

**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:**

The optimal value of alpha in Ridge and Lasso regression depends on the specific dataset and the problem you are trying to solve. Cross-validation is commonly used to find the optimal alpha value.

In Ridge regression, alpha is a regularization parameter that controls the strength of the penalty term. A higher alpha results in stronger regularization. The coefficients may be reduced close to zero but will never be zero.

In Lasso regression, alpha also controls the strength of the penalty term, but it has the additional effect of causing some coefficients to be exactly zero.

Doubling the value of alpha, increasing alpha in Ridge and Lasso regression will generally increase the strength of regularization. It the penalty on the coefficients. This can result in a sparser model, where more coefficients are shrunk towards zero. This might lead to a model with fewer features, especially in Lasso regression where some coefficients can be driven exactly to zero. In Ridge regression, all coefficients will still be present, but their magnitudes may be further reduced.

The non-zero coefficients in Lasso regression after doubling alpha are the most important predictor variables in the updated model. In Ridge regression, all coefficients will still be present, but their magnitudes may be further reduced.

In the assignment case alpha for Ridge and Lasso is 20 and 0.001 respectively.

Change in model is not too much is :

For Ridge:

R2 score on test and train before doubling

0.9167125610310837

0.8879704427284993

R2 score after doubling :

0.9091887698214963

0.8863063649417942

For Lasso:

R2 score on test and train before doubling

0.9088060355428663

0.8840774918446787

R2 score after doubling :

0.8932614411966211

0.8781591442944001

After doubling the predictor variables are :

Top 5 predictor variables in

Ridge :

LotFrontage

OverallCond

BsmtFullBath

Neighborhood\_Crawfor

Condition1\_Norm

Lasso:

BsmtUnfSF

RoofMatl\_WdShake

TotalBsmtSF

RoofMatl\_WdShngl

RoofMatl\_CompShg

### **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:**

Key Considerations for choosing Regression model:

Feature Importance:

Lasso is preferred to identify a subset of the most important features for model interpretation and potentially reducing data collection costs.

Ridge is preferred to retain all features with varying levels of importance and don't need explicit feature selection.

Multicollinearity:

If your dataset has high multicollinearity (highly correlated features), Ridge is often more effective as it shrinks coefficients more evenly, reducing the impact of correlated features without completely removing them.

Model Interpretability:

Lasso produces sparse models with fewer non-zero coefficients, making them easier to interpret and explain, especially for non-technical audiences.

Ridge models retain all features, potentially making interpretation more complex, especially with many features.

Predictive Performance:

Both models seems to achieve good predictive performance under different conditions.

However I picked top 5 from Lasso for simplicity. However in real world I would dive deep into discussions about problem/expected outcome with domain experts before taking final decision and tweaking the model accordingly .

**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:**

In real world scenarios I would look for alternate variables which could be used instead or they just could be excluded out of new model, else based on calculation with Python code:

After removing top 5 of the from Lasso and rebuilding model, new top 5 predictors are (along with their coefficients ):

```
RoofMatl_Roll 0.113475
BsmtUnfSF    0.091173
TotalBsmtSF  0.041387
RoofMatl_Tar&Grv 0.030704
MSZoning_FV 0.020243
```

Revised R2 for train and test are :  
0.9088936520670919  
0.8845694752217836

**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:**

By following these steps, we can develop models that are not only accurate but also reliable and effective in real-world applications.

1. Data Quality:
  - Clean and pre-process data to remove noise, outliers, and biases.
  - Use diverse and representative training data to capture real-world scenarios.
2. Regularization Techniques:
  - Implement techniques like Lasso, Ridge, or Dropout to prevent overfitting on training data.
  - This reduces model complexity and improves performance on unseen data.
3. Model Selection and Evaluation:
  - Use cross-validation and hold-out sets to measure performance on unseen data.
  - Don't solely rely on training accuracy, as it can be misleading for generalizability.

#### 4. Monitoring and Adaptation:

- Continuously monitor model performance in production and be ready to adapt.
- This might involve retraining or updating the model with new data over time.

#### Implications for accuracy:

- Higher quality data and regularization usually lead to improved generalization.
- Overfitting to training data might give high apparent accuracy but poor performance on unseen data.
- Robust and generalizable models offer consistent accuracy across different data sets and scenarios.