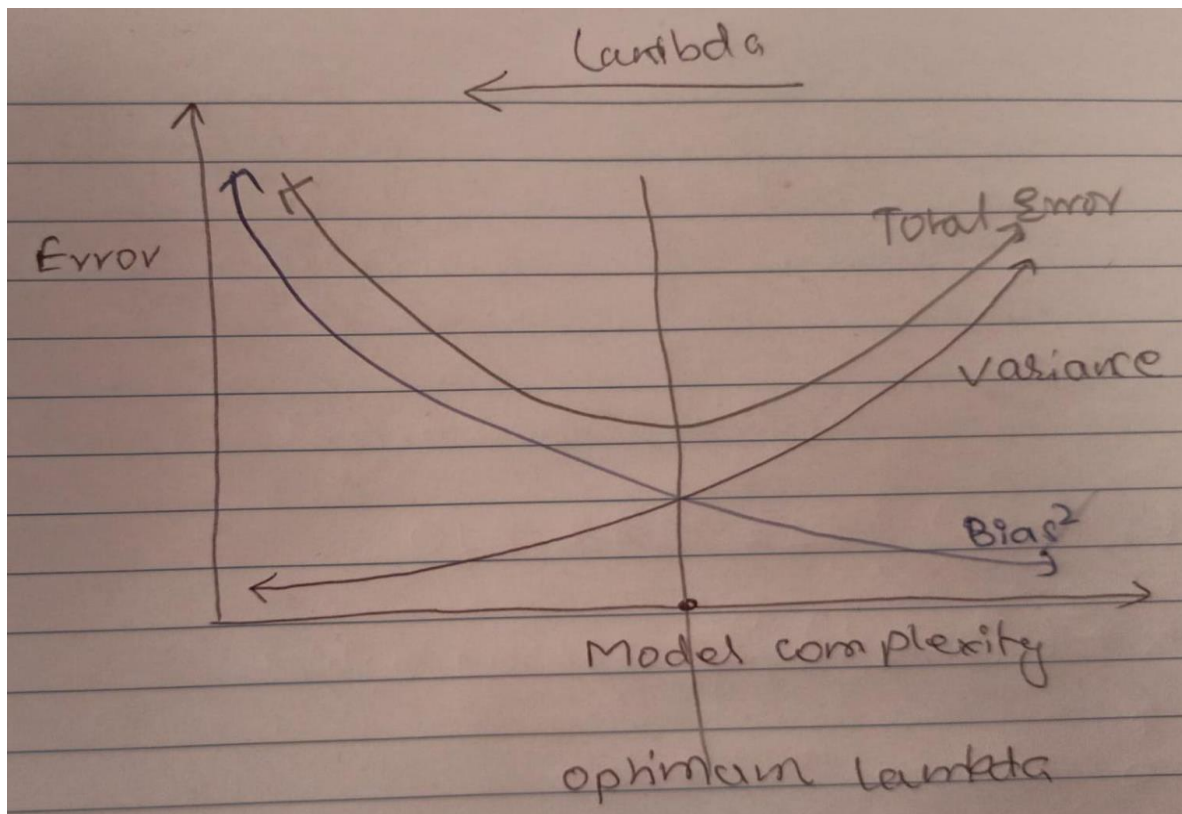# Question 1

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 1:

The optimal value of alpha ($\alpha$) for ridge and lasso regression isn't predetermined and should be found using model selection methods such as:

