Assignment-based Subjective Questions

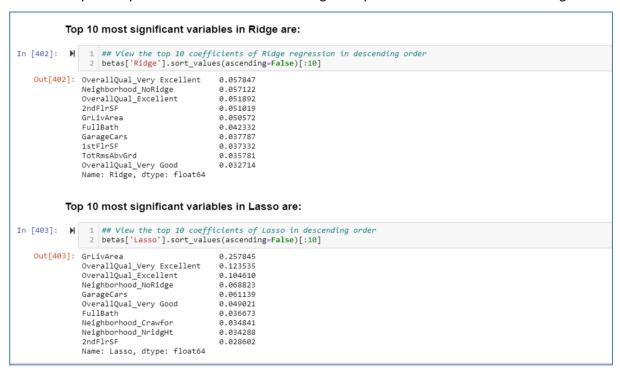
1. What is the optimal value of alpha for ridge and lasso regression? What will be the changes in the model if you choose double the value of alpha for both ridge and lasso? What will be the most important predictor variables after the change is implemented?

Optimal Value of alpha for ridge and lasso regression are:

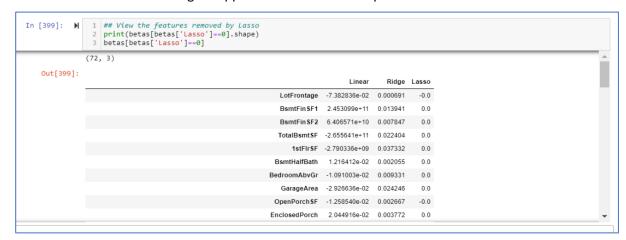
- Optimal Value of lambda for Ridge: 6
- Optimal Value of lambda for Lasso: 0.0001

If we choose to double the value of alpha for both ridge and lasso: in case of Ridge that will lower the coefficients and in case of Lasso, there would be more less important features coefficients turning 0.

The most important predictor variable after the change is implemented are those which are significant.



The number of zero coefficients increased from earlier 51 to 72. Also, there is a reduction in the high value of coefficients after Ridge is applied with double the alpha value.



