79.1425
GBM Regressor	5	2.0148	3.5724

## Section 6: Testing your models

A25:

Model	Section	Train MSE	Valid MSE	Test MSE
		Cross validated		Retrained
Dummy	2	99.13589	99.3029	108.1524
Linear	2	85.45293	91.10669	93.3292
Ridge Linear	3	86.07993	89.7909	93.9205
Polynomial Ridge Linear	4	76.9089	79.1425	83.9243
GBM Regressor	5	2.0148	3.5724	5.1886

The Gradient Boosting model performed the best on the training set, which was quite predictable since its train and validation MSE's were much smaller than those of the other models.

The target variable 'contamination\_level' does not behave linearly with all the features and therefore is not separable. All the other models (except the dummy and GBM) are trying to find a linear plain that will separate our data which is why they are not performing well. Whereas the gradient boosting model is nonlinear and more flexible due to the fact it is using decision trees, its train, validation, and test MSE are all much lower than the others thanks to those properties.

We can see that in all of the models, except gradient boosting, we have very high train, validation and test MSE, which can be seen as underfitting of the data caused by an unfit hypothesis class chosen to capture the connection between the features and the target variable as explained above.

All the models above have higher test than train MSE but we don't think the difference is high enough to conclude that they suffer from overfitting.

Additional insight is that in Q9 we were told to find the features with the top 5 largest coefficients and as we explained in A11 the magnitude of each coefficient can be seen as how correlated the corresponding feature with our target variable is. In section 5 we are using only a subset of the features and in that subset some of the features are the ones listed as the top 5 most correlated features with 'contamination\_level'. This might also be able to explain why the gradient boosting model performed so much better than all the other models that tried to capture the relationships of all the features instead of using only a subset of the most expressive and informative features.