1. **Cross-Validation**: Performing k-fold cross-validation to evaluate model performance across a range of $\alpha$ values and choosing the one that minimizes cross-validation error, like mean squared error in regression tasks.
2. **Grid Search**: Running a grid search over various $\alpha$ values to find the one that delivers the best model performance according to a chosen metric.
3. **Automated Techniques**: Utilizing automated methods like GridSearchCV class in sklearn.model_selection, which automatically conduct grid search and cross-validation to determine the optimal $\alpha$.
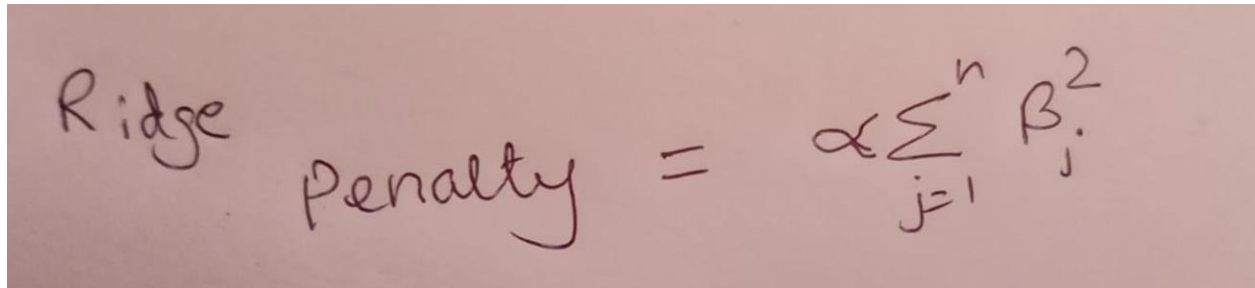
## Impact of Doubling the Value of Alpha

When α is doubled in ridge and lasso regression, the following changes typically occur:

**Ridge Regression**

Ridge regression adds a penalty proportional to the sum of squared coefficients to the loss function.
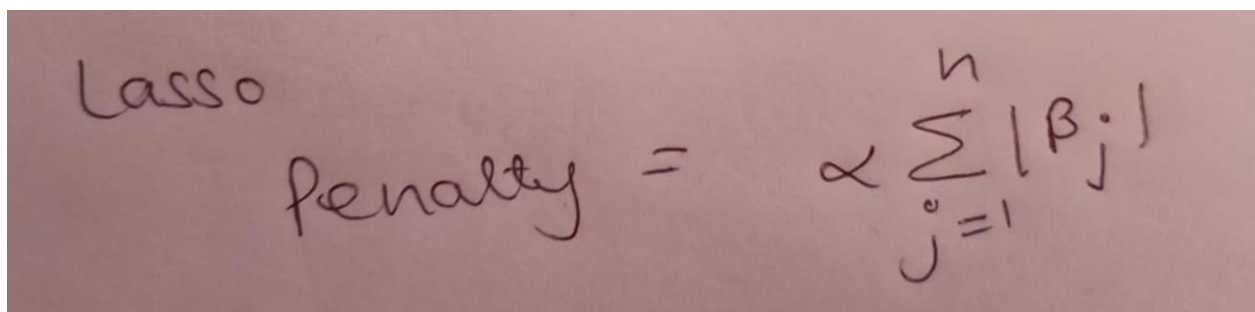
$$\text{Ridge Penalty} = \alpha \sum_{j=1}^{n} \beta_j^2$$

Increasing α results in:

- **Stronger Regularization**: Coefficients are more strongly shrunk toward zero, reducing model complexity.
- **Bias-Variance Tradeoff**: Higher α decreases variance but increases bias, which might improve model performance on unseen data if the original model was overfitting.
- **Coefficient Magnitudes**: Coefficients' magnitudes decrease, leading to a simpler model that might generalize better but could lose some precision.

**Lasso Regression**

Lasso regression adds a penalty proportional to the sum of absolute values of coefficients to the loss function.

$$\text{Lasso Penalty} = \alpha \sum_{j=1}^{n} |\beta_j|$$

Increasing α leads to:

- **Stronger Regularization**: Coefficients are pushed more strongly toward zero, with some becoming exactly zero, effectively performing feature selection.
- **Sparse Solution**: More coefficients are likely to be set to zero with a higher α, significantly simplifying the model by reducing the number of predictors.

- **Bias-Variance Tradeoff**: Similar to ridge regression, higher α results in increased bias and decreased variance.

## Key Predictor Variables After Doubling Alpha

The most important predictor variables after doubling α depend on the specific dataset and context, but generally:

- **Ridge Regression**: Important variables will still have the largest coefficients (in absolute terms), although these coefficients will be smaller compared to the original α.
- **Lasso Regression**: Many coefficients may become zero, leaving only a subset of the original predictors. The remaining non-zero coefficients indicate the most important predictor variables.