Below are the snapshots from Jupyter Notebook, capturing key information when the lambda for Ridge and Lasso are changed to double their initial value.

```
#Fitting Ridge model for optimal value of alpha and printing coefficients which have been penalised

alpha = 12

# Create a ridge regreesion instance with optimum value of alpha
ridge = Ridge(alpha=alpha)

# Fit the model on training data
ridge.fit(X_train, y_train)

# View the coefficients of ridge regression fitted model
print(ridge.coef_)
```

```
In [380]: | 1 ridge.score(X_train,y_train)

Out[380]: 0.8698606374816347

In [381]: | 1 ridge.score(X_test,y_test)

Out[381]: 0.8569085311544
```

```
R-Squared (Train) = 0.87
R-Squared (Test) = 0.86
RSS (Train) = 1.60
RSS (Test) = 0.78
MSE (Train) = 0.00
MSE (Test) = 0.00
```

```
[('constant', 0.18),
 ('LotFrontage', 0.001),
 ('LotArea', 0.019),
 ('MasVnrArea', 0.025),
 ('BsmtFinSF1', 0.014),
 ('BsmtFinSF2', 0.008),
 ('BsmtUnfSF', 0.02),
 ('TotalBsmtSF', 0.022),
 ('1stFlrSF', 0.037),
 ('2ndFlrSF', 0.051),
 ('LowQualFinSF', -0.007),
 ('GrLivArea', 0.051),
 ('BsmtFullBath', 0.023),
 ('BsmtHalfBath', 0.002),
 ('FullBath', 0.042),
 ('HalfBath', 0.017),
 ('BedroomAbvGr', 0.009),
 ('KitchenAbvGr', -0.012),
 ('TotRmsAbvGrd', 0.036),
 ('Fireplaces', 0.02),
 ('GarageCars', 0.038),
 ('GarageArea', 0.024),
 ('WoodDeckSF', 0.017),
 ('OpenPorchSF', 0.003),
 ('EnclosedPorch', 0.004),
 ('3SsnPorch', 0.008),
 ('ScreenPorch', 0.008),
 ('PoolArea', 0.003),
 ('MiscVal', -0.001),
 ('YearSinceRemodelAtSale', -0.017),
```

```
('AgeAtSale', -0.014),
 ('MSSubClass 1-1/2 STORY FINISHED ALL AGES', 0.006),
 ('MSSubClass 1-STORY 1945 & OLDER', -0.006),
 ('MSSubClass 1-STORY 1946 & NEWER ALL STYLES', 0.012),
 ('MSSubClass 1-STORY PUD (Planned Unit Development) - 1946 & NEWER', -0.0
1),
 ('MSSubClass 1-STORY W/FINISHED ATTIC ALL AGES', 0.003),
 ('MSSubClass 2 FAMILY CONVERSION - ALL STYLES AND AGES', -0.005),
 ('MSSubClass 2-1/2 STORY ALL AGES', 0.004),
 ('MSSubClass 2-STORY 1945 & OLDER', 0.005),
 ('MSSubClass 2-STORY 1946 & NEWER', 0.005),
 ('MSSubClass 2-STORY PUD - 1946 & NEWER', -0.009),
 ('MSSubClass DUPLEX - ALL STYLES AND AGES', -0.004),
 ('MSSubClass PUD - MULTILEVEL - INCL SPLIT LEV/FOYER', -0.005),
 ('MSSubClass SPLIT FOYER', -0.003),
 ('MSSubClass SPLIT OR MULTI-LEVEL', -0.001),
 ('MSZoning Others', -0.008),
 ('MSZoning RL', 0.003),
 ('MSZoning RM', -0.005),
 ('LotShape_IR2', 0.005),
 ('LotShape IR3', -0.012),
 ('LotShape Reg', -0.0),
 ('LotConfig CulDSac', 0.013),
 ('LotConfig FR2', -0.015),
 ('LotConfig FR3', -0.004),
 ('LotConfig Inside', -0.002),
 ('Neighborhood Blueste', -0.001),
 ('Neighborhood BrDale', -0.001),
 ('Neighborhood BrkSide', -0.0),
 ('Neighborhood_ClearCr', 0.006),
 ('Neighborhood CollgCr', -0.007),
 ('Neighborhood Crawfor', 0.026),
 ('Neighborhood Edwards', -0.026),
 ('Neighborhood Gilbert', -0.012),
 ('Neighborhood IDOTRR', -0.012),
 ('Neighborhood_MeadowV', -0.01),
 ('Neighborhood Mitchel', -0.013),
