

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?

The optimum values of alpha obtained are 1 for ridge and 0.0001 for Lasso

Original values				On doubling the alpha			
<div><div>Metric</div><div>Ridge Regression</div><div>Lasso Regression</div></div> <div><div>0</div><div>R2 Score (Train)</div><div>0.919236</div><div>0.919190</div></div> <div><div>1</div><div>R2 Score (Test)</div><div>0.882839</div><div>0.882459</div></div> <div><div>2</div><div>RSS (Train)</div><div>12.958133</div><div>12.965575</div></div> <div><div>3</div><div>RSS (Test)</div><div>8.477038</div><div>8.504503</div></div> <div><div>4</div><div>MSE (Train)</div><div>0.012692</div><div>0.012699</div></div> <div><div>5</div><div>MSE (Test)</div><div>0.019310</div><div>0.019372</div></div> <div><div>6</div><div>RMSE (Train)</div><div>0.112657</div><div>0.112689</div></div> <div><div>7</div><div>RMSE (Test)</div><div>0.138960</div><div>0.139185</div></div>				<div><div>Metric</div><div>Ridge Regression</div><div>Lasso Regression</div></div> <div><div>0</div><div>R2 Score (Train)</div><div>0.917914</div><div>0.917914</div></div> <div><div>1</div><div>R2 Score (Test)</div><div>0.883121</div><div>0.883121</div></div> <div><div>2</div><div>RSS (Train)</div><div>13.170236</div><div>13.170236</div></div> <div><div>3</div><div>RSS (Test)</div><div>8.456653</div><div>8.456653</div></div> <div><div>4</div><div>MSE (Train)</div><div>0.012899</div><div>0.012899</div></div> <div><div>5</div><div>MSE (Test)</div><div>0.019263</div><div>0.019263</div></div> <div><div>6</div><div>RMSE (Train)</div><div>0.113575</div><div>0.113575</div></div> <div><div>7</div><div>RMSE (Test)</div><div>0.138793</div><div>0.138793</div></div>			

Observation: There is a minute decrease (3rd order decimal) to R2 in the train and test set for both Ridge and Lasso.

