Submitted by:
Kiran Kumar Valluru (APFE21709855)

# Part II - Assignment-based Subjective Questions

**Question 1**

What is the optimal value of alpha for ridge and lasso regression? What will be the changes in the model if you choose double the value of alpha for both ridge and lasso? What will be the most important predictor variables after the change is implemented?

Answer:

Optimal value of alpha for Ridge and Lasso Regression is below:
- Alpha value for Ridge Regression: 0.6
- Alpha value for Lasso Regression: 0.0001

With these alphas the R2 score for the models is approximately 0.87.

After doubling the alpha values in the Ridge and Lasso, the prediction accuracy remains around 0.87 with minor reduction in both Train and Test R2 score and there is a small change in the co-efficient values. The new model is created and demonstrated in the Jupyter notebook.

Comparison of metrics for Ridge and Lasso Regression with Original and Double the alpha:

| | Metrics | Ridge Regression_1 | Lasso Regression_0.0001 | Ridge Regression_2 | Lasso Regression_0.0002 |
|---|---|---|---|---|---|
| 0 | R2 Score (Train) | 0.914965 | 0.913294 | 0.914965 | 0.909298 |
| 1 | R2 Score (Test) | 0.877569 | 0.879497 | 0.877569 | 0.876695 |
| 2 | RSS (Train) | 1.490569 | 1.519848 | 1.490569 | 1.589900 |
| 3 | RSS (Test) | 0.871262 | 0.857543 | 0.871262 | 0.877479 |
| 4 | MSE (Train) | 0.001466 | 0.001494 | 0.001466 | 0.001563 |
| 5 | MSE (Test) | 0.001998 | 0.001967 | 0.001998 | 0.002013 |
| 6 | RMSE (Train) | 0.038284 | 0.038658 | 0.038284 | 0.039539 |
| 7 | RMSE (Test) | 0.044702 | 0.044349 | 0.044702 | 0.044862 |

**Ridge:** Train R2 score reduces little bit from 91.63% to 91.50% and Test R2 score also reduce from 87.96% to 87.76%.

**Lasso:** Train R2 score reduces little bit from 91.33% to 90.93% and Test R2 score also reduce from 87.95% to 87.67%.

Below are the changes in the co-efficients:

| Ridge and Lasso Coefficients (Original) | Ridge and Lasso Coefficients (after doubling the Alpha) |
|---|---|

| | Linear | Ridge | Lasso |
|---|---|---|---|
| MSSubClass | -0.028241 | -0.024568 | -0.023338 |
| LotArea | 0.068456 | 0.070607 | 0.060161 |
| OverallQual | 0.134807 | 0.141014 | 0.154474 |
| OverallCond | 0.082265 | 0.080584 | 0.077553 |
| BsmtQual | 0.062567 | 0.059743 | 0.052758 |
| BsmtExposure | 0.024953 | 0.024192 | 0.025897 |
| BsmtFinSF1 | 0.041753 | 0.042101 | 0.042138 |
| TotalBsmtSF | 0.053048 | 0.059649 | 0.051086 |
| HeatingQC | 0.023433 | 0.023705 | 0.023485 |
| GrLivArea | 0.341041 | 0.311403 | 0.330820 |
| BsmtFullBath | 0.022312 | 0.020905 | 0.020791 |
| KitchenQual | 0.022818 | 0.026886 | 0.026134 |
| Fireplaces | 0.031911 | 0.035227 | 0.032656 |
| GarageYrBlt | 0.037946 | 0.034666 | 0.030036 |
| GarageCars | 0.079308 | 0.080197 | 0.079090 |
| HouseAge | -0.021935 | -0.023168 | -0.024063 |
| MSZoning_FV | 0.111259 | 0.081689 | 0.047202 |
| MSZoning_RH | 0.104383 | 0.074993 | 0.035674 |
| MSZoning_RL | 0.115058 | 0.087753 | 0.055071 |
| MSZoning_RM | 0.094593 | 0.066375 | 0.030368 |
| Neighborhood_Crawfor | 0.039942 | 0.040467 | 0.035943 |

| | Ridge_2 | Lasso_0002 |
|---|---|---|
| MSSubClass | -0.022128 | -0.021593 |
| LotArea | 0.070564 | 0.046789 |
| OverallQual | 0.143594 | 0.165880 |
| OverallCond | 0.078273 | 0.070143 |
| BsmtQual | 0.057629 | 0.042628 |
| BsmtExposure | 0.023726 | 0.025479 |
| BsmtFinSF1 | 0.042498 | 0.044130 |
| TotalBsmtSF | 0.063503 | 0.046764 |
| HeatingQC | 0.023912 | 0.023261 |
| GrLivArea | 0.289025 | 0.327614 |
| BsmtFullBath | 0.019772 | 0.019735 |
| KitchenQual | 0.030279 | 0.028143 |
| Fireplaces | 0.038009 | 0.033220 |
| GarageYrBlt | 0.032507 | 0.026290 |
| GarageCars | 0.081073 | 0.080845 |
| HouseAge | -0.024210 | -0.026039 |
| MSZoning_FV | 0.065488 | 0.015205 |
| MSZoning_RH | 0.058659 | 0.000000 |
| MSZoning_RL | 0.072246 | 0.027297 |
| MSZoning_RM | 0.050071 | 0.000000 |
| Neighborhood_Crawfor | 0.040846 | 0.032368 |

**Question 2**

You have determined the optimal value of lambda for ridge and lasso regression during the assignment. Now, which one will you choose to apply and why?