 ('Neighborhood NAmes', -0.012),
 ('Neighborhood NPkVill', -0.002),
 ('Neighborhood NWAmes', -0.006),
 ('Neighborhood NoRidge', 0.057),
 ('Neighborhood_NridgHt', 0.031),
 ('Neighborhood OldTown', -0.015),
 ('Neighborhood SWISU', -0.003),
 ('Neighborhood Sawyer', -0.015),
 ('Neighborhood SawyerW', -0.003),
 ('Neighborhood_Somerst', 0.015),
 ('Neighborhood_StoneBr', 0.017),
 ('Neighborhood Timber', -0.004),
 ('Neighborhood Veenker', 0.01),
 ('BldgType 2fmCon', -0.005),
 ('BldgType Duplex', -0.004),
 ('BldgType Twnhs', -0.012),
 ('BldgType TwnhsE', -0.01),
 ('HouseStyle 1Story', 0.004),
 ('HouseStyle 2Story', -0.005),
```

```
('HouseStyle Others', -0.002),
('HouseStyle SLvl', -0.005),
('OverallQual Average', -0.008),
('OverallQual Below Average', -0.009),
('OverallQual Excellent', 0.052),
('OverallQual Fair', -0.013),
('OverallQual Good', 0.01),
('OverallQual Poor', -0.007),
('OverallQual_Very Excellent', 0.058),
('OverallQual Very Good', 0.033),
('OverallQual Very Poor', -0.005),
('OverallCond Average', -0.009),
('OverallCond Below Average', -0.011),
('OverallCond Excellent', 0.016),
('OverallCond Fair', -0.017),
('OverallCond Good', 0.005),
('OverallCond Poor', -0.001),
('OverallCond Very Good', -0.0),
('OverallCond Very Poor', -0.005),
('RoofStyle Hip', 0.005),
('RoofStyle Others', 0.003),
('Exterior1st CemntBd', -0.003),
('Exterior1st HdBoard', -0.011),
('Exterior1st MetalSd', -0.007),
('Exterior1st_Others', -0.014),
('Exterior1st_Plywood', -0.004),
('Exterior1st VinylSd', -0.006),
('Exterior1st Wd Sdng', -0.003),
('Exterior2nd CmentBd', 0.001),
('Exterior2nd HdBoard', 0.004),
('Exterior2nd MetalSd', 0.001),
('Exterior2nd Others', 0.011),
('Exterior2nd Plywood', -0.004),
('Exterior2nd VinylSd', 0.002),
('Exterior2nd Wd Sdng', -0.002),
('Exterior2nd Wd Shng', -0.008),
('MasVnrType BrkFace', 0.004),
('MasVnrType None', 0.006),
('MasVnrType Stone', -0.0),
('ExterQual \overline{F}a', -0.017),
('ExterQual Gd', -0.004),
('ExterQual_TA', -0.016),
('Foundation CBlock', 0.003),
('Foundation Others', 0.003),
('Foundation PConc', 0.005),
('BsmtQual Fa', -0.019),
('BsmtQual_Gd', -0.033),
('BsmtQual No Basement', -0.015),
('BsmtQual TA', -0.028),
('BsmtExposure Gd', 0.029),
('BsmtExposure Mn', -0.004),
('BsmtExposure No', -0.013),
('BsmtExposure No Basement', -0.016),
('BsmtFinType1_BLQ', 0.001),
('BsmtFinType1 GLQ', 0.013),
('BsmtFinType1 LwQ', -0.004),
```

```
('BsmtFinType1 No Basement', -0.015),
 ('BsmtFinType1 Rec', -0.001),
('BsmtFinType1 Unf', -0.01),
 ('HeatingQC Fa', -0.005),
 ('HeatingQC_Gd', -0.004),
 ('HeatingQC Po', -0.002),
 ('HeatingQC TA', -0.004),
 ('KitchenQual Fa', -0.024),
 ('KitchenQual_Gd', -0.027),
 ('KitchenQual TA', -0.029),
 ('FireplaceQu None', -0.009),
 ('FireplaceQu Others', 0.005),
('FireplaceQu TA', -0.001),
('GarageType BuiltIn', -0.001),
('GarageType Detchd', -0.007),
 ('GarageType No Garage', -0.004),
 ('GarageType Others', -0.01),
 ('GarageFinish No Garage', -0.004),