Coefficients changes for Predictor variables

Original values		On doubling the alpha	
<div>OverallQual_Excellent0.213995</div> <div>Neighborhood_NridgHt0.149678</div> <div>Neighborhood_NoRidge0.138238</div> <div>Neighborhood_Crawfor0.132234</div> <div>OverallQual_Very Good0.131093</div> <div>Neighborhood_StoneBr0.119630</div> <div>Neighborhood_Somerst0.119350</div> <div>OverallCond_Excellent0.109616</div> <div>OverallQual_Very Excellent0.104470</div> <div>Neighborhood_ClearCr0.086901</div> <div>MSSubClass_2-1/2 STORY ALL AGES0.072413</div> <div>SaleCondition_Partial0.068016</div> <div>Neighborhood_Veenker0.066420</div> <div>GrLivArea0.063446</div> <div>BsmtExposure_Gd0.060882</div> <div>OverallCond_Very Good0.054710</div> <div>SaleCondition_Normal0.054673</div> <div>OverallQual_Good0.052662</div> <div>Exterior2nd_MetalSd0.049731</div> <div>2ndFlrSF0.047862</div> <div>Name: Ridge, dtype: float64</div>		<div>OverallQual_Excellent0.200143</div> <div>Neighborhood_NridgHt0.135041</div> <div>OverallQual_Very Good0.126250</div> <div>Neighborhood_Crawfor0.122665</div> <div>Neighborhood_NoRidge0.119018</div> <div>Neighborhood_Somerst0.109239</div> <div>OverallCond_Excellent0.104502</div> <div>Neighborhood_StoneBr0.102728</div> <div>OverallQual_Very Excellent0.083833</div> <div>Neighborhood_ClearCr0.073615</div> <div>MSSubClass_2-1/2 STORY ALL AGES0.066412</div> <div>SaleCondition_Partial0.066247</div> <div>GrLivArea0.064248</div> <div>BsmtExposure_Gd0.060328</div> <div>OverallCond_Very Good0.056056</div> <div>Neighborhood_Veenker0.055531</div> <div>SaleCondition_Normal0.054291</div> <div>OverallQual_Good0.049735</div> <div>2ndFlrSF0.048289</div> <div>MSSubClass_1-STORY 1946 & NEWER ALL STYLES0.041306</div> <div>Name: Ridge, dtype: float64</div>	
<div>OverallQual_Excellent0.225407</div> <div>Neighborhood_NridgHt0.152167</div> <div>Neighborhood_NoRidge0.143192</div> <div>Neighborhood_Crawfor0.141590</div> <div>OverallQual_Very Good0.136055</div> <div>Neighborhood_Somerst0.134324</div> <div>Neighborhood_StoneBr0.121012</div> <div>OverallQual_Very Excellent0.109521</div> <div>OverallCond_Excellent0.109516</div> <div>Neighborhood_ClearCr0.092407</div> <div>Neighborhood_Veenker0.068806</div> <div>SaleCondition_Partial0.065214</div> <div>BsmtExposure_Gd0.061631</div> <div>2ndFlrSF0.059537</div> <div>MSSubClass_2-1/2 STORY ALL AGES0.058001</div> <div>OverallQual_Good0.055116</div> <div>OverallCond_Very Good0.053875</div> <div>SaleCondition_Normal0.052154</div> <div>GrLivArea0.047632</div> <div>1stFlrSF0.046092</div> <div>Name: Lasso, dtype: float64</div>		<div>OverallQual_Excellent0.219288</div> <div>OverallQual_Very Good0.132459</div> <div>Neighborhood_NridgHt0.130515</div> <div>Neighborhood_Crawfor0.129373</div> <div>Neighborhood_NoRidge0.118820</div> <div>Neighborhood_Somerst0.114878</div> <div>GrLivArea0.108221</div> <div>OverallCond_Excellent0.105081</div> <div>Neighborhood_StoneBr0.096759</div> <div>OverallQual_Very Excellent0.085727</div> <div>Neighborhood_ClearCr0.075824</div> <div>SaleCondition_Partial0.060788</div> <div>BsmtExposure_Gd0.060766</div> <div>OverallCond_Very Good0.056168</div> <div>OverallQual_Good0.052386</div> <div>SaleCondition_Normal0.050059</div> <div>Neighborhood_Veenker0.049828</div> <div>MSSubClass_2-1/2 STORY ALL AGES0.046595</div> <div>GarageCars0.040859</div> <div>LotConfig_CulDSac0.040108</div> <div>Name: Lasso, dtype: float64</div>	

Predictor variables are still the same but the order has changed.

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?

In this case where we have too many variables it's better to use Lasso.

Question 3

After building the model, you realized 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?

```
OverallCond_Excellent      0.127606
MSSubClass_2-1/2 STORY ALL AGES  0.090685
2ndFlrSF                   0.072037
BsmtExposure_Gd            0.067025
1stFlrSF                   0.058205
Name: Lasso, dtype: float64
```

Now we have overall condition to be excellent, MSSubCalss- 2-½ Story, 2nd floor Square ft,Good Basement Exposure, 1st floor SF to be Top 5 parameters

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?

We can ensure to keep the model robust and generalisable by ensuring that overfitting is avoided to maximum extent. This is possible by evaluating residual error of the model against predictors to check whether they fit the linear regression conditions (heteroskedasticity, no noticeable pattern between predictors and residuals). Further by employing strategies like Cross validate (GridSearchCV etc.), robust feature engineering (here we introduced Age as feature) and ensuring R2 scores are nominal in both train and test. Too complex model will have high accuracy but will have decreased variance and thereby low bias. By generalizing models, we are to a few extent compromising on accuracy by introducing bias. However, we ensure as we saw here that RMSEs are not overly increased and are fairly regulated to keep the model generalized and thereby avoiding overfitting.