Answer:

The optimum lambda value in case of Ridge and Lasso is as follows:-
- Ridge – 0.6
- Lasso – 0.0001

The Mean Squared Error in case of Ridge and Lasso are:
- Ridge - 0.001965
- Lasso - 0.001967

The Mean Squared Error of both the models are almost same. R2 Score of Lasso is better than Ridge for Test Data.

Since Lasso helps in feature reduction (as the coefficient value of some of the features become zero), **Lasso** has a better edge over Ridge and should be used as the final model.

**Question 3**

After building the model, you realised that the five most important predictor variables in the lasso model are not available in the incoming data. You will now have to create another model excluding the five most important predictor variables. Which are the five most important predictor variables now?

Answer:

The five most important predictor variables in the current lasso model is:-
1. GrLivArea: Above grade (ground) living area square feet
2. OverallQual: Rates the overall material and finish of the house
3. GarageCars: Size of garage in car capacity
4. OverallCond: Rates the overall condition of the house
5. LotArea: Lot size in square feet

Have build a Lasso model in the Jupyter notebook after removing these attributes from the dataset.

Train R2 score reduces drastically from 91.33% to 82.19% and Test R2 score also reduces from 87.95% to 73.65% after removal of Top 5 predictive variables.

The Mean Squared Error increases to 0.004300652262493751

The new Top 5 predictos after removing the Top 5 predictive variables are:-

1. TotalBsmtSF: Total square feet of basement area
2. FirePlaces: Number of fireplaces

3. KitchenQual: Kitchen quality
4. Neighborhood_NoRidge: Physical locations within Ames city limits - Northridge
5. BsmtQual: Evaluates the height of the basement

| | Features | Coeficient_drop |
|---|---|---|
| **4** | TotalBsmtSF | 0.190762 |
| **8** | Fireplaces | 0.103577 |
| **7** | KitchenQual | 0.101346 |
| **18** | Neighborhood_NoRidge | 0.099565 |
| **1** | BsmtQual | 0.090147 |

**Question 4**

How can you make sure that a model is robust and generalisable? What are the implications of the same for the accuracy of the model and why?

Answer:

As Per, Occam's Razor— given two models that show similar 'performance' in the finite training or test data, we should pick the one that makes fewer on the test data due to following reasons:

▪ Simpler models are usually more 'generic' and are more widely applicable.
▪ Simpler models require fewer training samples for effective training than the more complex ones and hence are easier to train.
▪ Simpler models are more robust.
    ○ Complex models tend to change wildly with changes in the training data set
    ○ Simple models have low variance, high bias and complex models have low bias, high variance
    ○ Simpler models make more errors in the training set. Complex models lead to overfitting - they work very well for the training samples, fail miserably when applied to other test samples.

Therefore, to make the model more robust and generalizable, make the model simple but not simpler which will not be of any use.

Regularization can be used to make the model simpler. Regularization helps to strike the delicate balance between keeping the model simple and not making it too naive to be of any use. For regression, regularization involves adding a regularization term to the cost that adds up the absolute values or the squares of the parameters of the model.
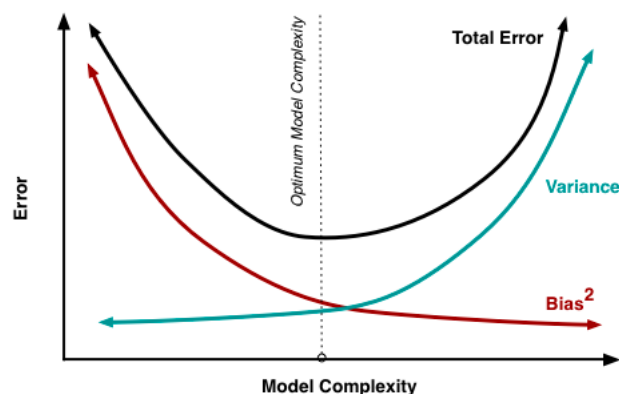
Also, Making a model simple leads to Bias-Variance Trade-off:
- A complex model will need to change for every little change in the dataset and hence is very unstable and extremely sensitive to any changes in the training data.
- A simpler model that abstracts out some pattern followed by the data points given is unlikely to change wildly even if more points are added or removed.

Bias quantifies how accurate is the model likely to be on test data. A complex model can do an accurate job prediction provided there is enough training data. Models that are too naïve, for e.g., one that gives same answer to all test inputs and makes no discrimination whatsoever has a very large bias as its expected error across all test inputs are very high.

Variance refers to the degree of changes in the model itself with respect to changes in the training data.

Thus, <u>accuracy of the model can be maintained by keeping the balance between Bias and Variance</u> as it minimizes the total error as shown in the below graph.



Below are the pointers learned during the course module as key items:
- Robust refers the model works for a broad range of inputs. If the model gets really good results at training time (it seems "more accurate") but won't generalize to out-of-sample data (i.e., it isn't robust) then we call it overfitting.
- The model should be generalized so that the test accuracy is not lesser than the training score.
- Here in our case, based on all data and modelling both Ridge and Lasso performed good on Train and Test Data which shows our model with Alpha value "0.6" for Ridge and "0.0001" for Lasso is Robust and more Generalized model.
- Too much importance should not be given to the outliers so that the accuracy predicted by the model is high. But outlier's analysis needs to be done and only those which are relevant to the dataset need to be retained and rest should be dropped.