# Assignment - Advanced Regression

**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?**

From Lasso regression or L1 regularization the optimal value of alpha is found using best_params_ as (Alpha= 0.001) for which the training accuracy is 0.90 and testing accuracy is 0.78.
From Ridge regression or L2 regularization the optimal value of alpha is found using best_params_ as (Alpha= 100) for which the training accuracy is 0.90 and testing accuracy is 0.79.

When we double the value of alpha the neg_mean_absolute error will become more negative and the training accuracy will further drop down.

For Lasso regression on doubling the value of Alpha in model we see a drop in training and testing accuracy to 0.89 and 0.76 respectively.

And the most important predictor variables after the change are **MSZoning, HouseStyle, YearRemodAdd, GarageQual and Electrical.**

For Ridge regression on doubling the value of Alpha in model we see drop in training and testing accuracy to 0.88 and 0.76 respectively.

The most important predictor variables after the change are **MSZoning, YearRemodAdd, HouseStyle, Fireplaces and GarageQual.**


**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?**

Out of L1(Lasso) and L2(Ridge) regression model ran in the assignment, we are getting the values of test accuracy as nearly same with value for Ridge regression model higher than Lasso regression.
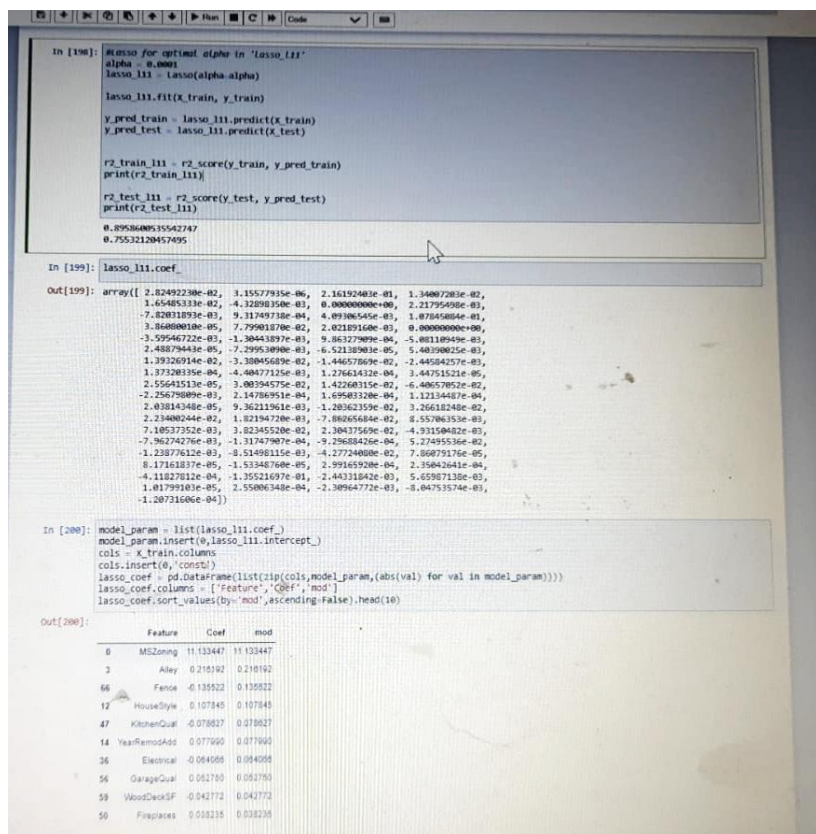
However if we look at the most important predictor variables they are almost same for both. We can see that lasso regression eliminates a lot of feature that are related to each other and ridge regression does not, it adds bias to reach optimal model. Hence we can choose Ridge regression as it is performing better on the unseen test data and comparatively more robust.

## 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?**

After building the model if most important top 5 predictors in the lasso model are not available in incoming data then we need create another model without these predictors which are **MSubClass, BldgType, OverallCond, YearBuilt and GarageArea.**

After recreating the lasso model as done in the python code below the most important predictors now are **MSZoning, Alley, Fence, HouseStyle and KitchenQual.**



## 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?**
A robust and generalizable machine learning model tends to be better working on unseen or test data.
One way is that we can simply split your original data set into 3: training, validation and test; use training and validation for the model development and keep test data unseen for a final check after cross validation. This will check your model's generalizability.

We can also make our machine learning model more robust and generalizable by using the following options:

**Regularization:** With the help of regularization we can reduce over-fitting and variance. This will give us a model which will work well on unseen data using the optimal values of hyperparameter.

Also we can choose a more practical error metric for eg Mean absolute error or Huber loss.

**Data preparation:** When the data is heavily skewed for eg target variable in our House prediction dataset is heavily positive skewed we can transform the data into Log values to get it near to normal distribution as machine learning models tend to work much better with normally distributed target variables.

In case there are extreme values present in the dataset we can cap the data at various thresholds from top and bottom percentiles. For eg we can cap the bottom and top 5% values at 5th and 95th percentiles. This is called Winsorization.

Removing/ handling outlier values and Winsorize basically means removing information from dataset and are generally not preferred.

**Choosing different models:** Model selection is the best way to look for a machine learning model that is best suitable to dataset in our problem and performs better with unseen data.