 ('GarageFinish RFn', -0.008),
 ('GarageFinish_Unf', -0.011),
 ('Fence None', -0.001),
 ('Fence Others', -0.007),
 ('SaleCondition Normal', 0.002),
 ('SaleCondition Others', 0.002),
('SaleCondition Partial', 0.011)]
   In [388]: N 1 #optimum alpha
             3 alpha = 0.0002
                Create a lasso regreesion instance with optimum value of alpha
             5 lasso = Lasso(alpha=alpha)
             7 # Fit the model on training data
            8 lasso.fit(X_train, y_train)
            Lasso(alpha=0.0002)
In [390]: | 1 lasso.score(X_train,y_train)
   Out[390]: 0.8767438395130009
In [391]: | 1 lasso.score(X_test,y_test)
   Out[391]: 0.8594790000453509
  R-Squared (Train) = 0.88
  R-Squared (Test) = 0.86
  RSS (Train) = 1.52
  RSS (Test) = 0.76
  MSE (Train) = 0.00
  MSE (Test) = 0.00
[('constant', 0.127),
('LotFrontage', -0.0),
 ('LotArea', 0.006),
 ('MasVnrArea', 0.004),
('BsmtFinSF1', 0.0),
 ('BsmtFinSF2', 0.0),
 ('BsmtUnfSF', 0.002),
```

```
('TotalBsmtSF', 0.0),
('1stFlrSF', 0.0),
('2ndFlrSF', 0.029),
('LowQualFinSF', -0.008),
('GrLivArea', 0.258),
('BsmtFullBath', 0.024),
('BsmtHalfBath', 0.0),
('FullBath', 0.037),
('HalfBath', 0.012),
('BedroomAbvGr', 0.0),
('KitchenAbvGr', -0.024),
('TotRmsAbvGrd', 0.001),
('Fireplaces', 0.016),
('GarageCars', 0.061),
('GarageArea', 0.0),
('WoodDeckSF', 0.007),
('OpenPorchSF', -0.0),
('EnclosedPorch', 0.0),
('3SsnPorch', 0.0),
('ScreenPorch', 0.005),
('PoolArea', -0.0),
('MiscVal', -0.0),
('YearSinceRemodelAtSale', -0.013),
('AgeAtSale', -0.023),
('MSSubClass 1-1/2 STORY FINISHED ALL AGES', 0.003),
('MSSubClass 1-STORY 1945 & OLDER', -0.0),
('MSSubClass 1-STORY 1946 & NEWER ALL STYLES', 0.019),
('MSSubClass 1-STORY PUD (Planned Unit Development) - 1946 & NEWER', -0.0
('MSSubClass 1-STORY W/FINISHED ATTIC ALL AGES', 0.0),
('MSSubClass_2 FAMILY CONVERSION - ALL STYLES AND AGES', -0.0),
('MSSubClass 2-1/2 STORY ALL AGES', -0.0),
('MSSubClass 2-STORY 1945 & OLDER', 0.0),
('MSSubClass 2-STORY 1946 & NEWER', 0.0),
('MSSubClass 2-STORY PUD - 1946 & NEWER', -0.007),
('MSSubClass DUPLEX - ALL STYLES AND AGES', -0.0),
('MSSubClass PUD - MULTILEVEL - INCL SPLIT LEV/FOYER', -0.0),
('MSSubClass SPLIT FOYER', 0.0),
('MSSubClass SPLIT OR MULTI-LEVEL', 0.0),
('MSZoning Others', -0.002),
('MSZoning RL', 0.007),
('MSZoning_RM', -0.0),
('LotShape IR2', 0.003),
('LotShape_IR3', -0.005),
('LotShape Reg', 0.0),
('LotConfig CulDSac', 0.013),
('LotConfig_FR2', -0.007),
('LotConfig FR3', -0.0),
('LotConfig Inside', -0.0),
('Neighborhood Blueste', -0.0),
('Neighborhood BrDale', 0.0),
('Neighborhood BrkSide', 0.0),
('Neighborhood ClearCr', 0.011),
('Neighborhood CollgCr', 0.0),
('Neighborhood Crawfor', 0.035),
('Neighborhood Edwards', -0.024),
```

```
('Neighborhood Gilbert', -0.0),
('Neighborhood IDOTRR', -0.004),
('Neighborhood MeadowV', -0.0),
('Neighborhood Mitchel', -0.004),
('Neighborhood NAmes', -0.005),
('Neighborhood NPkVill', -0.0),
