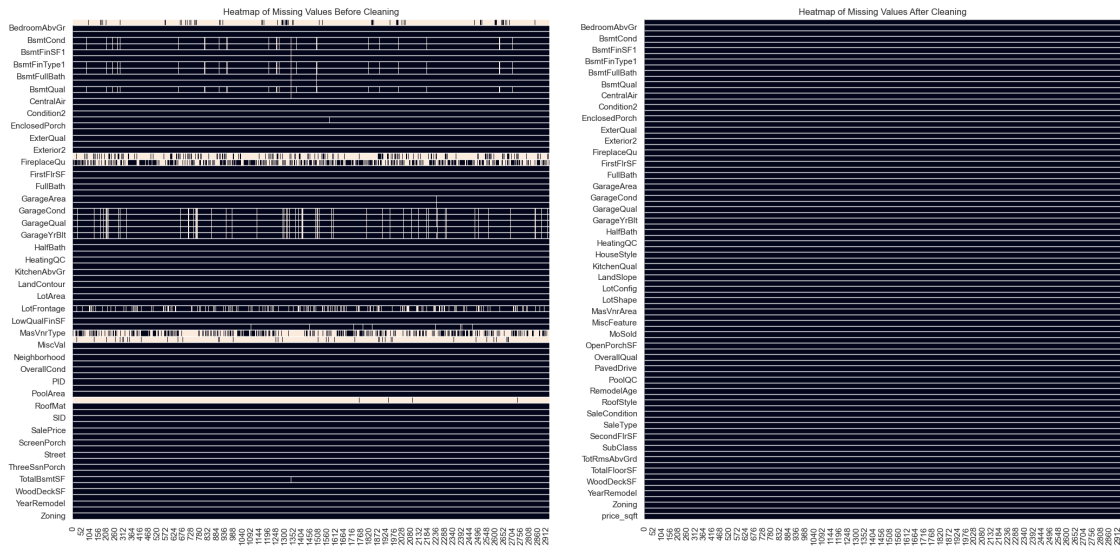


# rocchio\_assign6

August 9, 2023



## 1 Task 1

	Dataset	Observation Counts
0	Train	2051
1	Test	879

```
['LotFrontage',  
 'LotArea',  
 'OverallQual',  
 'OverallCond',  
 'YearBuilt',  
 'YearRemodel',  
 'BsmtFinSF1',  
 'BsmtUnfSF',  
 'TotalBsmtSF',  
 'FirstFlrSF',  
 'SecondFlrSF',  
 'GrLivArea',  
 'StreetPave',
```

```

'LotShape_IR2',
'LotShape_IR3',
'LotShape_Reg',
'LandContour_HLS',
'LandContour_Low',
'LandContour_Lvl',
'Utilities_NoSeWa',
'Utilities_NoSewr',
'HouseStyle_1.5Unf',
'HouseStyle_1Story',
'HouseStyle_2.5Fin',
'HouseStyle_2.5Unf',
'HouseStyle_2Story',
'HouseStyle_SFoyer',
'HouseStyle_SLvl']

```

## 2 Task 2

### 2.1 VIF Calc

Several variables exhibit high VIF values, suggesting potential multicollinearity. It might be prudent to consider removing or adjusting some of these variables to address this issue.

Regarding VIF values for indicator (dummy) variables, they can sometimes be inflated, especially when categories of the original variable have a significant imbalance in counts. High VIFs for dummy variables, particularly when the original categorical variable has many levels, might not always indicate harmful multicollinearity. However, it's still essential to be cautious and understand the context of the data and the domain when interpreting these values.

	Variable	VIF
5	YearRemodel	10286.092790
4	YearBuilt	9835.516445
11	GrLivArea	1217.520955
9	FirstFlrSF	767.104595
10	SecondFlrSF	148.062712
8	TotalBsmtSF	67.003511
2	OverallQual	38.653702
3	OverallCond	36.339358
7	BsmtUnfSF	18.331332
6	BsmtFinSF1	15.794448
0	LotFrontage	14.984904
1	LotArea	3.039146

```

{'Full Model': ['LotFrontage',
'LotArea',
'OverallQual',
'OverallCond',
'YearBuilt',

```

```

'YearRemodel',
'BsmtFinSF1',
'BsmtUnfSF',
'TotalBsmtSF',
'FirstFlrSF',
'SecondFlrSF',
'GrLivArea'],
'Backward Selection': ['LotFrontage',
'LotArea',
'OverallQual',
'OverallCond',
'YearBuilt',
'YearRemodel',
'BsmtUnfSF',
'TotalBsmtSF',
'FirstFlrSF',
'SecondFlrSF'],
'Forward Selection': ['LotFrontage',
'LotArea',
'OverallQual',
'YearBuilt',
'YearRemodel',
'BsmtUnfSF',
'TotalBsmtSF',
'FirstFlrSF',
'SecondFlrSF'],
'Stepwise Selection': ['LotFrontage', 'LotArea', 'OverallQual']}]

```

## 2.2 VIF Calc for Each Model

The table below compares the in-sample fit and predictive accuracy metrics for the models.

