

## Assignment-based Subjective Questions

### Question 1

**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?**

Ans:

The optimal value of alpha for

- Ridge is 2
- Lasso is 0.01

R2 value with the above alphas is 0.85

After doubling the alpha values in the Ridge and Lasso, the prediction accuracy remains around 0.84, but there is a small change in the co-efficient values. Below are the changes in the co-efficient.

### Lasso Regression:

Alpha coefficient

```
[('constant', 12.018),
 ('MSSubClass', -0.02),
 ('LotFrontage', 0.0),
 ('LotArea', 0.01),
 ('OverallQual', 0.131),
 ('OverallCond', 0.041),
 ('MasVnrArea', 0.0),
 ('BsmtFinSF1', 0.0),
 ('BsmtFinSF2', 0.0),
 ('TotalBsmtSF', 0.021),
 ('1stFlrSF', 0.006),
 ('2ndFlrSF', 0.0),
 ('LowQualFinSF', -0.0),
 ('GrLivArea', 0.108),
 ('BsmtFullBath', 0.033),
 ('BsmtHalfBath', 0.0),
 ('FullBath', 0.017),
```

```
('HalfBath', 0.001),
('BedroomAbvGr', 0.0),
('KitchenAbvGr', -0.003),
('Fireplaces', 0.032),
('GarageArea', 0.045),
('WoodDeckSF', 0.014),
('OpenPorchSF', 0.0),
('EnclosedPorch', 0.0),
('3SsnPorch', 0.0),
('ScreenPorch', 0.004)]
```

### Alpha Doubled :

```
[('constant', 12.022),
 ('MSSubClass', -0.011),
 ('LotFrontage', 0.0),
 ('LotArea', 0.005),
 ('OverallQual', 0.137),
 ('OverallCond', 0.026),
 ('MasVnrArea', 0.0),
 ('BsmtFinSF1', 0.0),
 ('BsmtFinSF2', 0.0),
 ('TotalBsmtSF', 0.019),
 ('1stFlrSF', 0.009),
 ('2ndFlrSF', 0.0),
 ('LowQualFinSF', -0.0),
 ('GrLivArea', 0.1),
 ('BsmtFullBath', 0.024),
 ('BsmtHalfBath', 0.0),
 ('FullBath', 0.009),
 ('HalfBath', 0.0),
 ('BedroomAbvGr', 0.0),
 ('KitchenAbvGr', -0.0),
 ('Fireplaces', 0.031),
 ('GarageArea', 0.046),
 ('WoodDeckSF', 0.011),
 ('OpenPorchSF', 0.0),
 ('EnclosedPorch', -0.0),
 ('3SsnPorch', 0.0),
 ('ScreenPorch', 0.0),
 ('PoolArea', -0.004)]
```

### Ridge Regression Model

#### Alpha coefficients:

```
[('constant', 11.459),
 ('MSSubClass', -0.024),
 ('LotFrontage', -0.01),
```

```
( 'LotArea', 0.019),
( 'OverallQual', 0.072),
( 'OverallCond', 0.041),
( 'MasVnrArea', -0.004),
( 'BsmtFinSF1', -0.012),
( 'BsmtFinSF2', 0.005),
( 'TotalBsmtSF', 0.007),
( '1stFlrSF', 0.044),
( '2ndFlrSF', 0.047),
( 'LowQualFinSF', 0.003),
( 'GrLivArea', 0.071),
( 'BsmtFullBath', 0.028),
( 'BsmtHalfBath', 0.004),
( 'FullBath', 0.019),
( 'HalfBath', 0.013),
( 'BedroomAbvGr', 0.017)]
```

### Alpha doubled:

```
[ ('constant', 11.576),
  ('MSSubClass', -0.025),
  ('LotFrontage', -0.01),
  ('LotArea', 0.018),
  ('OverallQual', 0.076),
  ('OverallCond', 0.043),
  ('MasVnrArea', -0.005),
  ('BsmtFinSF1', -0.012),
  ('BsmtFinSF2', 0.005),
  ('TotalBsmtSF', 0.007),
  ('1stFlrSF', 0.043),
  ('2ndFlrSF', 0.047),
  ('LowQualFinSF', 0.003),
  ('GrLivArea', 0.071),
  ('BsmtFullBath', 0.029),
  ('BsmtHalfBath', 0.004),
  ('FullBath', 0.02),
  ('HalfBath', 0.013),
  ('BedroomAbvGr', 0.018)]
```

Since the alpha values are small, change in the model after doubling the alpha is not seen.

### 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?**

- The optimum lambda value in case of Ridge and Lasso is 2 for Ridge and 0.01 for Lasso
- The Mean Squared Error in case of Ridge and Lasso are:
  - Ridge - 0.13515950395853346
  - Lasso - 0.15188950651109254
- The Mean Squared Error of both the models are almost same.
- Since Lasso helps in feature reduction (as coefficient of some features become zero), Lasso is preferred over Ridge and should be used as the final model.

### 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?**

Ans: The five most important predictor variables in the current lasso model is:-

Coeff	
OverallQual	0.131
GrLivArea	0.108
GarageArea	0.045
OverallCond	0.041
BsmtFullBath	0.033
Fireplaces	0.032

The R2 score without the top 5 predictors drops to .78

The Mean Squared Error increases to 0.0028575670906482538

The new Top 5 predictors are:-

BsmtFullBath 0.033

Fireplaces 0.032

TotalBsmtSF 0.021

FullBath 0.017

WoodDeckSF 0.014

#### 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?**

As Per, Occam's Razor— given two models that show similar 'performance' in the finite training or test data, we should pick the one that makes fewer on the test data due to following reasons:-

- Simpler models are usually more 'generic' and are more widely applicable
- Simpler models require fewer training samples for effective training than the more complex ones and hence are easier to train.
- Simpler models are more robust.
- Complex models tend to change wildly with changes in the training data set
- Simple models have low variance, high bias and complex models have low bias, high variance.
- Simpler models make more errors in the training set. Complex models lead to overfitting. They work very well for the training samples, fail miserably when applied to other test samples

Therefore, to make the model more robust and generalizable, make the model simple but not simpler which will not be of any use.

Regularization can be used to make the model simpler. Regularization helps to strike the delicate balance between keeping the model simple and not making it too naive to be of any use. For regression, regularization involves adding a

regularization term to the cost that adds up the absolute values or the squares of the parameters of the model.

Also, Making a model simple leads to Bias-Variance Trade-off:

- A complex model will need to change for every little change in the dataset and hence is very unstable and extremely sensitive to any changes in the training data.
- A simpler model that abstracts out some pattern followed by the data points given is unlikely to change wildly even if more points are added or removed.

Bias quantifies how accurate is the model likely to be on test data. A complex model can do an accurate job prediction provided there is enough training data. Models that are too naïve, for e.g., one that gives same answer to all test inputs and makes no discrimination whatsoever has a very large bias as its expected error across all test inputs are very high.

Variance refers to the degree of changes in the model itself with respect to changes in the training data.

Thus accuracy of the model can be maintained by keeping the balance between Bias and Variance as it minimizes the total error