('Neighborhood NWAmes', -0.0),
('Neighborhood NoRidge', 0.069),
('Neighborhood NridgHt', 0.034),
('Neighborhood OldTown', -0.014),
('Neighborhood SWISU', -0.0),
('Neighborhood Sawyer', -0.01),
('Neighborhood SawyerW', 0.0),
('Neighborhood Somerst', 0.023),
('Neighborhood_StoneBr', 0.008),
('Neighborhood Timber', 0.0),
('Neighborhood Veenker', 0.005),
('BldgType 2fmCon', -0.0),
('BldgType Duplex', -0.0),
('BldgType_Twnhs', -0.011),
('BldgType TwnhsE', -0.005),
('HouseStyle 1Story', 0.002),
('HouseStyle 2Story', -0.005),
('HouseStyle Others', -0.0),
('HouseStyle SLvl', -0.0),
('OverallQual Average', -0.003),
('OverallQual Below Average', -0.004),
('OverallQual Excellent', 0.105),
('OverallQual Fair', -0.001),
('OverallQual Good', 0.017),
('OverallQual Poor', -0.0),
('OverallQual Very Excellent', 0.124),
('OverallQual Very Good', 0.049),
('OverallQual Very Poor', -0.0),
('OverallCond Average', -0.011),
('OverallCond Below Average', -0.012),
('OverallCond Excellent', 0.009),
('OverallCond Fair', -0.017),
('OverallCond Good', 0.004),
('OverallCond Poor', -0.0),
('OverallCond Very Good', 0.001),
('OverallCond_Very Poor', -0.0),
('RoofStyle Hip', 0.0),
('RoofStyle Others', 0.0),
('Exterior1st CemntBd', -0.0),
('Exterior1st HdBoard', -0.0),
('Exterior1st_MetalSd', -0.0),
('Exterior1st Others', -0.011),
('Exterior1st Plywood', -0.0),
('Exterior1st VinylSd', 0.0),
('Exterior1st Wd Sdng', -0.0),
('Exterior2nd CmentBd', 0.0),
('Exterior2nd HdBoard', 0.0),
('Exterior2nd MetalSd', -0.0),
('Exterior2nd Others', 0.008),
('Exterior2nd Plywood', -0.0),
```

```
('Exterior2nd VinylSd', 0.004),
('Exterior2nd Wd Sdng', -0.0),
('Exterior2nd Wd Shng', -0.002),
('MasVnrType BrkFace', 0.0),
('MasVnrType None', 0.0),
('MasVnrType Stone', -0.002),
('ExterQual Fa', -0.004),
('ExterQual Gd', 0.006),
('ExterQual_TA', -0.0),
('Foundation CBlock', 0.0),
('Foundation Others', -0.0),
('Foundation PConc', 0.0),
('BsmtQual Fa', -0.014),
('BsmtQual Gd', -0.026),
('BsmtQual No Basement', -0.029),
('BsmtQual TA', -0.021),
('BsmtExposure Gd', 0.027),
('BsmtExposure Mn', -0.0),
('BsmtExposure No', -0.01),
('BsmtExposure No Basement', -0.016),
('BsmtFinType1 BLQ', 0.0),
('BsmtFinType1 GLQ', 0.009),
('BsmtFinType1 LwQ', -0.001),
('BsmtFinType1 No Basement', -0.0),
('BsmtFinType1 Rec', -0.0),
('BsmtFinType1 Unf', -0.008),
('HeatingQC_Fa', -0.0),
('HeatingQC Gd', -0.001),
('HeatingQC Po', -0.0),
('HeatingQC TA', -0.001),
('KitchenQual_Fa', -0.025),
('KitchenQual Gd', -0.021),
('KitchenQual TA', -0.025),
('FireplaceQu None', -0.006),
('FireplaceQu Others', 0.0),
('FireplaceQu TA', -0.0),
('GarageType BuiltIn', -0.002),
('GarageType Detchd', -0.004),
('GarageType No Garage', -0.0),
('GarageType Others', -0.007),
('GarageFinish No Garage', -0.0),
('GarageFinish_RFn', -0.003),
('GarageFinish Unf', -0.008),
('Fence None', -0.0),
('Fence Others', -0.003),
('SaleCondition Normal', 0.003),
('SaleCondition_Others', 0.0),
('SaleCondition Partial', 0.013)]
```