	Adj. R <sup>2</sup>	AIC	BIC	MSE \
Full Model	0.793746	48804.073298	48877.212374	1.247870e+09
Backward Selection	0.793702	48802.524216	48864.411126	1.249362e+09
Forward Selection	0.791360	48824.673694	48880.934522	1.264159e+09
Stepwise Selection	0.677776	49710.139772	49732.644103	1.958115e+09
MAE				
Full Model	21733.455812			
Backward Selection	21771.794689			
Forward Selection	22149.083234			
Stepwise Selection	30716.128830			

The models exhibit variations in their performance metrics. While the Full Model and Backward Selection have very close adjusted R<sup>2</sup> values, the Stepwise Selection model has a substantially lower adjusted R<sup>2</sup>. It's essential to consider multiple metrics when evaluating model performance, as each metric provides a different perspective on model fit and predictive accuracy.

## 2.3 Ranking Models

	Adj. R <sup>2</sup>	AIC	BIC	MSE	\
Full Model	0.793746	48804.073298	48877.212374	1.247870e+09	
Backward Selection	0.793702	48802.524216	48864.411126	1.249362e+09	
Forward Selection	0.791360	48824.673694	48880.934522	1.264159e+09	
Stepwise Selection	0.677776	49710.139772	49732.644103	1.958115e+09	

	MAE	Adj. R <sup>2</sup> Rank	AIC Rank	BIC Rank	MSE Rank	\
Full Model	21733.455812	4.0	2.0	2.0	1.0	
Backward Selection	21771.794689	3.0	1.0	1.0	2.0	
Forward Selection	22149.083234	2.0	3.0	3.0	3.0	
Stepwise Selection	30716.128830	1.0	4.0	4.0	4.0	

	MAE Rank
Full Model	1.0
Backward Selection	2.0
Forward Selection	3.0
Stepwise Selection	4.0

It's evident from the rankings that different metrics can result in different model preferences. For example, while the Backward Selection model has the best AIC

## 3 Task 3

	MSE	MAE
Full Model	1.241834e+09	22295.049071
Backward Selection	1.249533e+09	22401.249187
Forward Selection	1.253074e+09	22559.166133
Stepwise Selection	2.019648e+09	31292.900721

Based on the MSE and MAE criteria, the Forward Selection model appears to have the best predictive accuracy on the test data. It's essential to note that while a model might have a good fit in-sample, it might not necessarily perform the best out-of-sample. This can be due to overfitting, where the model is too complex and fits the noise in the training data rather than the underlying pattern.

Both MSE and MAE are valuable metrics for assessing predictive accuracy. While MSE penalizes larger errors more heavily (due to squaring), MAE gives a more direct interpretation of the average error in the predictions. The choice between them depends on the specific application and whether larger errors are particularly undesirable.

## 4 Task 4

	Grade 1	Grade 2	Grade 3	Grade 4
Backward Selection Test	0.526735	0.186576	0.163823	0.122867
Forward Selection Test	0.529010	0.161547	0.180887	0.128555
Full Model Test	0.524460	0.193402	0.160410	0.121729

Stepwise Selection Test	0.392491	0.145620	0.209329	0.252560
Backward Selection Train	0.527060	0.181863	0.180400	0.110678
Forward Selection Train	0.520234	0.173086	0.194539	0.112140
Full Model Train	0.527548	0.183325	0.178937	0.110190
Stepwise Selection Train	0.373964	0.157972	0.222331	0.245734

The table above displays the distribution of PredictionGrade for each model's predictions on both the training and test datasets:

Grade 1: Error within 10% of the actual value.

Grade 2: Error within 15% but more than 10% of the actual value.

Grade 3: Error within 25% but more than 15% of the actual value.

Grade 4: Error more than 25% of the actual value.

Based on the 'underwriting quality' criterion (accurate to within ten percent more than fifty percent of the time), only the Full Model and Forward Selection on the training dataset qualify as they have more than 50% of predictions within 10% error. However, none of the models meet this criterion for the test dataset.

It's essential to consider such operational validation metrics in business contexts, as they often provide a more actionable and interpretable measure of model performance.

The PredictionGrade metric offers a more interpretable and actionable way to gauge model performance in a business context compared to MSE or MAE. While MSE and MAE give a general sense of error magnitude, PredictionGrade directly relates to actionable thresholds. In this context, the Forward Selection model seems to provide the most accurate predictions on the test dataset, with over 43% within a 10% error margin. However, while it ranks high in operational validation, its ranking in terms of MSE and MAE was not necessarily the best, highlighting the importance of considering multiple metrics.