# 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 2:

This is metric generated after the model run.

| | Metric | Linear Regression | Ridge Regression | Lasso Regression |
|---|---|---|---|---|
| 0 | R2 Score (Train) | 1.000000e+00 | 0.836155 | 0.838991 |
| 1 | R2 Score (Test) | 1.000000e+00 | 0.839991 | 0.839198 |
| 2 | RSS (Train) | 8.647749e-29 | 2.016121 | 1.981222 |
| 3 | RSS (Test) | 3.269522e-29 | 0.869780 | 0.874089 |
| 4 | MSE (Train) | 2.910306e-16 | 0.044437 | 0.044051 |
| 5 | MSE (Test) | 2.732153e-16 | 0.044562 | 0.044673 |

When deciding whether to apply ridge regression or lasso regression based on the provided metrics, several factors should be considered:

1. **Overfitting and Generalization**:
   - **Linear Regression**: It shows perfect fit on training data with R2 of 1.0, but this likely indicates overfitting since the performance on the test data might not generalize well.

- **Ridge and Lasso Regression**: Both have significantly lower R2 scores on the training data compared to linear regression, indicating that they are applying regularization and preventing overfitting.
2. **Test Performance**:
   - **Ridge Regression**: R2 on test data is 0.839991, with RSS of 0.869780 and MSE of 0.044562.
   - **Lasso Regression**: R2 on test data is 0.839198, with RSS of 0.874089 and MSE of 0.044673.
3. **Comparative Analysis**:
   - **R2 Scores**: Ridge regression has a slightly higher R2 score on the test data (0.839991) compared to lasso regression (0.839198), suggesting a marginally better fit to the test data.
   - **RSS and MSE**: Ridge regression has a slightly lower RSS (0.869780) and MSE (0.044562) on the test data compared to lasso regression (RSS of 0.874089 and MSE of 0.044673), indicating that ridge regression has a marginally better prediction accuracy.

## Choosing Between Ridge and Lasso Regression

Given the above analysis, **ridge regression** seems to be the preferable choice based on the following reasons:

1. **Test Performance**: Ridge regression demonstrates marginally better performance on the test data across all metrics (R2, RSS, and MSE), suggesting that it generalizes slightly better to unseen data compared to lasso regression.
2. **Regularization Preference**: Ridge regression is preferable when all predictors are expected to have a small but non-zero effect because it shrinks coefficients but does not set them exactly to zero. In contrast, lasso regression might be preferred if you expect that only a few predictors are truly important and want a model that performs automatic feature selection by setting some coefficients to zero.

Given the metrics provided, ridge regression offers a slightly better generalization performance and would be the recommended model to apply in this scenario.

# 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 3:

The five most important predictors variables are GrLivArea, OverallQual, GarageCars, KitchenQual and BsmtQual. Will run all three models after excluding the variables.

Given that the five most important predictor variables identified by the lasso model are no longer available in the incoming data, you'll need to build a new model excluding these variables. The variables and their correlations with the SalePrice were:

- GrLivArea with SalesPrice:0.7086244776126515
- OverAllQual with SalePrice:0.7909816005838053
- GarageCars with SalePrice:0.6404091972583519
- KitchenQual with SalesPrice:-0.589188778299423
- BsmtQual with SalePrice:-0.6208861300191687

With these variables removed, you will need to identify the next set of most important predictor variables. This can be done by refitting the lasso model on the remaining predictors. Here's a step-by-step approach to determine the next most important variables:

## Steps to Identify the Next Most Important Predictor Variables

1. **Remove the Important Predictors**: Exclude GrLivArea, OverAllQual, GarageCars, KitchenQual, and BsmtQual from the dataset.
2. **Re-fit the Lasso Model**: Fit a new lasso regression model on the dataset without these variables.
3. **Determine the New Important Variables**: Extract the coefficients of the new model to identify the next most important predictors.

Steps are performed in the notebook. Following five predictors' variables can be concluded.