Out[394]:

	Metric	Linear Regression	Ridge Regression	Lasso Regression
0	R2 Score (Train)	8.935476e-01	0.869861	0.876744
1	R2 Score (Test)	-9.550277e+19	0.856909	0.859479
2	RSS (Train)	1.309900e+00	1.601369	1.516671
3	RSS (Test)	5.191360e+20	0.777820	0.763847
4	MSE (Train)	3.581840e-02	0.039603	0.038542
5	MSE (Test)	1.088688e+09	0.042141	0.041761

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?

Optimal Value of alpha for ridge and lasso regression are:

- Optimal Value of lambda for ridge: 6
- Optimal Value of lambda for Lasso: 10

As we got good score for both the models so we can go with Lasso regression as it results in model parameters such that lesser important features coefficients become zero.

Also, the model we choose to apply will depend on the use case. If we have too many variables and one of our primary goals is feature selection, then we will use Lasso. If we don't want to get too large coefficients and reduction of coefficient magnitude is one of our prime goals, then we will use Ridge Regression.

Ridge: Train: 88.0 Test: 86.5 and Lasso: Train: 86.5 Test: 86.7

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?

Below is the list of predictor variables after the first execution of Lasso model.

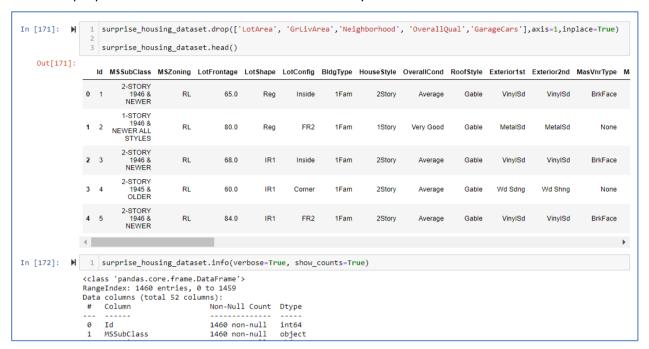
Top 6 parameters – Lasso	
GrLivArea	0.250554
OverallQual_Very Excellent	0.118687
OverallQual_Excellent	0.099083
Neighborhood_NoRidge	0.071331
LotArea	0.064057
GarageCars	0.059015

Assuming that completely new set of 5 predictor variables is now expected, we remove OverallQual and Neighborhood altogether rather than only removing the specific dummy variable as listed above. Also, the above list contains 6 variables instead of 5 as we consider the two categories of OverquallQual as one single variable.

On running the same notebook and removing the top 5 significant variables, we found below variables as next 5 significant.

Below are the snapshots from Jupyter Notebook, capturing key information when the top 5 predictor variables obtained from Lasso regression are removed from the dataset and Lasso regression is again performed.

Deleted top 5 predictor variables before creation of dummy variables.



Latest set of variables.

Optimum alpha remains same.

```
# Printing the best hyperparameter alpha
print(model_cv.best_params_)
In [207]: H
                {'alpha': 0.0001}
In [208]: ▶
                 1 #optimum alpha
                     alpha = 0.0001
                    # Create a lasso regreesion instance with optimum value of alpha lasso = Lasso(alpha=alpha)
                     # Fit the model on training data
                  8 lasso.fit(X_train, y_train)
   Out[208]:
                Lasso(alpha=0.0001)
                 1 # View the coefficients of lasso fitted model
                 2 lasso.coef
    Out[209]: array([-6.68193272e-03, 7.85489070e-02, 0.00000000e+00, 7.35836887e-03,
                        2.50475467e-02, 0.00000000e+00, 2.47629817e-01, -0.00000000e+00, 2.37966926e-02, 0.00000000e+00,
                                                                                    1.73814357e-01.
                                                                0.00000000e+00,
                                                                                    4.16342480e-02,
                         1.04122321e-02, -0.00000000e+00, -5.50291601e-02,
                                                                                   7.25458338e-03,
                         3.16165102e-02, 7.44516404e-02, 1.32349528e-02, -4.89352646e-03,
                         0.00000000e+00, 0.0000000e+00,
                                                                3.21912156e-03, -0.00000000e+00,
                        -0.00000000e+00, -7.80576027e-03, -4.45578194e-02, 0.00000000e+00,
                        -0.00000000e+00, 9.23039847e-03, -0.00000000e+00, 0.00000000e+00, -0.00000000e+00, 0.00000000e+00, 1.66087289e-02, 0.00000000e+00,
                         -0.00000000e+00, -1.68764757e-03, -0.00000000e+00, -2.35553316e-03
```

Out[214]:					
		Metric	Linear Regression	Ridge Regression	Lasso Regression
	0	R2 Score (Train)	8.495093e-01	0.840442	0.841459
	1	R2 Score (Test)	-2.631968e+18	0.840677	0.842775
	2	RSS (Train)	1.851792e+00	1.963365	1.950856
	3	RSS (Test)	1.430691e+19	0.866051	0.854649
	4	MSE (Train)	4.258761e-02	0.043852	0.043712
	5	MSE (Test)	1.807323e+08	0.044467	0.044173