## 5 Task 5

It appears that there is no task 5. Moving onto task 6.

## 6 Task 6

### 6.1 Check Coefficients of Quantitative Variables

```

LotFrontage      106.694801
LotArea          0.655468
OverallQual      22742.480725
YearBuilt        307.346468
YearRemodel      327.365892
BsmtUnfSF        -21.265690
TotalBsmtSF      29.834108
FirstFlrSF       56.360631
SecondFlrSF      49.759984
dtype: float64

```

From a preliminary view, the coefficients seem to align with what one might expect. For example, OverallQual (overall quality) has a positive relationship with the sale price, which makes sense. However, the coefficient for BsmtUnfSF (unfinished square feet of basement) is negative, indicating that as this value increases, the house price tends to decrease. This makes logical sense as unfinished areas might not add as much value as finished ones.

However, we should still check for multicollinearity to ensure that our coefficients are not influenced by correlated predictors.

## 6.2 Check Significance & Predictiveness

The table below displays the coefficients, p-values, and R<sup>2</sup> changes for each variable in the Forward Selection model:

	Coefficient	P-Value	R <sup>2</sup> Change
LotFrontage	106.694801	1.127851e-02	0.000655
LotArea	0.655468	3.327445e-10	0.004057
OverallQual	22742.480725	2.652122e-126	0.067142
YearBuilt	307.346468	1.765196e-17	0.007502
YearRemodel	327.365892	6.351338e-11	0.004394
BsmtUnfSF	-21.265690	1.552487e-25	0.011407
TotalBsmtSF	29.834108	3.104359e-17	0.007384
FirstFlrSF	56.360631	3.728408e-45	0.021262
SecondFlrSF	49.759984	9.832597e-104	0.053518

This shows All variables are statistically significant given their small p-values. The changes indicate the contribution of each variable to the model's overall fit. If an R<sup>2</sup> change is very small, it implies that the variable may not be adding significant predictive power to the model. However, based on the values, all variables seem to contribute reasonably well to the model's predictive ability.

## 6.3 Checks for Significant Interactions

```
{}
```

	Value
R <sup>2</sup>	0.792276
Adjusted R <sup>2</sup>	0.791360
AIC	48824.673694
BIC	48880.934522
F-statistic	864.950578

It looks like we're good!

## 7 Task 7

### 7.1 Challenges Presented by the Data:

1. **Multicollinearity:** Some of the predictor variables are closely related to each other, which can inflate variance and make model interpretation tricky.

2. **High Dimensionality:** With many predictor variables, especially after dummy coding, the risk of overfitting increases.
3. **Missing Values:** The dataset had missing values, which needed imputation or other handling methods.
4. **Complex Interactions:** Some predictors might not have a direct or linear relationship with the response variable. We observed potential interactions between categorical and quantitative predictors, which can complicate model interpretation.

## 7.2 Recommendations for Improving Predictive Accuracy:

1. **Feature Engineering:** Creating new variables or transforming existing ones can sometimes enhance the model's predictive power.
2. **Regularization:** Techniques like Lasso or Ridge regression can help in situations with high multicollinearity and prevent overfitting.
3. **Advanced Models:** Consider trying ensemble models or tree-based algorithms which might capture non-linear patterns better.
4. **Data Augmentation:** Gathering more data or utilizing external datasets to augment the existing data might offer more insights.

## 7.3 Parsimony and Model Complexity:

Parsimony, in the context of modeling, refers to the principle that simpler models with fewer variables are preferable if they provide similar predictive power as more complex models. The benefits of parsimonious models are:

- **Interpretability:** Simpler models are easier to understand and explain.
- **Generalizability:** They tend to generalize better to new, unseen data.
- **Reduced Overfitting:** Fewer variables mean lesser chances of fitting to noise.

However, the goal should always be a balance between simplicity and accuracy. While we should strive for parsimony, we should not oversimplify to the point where we lose essential predictive information.

## 7.4 Max Fit Model vs. Simpler Model:

In many real-world scenarios, interpretability is as crucial, if not more so, than raw predictive power. Especially in sectors like healthcare, finance, and public policy, being able to explain why a model makes a particular prediction can be essential for trust and decision-making.

That said, the choice between a max fit model and a simpler model depends on the objective:

- If the primary goal is prediction, and the model's inner workings are less relevant, a max fit model might be more appropriate.
- However, if the goal is understanding relationships between variables or making decisions based on model outputs where stakeholders require explanations, a simpler but interpretable model might be better. I often have to make this decision at work when conforming to regulatory standards.

In conclusion, the journey through this dataset has highlighted the complexities and nuances of predictive modeling. While automated approaches provide a good starting point, human judgment, domain knowledge, and iterative refinement are critical to building robust and useful models.