# 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 to double the value of alpha for both ridge and lasso? What will be the most important predictor variables after the change is implemented?**

    The optimal alpha value for **Ridge Regression** is 8 while the optimal value for **Lasso Regression** is 0.001

    If we double the alpha value for both Ridge and Lasso

    - *Alpha value of 16 for Ridge*

| | Before | | | After | | |
|---|---|---|---|---|---|---|
| **Model Score** | Score | RSS | MSE | Score | RSS | MSE |
| | Train 0.928075 | 13.217180 | 0.011316 | Train 0.921735 | 14.382278 | 0.012314 |
| | Test 0.875518 | 6.091068 | 0.020860 | Test 0.873962 | 6.167181 | 0.021120 |

| | Before | | | After | | |
|---|---|---|---|---|---|---|
| | Feature | Coefficient | Affect to Housing Price (%) | Feature | Coefficient | Affect to Housing Price (%) |
| **Most important features** | 0 OverallQual_9 | 1.116580 | 11% | 0 OverallQual_9 | 1.094483 | 9% |
| | 1 OverallCond_3 | 0.885644 | -11% | 1 OverallCond_3 | 0.905350 | -9% |
| | 2 Neighborhood_StoneBr | 1.107288 | 10% | 2 CentralAir_Y | 1.090714 | 9% |
| | 3 Neighborhood_Crawfor | 1.097302 | 9% | 3 Neighborhood_Crawfor | 1.086704 | 8% |
| | 4 OverallCond_9 | 1.098291 | 9% | 4 Neighborhood_StoneBr | 1.074609 | 7% |
| | 5 CentralAir_Y | 1.095559 | 9% | 5 GrLivArea | 1.067843 | 6% |
| | 6 Neighborhood_MeadowV | 0.916628 | -8% | 6 Condition1_Norm | 1.060987 | 6% |
| | 7 MSZoning_RL | 1.070519 | 7% | 7 OverallQual_3 | 0.937040 | -6% |
| | 8 OverallQual_2 | 0.924474 | -7% | 8 OverallQual_8 | 1.060895 | 6% |
| | 9 OverallQual_3 | 0.925317 | -7% | 9 OverallCond_7 | 1.063464 | 6% |

*Insights*:

  - Train score decreases from 0.928 to 0.922, Test score decreases from 0.876 to 0.874
  - Train RSS increases from 13.217 to 14.382, Test RSS increases from 6.091 to 6.167
  - Train MSE increases from 0.011 to 0.012, Test MSE in increases from 0.020 to 0.021
  - Most important features after the change:
      - OverallQual_9
      - OverallCond_3
      - CentralAir_Y
      - Neighborhood_Crawfor
      - Neighborhood_StoneBr
      - GrLivArea
      - Condition1_Norm
      - OverallQual_3

- OverallQual_8
- OverallCond_7

- *Alpha value of 0.002 for Lasso*

| | Before | After |
|---|---|---|
| **Model Score** | | |

**Before — Model Score**

| | Score | RSS | MSE |
|---|---|---|---|
| Train | 0.909501 | 16.630436 | 0.014238 |
| Test | 0.876824 | 6.027184 | 0.020641 |

**After — Model Score**

| | Score | RSS | MSE |
|---|---|---|---|
| Train | 0.893128 | 19.639088 | 0.016814 |
| Test | 0.856143 | 7.039094 | 0.024106 |

**Most important features — Before**

| | Feature | Coefficient | Affect to Housing Price (%) |
|---|---|---|---|
| 0 | OverallQual_9 | 1.168468 | 16% |
| 1 | OverallCond_3 | 0.853955 | -14% |
| 2 | CentralAir_Y | 1.122521 | 12% |
| 3 | Neighborhood_Crawfor | 1.119126 | 11% |
| 4 | GrLivArea | 1.106896 | 10% |
| 5 | Neighborhood_Somerst | 1.079177 | 7% |
| 6 | Neighborhood_StoneBr | 1.071337 | 7% |
| 7 | OverallQual_8 | 1.076084 | 7% |
| 8 | Functional_Typ | 1.070031 | 7% |
| 9 | SaleType_New | 1.072012 | 7% |

**Most important features — After**

| | Feature | Coefficient | Affect to Housing Price (%) |
|---|---|---|---|
| 0 | GrLivArea | 1.128568 | 12% |
| 1 | OverallQual_9 | 1.116716 | 11% |
| 2 | CentralAir_Y | 1.118783 | 11% |
| 3 | OverallCond_3 | 0.896748 | -10% |
| 4 | Neighborhood_Crawfor | 1.092824 | 9% |
| 5 | TotalBsmtSF | 1.066509 | 6% |
| 6 | Functional_Typ | 1.066514 | 6% |
| 7 | Condition1_Norm | 1.057364 | 5% |
| 8 | OverallQual_8 | 1.054059 | 5% |
| 9 | OverallCond_4 | 0.947241 | -5% |

| Feature Eliminated | 191 | 221 |
|---|---|---|

**Insights:**

- Train score decreases from 0.91 to 0.893, Test score decreases from 0.877 to 0.856
- Train RSS increases from 16.630 to 19.639, Test RSS increases from 6.027 to 7.039
- Train MSE increases from 0.014 to 0.017, Test MSE in increases from 0.021 to 0.024
- Number of eliminated features increases from 191 to 221
- Most important features after the change:
    - GrLivArea
    - OverallQual_9
    - CentralAir_Y
    - OverallCond_3
    - Neighborhood_Crawfor
    - TotalBsmtSF
    - Functional_Typ
    - Condition1_Norm
    - OverallQual_8
    - OverallCond_4

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?**

   The optimal value of lambda for **Ridge** is 8, and for **Lasso** is 0.001

   Both models score is good, but the business goal is to find the most important features so feature selection should be performed. Also, feature elimination helps making the model simple and robust. So, **we should choose Lasso in this scenario to eliminate less important features**

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?**

Top 5 variables to be removed: **OverallQual_9, CentralAir_Y, Neighborhood_Crawfor, GrLivArea, Neighborhood_StoneBr**

After using cross validation, optimal alpha value for Lasso Model is now **0.0001**

|  | Score | RSS | MSE |
|---|---|---|---|
| Train | 0.937855 | 11.419959 | 0.009777 |
| Test | 0.826455 | 8.491802 | 0.029082 |

Top 5 variables and its effect to housing price per 1 unit are:

|  | Feature | Coefficient | Affect to Housing Price (%) |
|---|---|---|---|
| 0 | Condition2_PosN | 0.417991 | -58% |
| 1 | MSZoning_FV | 1.523940 | 52% |
| 2 | MSZoning_RH | 1.521526 | 52% |
| 3 | MSZoning_RL | 1.497821 | 49% |
| 4 | MSZoning_RM | 1.444594 | 44% |

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

To make sure a model is robust and generalizable, we should:

- The robustness and generalization of a model can be achieved by removing the impact of outliers. For example, regression-based models are usually affected by outliers, while tree-based models are not

- Switching from mean squared error to mean absolute difference also helps reducing the impact of outliers.

- Keeping the model as simple as possible to avoid overfitting, using Lasso instead of Ridge to make use of its feature selection.

- Transforming the data to reduce skewness, for example, use a log transform for a skewed distribution of data.

- Removing outliers in training data, those outliers could badly affect the accuracy of the model

There are trade-offs between accuracy and robustness, we need to keep balance between them like the balance between bias-variance trade-offs, keep the model as simple as possible but don't make it underfitting.