Number of zero coefficients changed to 43.

```
In [219]:
                 1 ## View the features removed by Lasso
                    print(betas[betas['Lasso']==0].shape)
betas[betas['Lasso']==0]
                (43, 3)
    Out[219]:
                                                                                               Ridge Lasso
                                                                                     Linear
                                                                   BsmtFin SF1 -4.808405e+11 0.024809
                                                                                                       0.0
                                                                   TotalBsmtSF 5.205414e+11 0.042899
                                                                 LowQualFinSF 2.943180e-03 -0.008406 -0.0
                                                                  BsmtHalfBath 3.595352e-04 -0.000859
                                                                BedroomAbvGr -2.507305e-02 -0.002695 -0.0
                                                                EnclosedPorch 1.287866e-02 0.008370
                                                                   3SsnPorch 2.249336e-02 0.013900 0.0
                                                                     PoolArea -2.987289e-02 -0.001747
                                                                      MiscVal -8.257151e-03 -0.005663 -0.0
                                     MSSubClass_1-1/2 STORY FINISHED ALL AGES -4.489708e-02 -0.000978 0.0
                                             MSSubClass_1-STORY 1945 & OLDER -4.700518e-02 -0.009014 -0.0
                MSSubClass_1-STORY PUD (Planned Unit Development) - 1946 & NEWER -6.703186e-02 -0.011601 -0.0
                                MSSubClass_1-STORY W/FINISHED ATTIC ALL AGES -5.088830e-02 0.003146 0.0
                        MSSubClass_2 FAMILY CONVERSION - ALL STYLES AND AGES -3.782447e+09 -0.006660
                                              MSSubClass_2-1/2 STORY ALL AGES -4.025841e-02 -0.001830 0.0
                                             MSSubClass_2-STORY 1946 & NEWER -5.986023e-02 -0.003034
                                        MSSubClass_2-STORY PUD - 1946 & NEWER -7.008362e-02 -0.008500 -0.0
                            MSSubClass_PUD - MULTILEVEL - INCL SPLIT LEV/FOYER -7.228947e-02 -0.015231 -0.0
                                              MSSubClass SDLIT OR MILITLI EVEL _4 5210420_02 _0 006500
```

```
Top 10 most significant variables in Ridge are:
In [221]: N
              1 ## View the top 10 coefficients of Ridge regression in descending order
              2 betas['Ridge'].sort_values(ascending=False)[:10]
   Out[221]: 2ndFlrSF
                              0.117354
             1stFlrSF
                              0.097380
             GarageArea
                              0.072756
             MasVnrArea
                              0.072739
             FullBath
                              0.056811
             TotRmsAbvGrd
                              0.050789
             TotalBsmtSF
                              0.042899
             BsmtUnfSF
                              0.040652
             BsmtExposure_Gd 0.032728
             Fireplaces
                              0.031864
             Name: Ridge, dtype: float64
         Top 5 most significant variables in Lasso are:
In [223]: N 1 ## View the top 10 coefficients of Lasso in descending order
              2 betas['Lasso'].sort_values(ascending=False)[:5]
   Out[223]: 1stFlrSF
                          0.247630
             2ndFlrSF
                          0.173814
             MasVnrArea 0.078549
             GarageArea 0.074452
             FullBath
                          0.041634
             Name: Lasso, dtype: float64
```

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 model is robust when any variation in the predictor variables data does not affect its performance considerably.
- A generalized model can adapt properly to new, previously unseen data, drawn from the same distribution as the one used to create the model. This is the test data out of the entire dataset.
- To ensure that a model is robust and generalizable, we must ensure that it does not overfit. This is
 because an overfitting model has very high variance due to which a slight change in data and its
 underlying pattern affects the model prediction heavily. Though such a model can properly identify
 all the patterns of a training data, but it fails to identify the patterns when made to run on unseen
 test data.
- To state in a different way, a model should not be too complex to be robust and generalizable. The
 model should be as simple as possible, though this may reduce its accuracy as a trade-off, but it
 will be more robust and generalisable.
- From accuracy point of view, a too complex model will have a very high accuracy. This can be also understood using the Bias-Variance trade-off.
- So, to make a model more robust and generalizable, we will have to decrease variance which will lead to some amount of bias. Addition of bias will consequently reduce the accuracy of the model.
- The simpler the model the more the bias it will acquire but lesser will be its variance and hence more generalizable. Its implication in terms of accuracy is that a robust and generalisable model

- will perform equally well on both training and test data i.e. the accuracy does not change much for training and test data.
- In general, we must find some balance between model accuracy and complexity. It is important to have balance in Bias and Variance to avoid overfitting and under-fitting of data. This can be achieved by Regularization techniques like Ridge Regression and Lasso.
- **Bias:** Bias is the error in a model due to which the model is weak to learn from the dataset. High bias means model is unable to learn details of the underlying pattern in the data. As a result, the model will perform poorly on training and testing data.
- Variance: Variance is the error in a model due to which the model tries to over learn from the data. High variance means model performs exceptionally well on training data as it has very well trained on the training dataset and its underlying patterns, but it performs poorly on test data which is unknown for the model.