- 1stFlrSF with SalePrice: 0.6058521846919153
- 2ndFlrSF with SalePrice: 0.31933380283206736
- TotalBsmtSF with SalePrice: 0.6135805515591943
- ExterQual with SalePrice: -0.6368836943991116
- KitchenAbvGr with SalePrice: -0.13590737084214105

# Question 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?

Answer 4:

To make sure that a model is robust and generalizable, it is important to follow a set of best practices and validation techniques throughout the model development process. Here are key steps and their implications for the model's accuracy:

1. **Data Preprocessing and Cleaning**:
   - **Handle Missing Values**: Properly impute or remove missing data to avoid biases.
   - **Outlier Treatment**: Detect and handle outliers which can skew the model.

- o **Feature Scaling**: Standardize or normalize features to ensure consistent scale, especially important for regularization methods.

**Implication**: Ensures that the model learns from clean, well-processed data, improving its ability to generalize.

2. **Cross-Validation**:
   - o **K-Fold Cross-Validation**: Split the dataset into k subsets and train the model k times, each time using a different subset as the validation set and the remaining data as the training set.

**Implication**: Provides a better estimate of the model's performance on unseen data, reducing the risk of overfitting and giving a more reliable measure of accuracy.

3. **Regularization**:
   - o **L1 (Lasso)** and **L2 (Ridge) Regularization**: Penalize large coefficients to prevent overfitting and enhance model generalizability.

**Implication**: Reduces the complexity of the model, improving its ability to generalize to new data at the cost of potentially slight reductions in training accuracy.

4. **Hyperparameter Tuning**:
   - o **Grid Search** and **Random Search**: Systematically explore the hyperparameter space to find the optimal settings.

**Implication**: Fine-tunes the model to balance bias and variance, leading to a more generalizable model.

5. **Model Validation**:
   - o **Holdout Validation**: Reserve a separate validation set that is not used during training for final model evaluation.
   - o **Nested Cross-Validation**: Useful for hyperparameter tuning and model assessment, providing an unbiased evaluation of the model's performance.

**Implication**: Ensures the model's performance metrics are realistic and not overly optimistic due to overfitting on the training data.

6. **Feature Engineering and Selection**:
   - o **Automatic Feature Selection**: Techniques like recursive feature elimination (RFE) or regularization-based methods (lasso).
   - o **Domain Knowledge**: Incorporate domain-specific insights to engineer relevant features.

**Implication**: Improves the model's predictive power and interpretability, focusing on the most important features.

7. **Model Evaluation Metrics**:
   o Use appropriate evaluation metrics (e.g., precision, recall, F1-score for classification; RMSE, MAE for regression) beyond simple accuracy to assess model performance.

   **Implication**: Provides a more comprehensive understanding of model performance, especially in imbalanced datasets.

8. **Robustness Checks**:
   o **Sensitivity Analysis**: Assess how changes in input variables affect the model's output.
   o **Stress Testing**: Evaluate model performance under extreme or edge cases.

   **Implication**: Ensures the model remains reliable and stable under a variety of conditions and inputs.

## Implications for Model Accuracy and Generalizability

1. **Bias-Variance Trade-off**:
   o **High Bias (Underfitting)**: Simplistic models that don't capture the underlying patterns well. These models have poor performance on both training and test data.
   o **High Variance (Overfitting)**: Complex models that perform well on training data but poorly on test data due to capturing noise instead of the underlying pattern.

   **Implication**: Regularization and cross-validation help find a balance, improving test accuracy even if it slightly reduces training accuracy.

2. **Generalization**:
   o Models that generalize well perform consistently on both training and unseen test data. Techniques like cross-validation, regularization, and proper preprocessing are crucial for this.

   **Implication**: A generalizable model maintains accuracy across different datasets, ensuring reliable predictions in real-world applications.

3. **Model Complexity**:
   o Simpler models tend to generalize better, while overly complex models risk overfitting.

   **Implication**: Balancing complexity helps maintain or improve test accuracy, even if it means a slight reduction in training accuracy due to reduced overfitting.

## Conclusion

By following these best practices, you ensure your model is robust and generalizable. This approach might lead to a slight decrease in training accuracy but results in improved

performance on unseen data, making the model more reliable and useful in real-world applications. Regularization, proper validation techniques, and careful preprocessing are paramount to achieving a balance that maximizes generalizability and overall accuracy.