

SOUTHERN METHODIST UNIVERSITY
MSDS 6371(401)

Kaggle Project

Michael J Wolfe
Sandesh Ojha
Carl Walenciak
Travis Daun

Github Repository

https://github.com/mjwolfe91/SFDS_401_Team3_Kaggle_Project

February 23, 2019

CONTENTS

1. Introduction	3
2. Data Description	3
3. Analysis Question 1	4
3.1. Restatement of Problem	4
3.2. Build and Fit the Model	4
3.3. Checking Assumptions	4
3.3.1. Assumptions	4
3.3.2. Outliers and Influential Points	5
3.3.3. Effect by Neighborhood	6
3.4. Model Metrics	6
3.4.1. Adj R^2	7
3.4.2. Internal Press	7
3.5. Parameters	7
3.5.1. Estimates	7
3.5.2. Interpretation	8
3.5.3. Confidence Intervals	8
3.6. Conclusion	8
4. Analysis Question 2	8
4.1. Restatement of Problem	8
4.2. Model Selection	9
4.2.1. Stepwise	9
4.2.2. Forward	9
4.2.3. Backward	9
4.2.4. Custom	9
4.3. Checking Assumptions	9
4.3.1. Residual Plots	9
4.3.2. Influential point analysis (Cook's D and Leverage)	9
4.3.3. Make sure to address each assumption	9
4.4. Comparing Competing Models	9
4.4.1. Adj R^2	9
4.4.2. Internal CV Press	9
4.4.3. Kaggle Score	9
4.5. Conclusion	9
A. Source Code for Analysis 1	10
B. Source Code for Analysis 2	16
B.1. Forward Selection	16
B.2. Backward Selection	18
B.3. Stepwise Selection	20
C. High Level Summary of Data	24
D. Detailed Data Description	25

1. INTRODUCTION

Many factors can impact the sale price of residential real estate. This report will explore those various aspects to help define what factors tend to impact home prices. We start by looking at the three distinct neighborhoods (North Ames, Edwards, and Brookside) in Ames, Iowa that Century 21 Ames sells houses to help better understand how above ground living space relates to sales price in each of these neighborhoods and if the neighborhood has an impact on the sales price. At the conclusion of this analysis, we provide a conclusion that quantifies the relationship between living area and sales price with respect to these neighborhoods.

After examining the impact of above ground living space on sales price for these three neighborhoods, we provide a much more complex analysis of factors that can be used to determine sales price for Ames as a whole. While we will use various methods to determine what factors are important for our predictive model, we will test this model on a blind dataset to show how it performs. The results of this model on the test data will be provided to further build the level of confidence in this tool for predicting future home sale prices.

2. DATA DESCRIPTION

The dataset used for this analysis was retrieved from www.kaggle.com/c/house-prices-advanced-regression-techniques.

DATA The data for this evaluation contained 79 explanatory variables describing (almost) every aspect of residential homes in Ames, Iowa. A complete description of these data can be found in Appendix D. The explanatory variables contain both categorical and numeric attributes. Appendix C.1, provides a high level summary of the variables and variable types contained in this dataset.

PREPROCESSING The variable `LotFrontage` posed a challenge since it was a continuous numerical variable that contained NA. This is because the variable was for the linear feet of street connected to property. In many cases (259 of 1460) this was either unknown or unrecorded. Our team made the decision to convert these NA values to 0. We believe this is an acceptable practice since we performed a sensitivity analysis by replacing the NAs with values from 0 to 140 (mean value is 70) with no impact on the linear regression model selection process and this factor was not utilized in any of our final models.

TRAINING Training of the linear regression models will be done utilizing the `training.csv` data obtained from above. A five fold cross validation will be employed for model selection.

TESTING Testing of the linear regression models will be done utilizing the `test.csv` data obtained from above.

RESULTS Datafiles containing the test results of the linear regression models can be found in our Github repository at https://github.com/mjwolfe91/SFDS_401_Team3_Kaggle_Project

3. ANALYSIS QUESTION 1

3.1. Restatement of Problem

REQUEST Century 21 Ames only sells houses in the North Ames, Edwards and Brookside neighborhoods and wants an estimate of how the sale price of the house is related to the square footage of the living area of the house. Additionally, Century 21 Ames would like to know if the sale price (and its relationship to square footage) depends on which neighborhood the house is located in. A fit a model will be used to answer this question.

DELIVERABLE Provide the estimate (or estimates if it varies by neighborhood) as well as confidence intervals for any estimate(s) you provide. Provide evidence that the model assumptions are met and that any suspicious observations (outliers / influential observations) have been identified and addressed. Finally, a conclusion that quantifies the relationship between living area and sale price with respect to these three neighborhoods.

3.2. Build and Fit the Model

Appendix A. Source Code For Analysis 1 contains the SAS code used to check the assumptions, clean the data, and run the model to determine the best estimate for sale price based on square footage for the North Ames, Edwards, and Brookside neighborhoods. We removed four data points from the original dataset due to being outliers. First sale prices greater than \$300,000 were determined to not be representative of the total population of these three neighborhoods. Second, sale condition was limited to only those that were normal sales. Again we feel that home sales that were not normal sales such as foreclosures, linked properties, and land purchases were not representative of the population of interest.



Figure 3.1: Simple Linear Regression Models.

3.3. Checking Assumptions

3.3.1. Assumptions

LINEARITY The linearity assumption is met by the reviewing the scatter plots associated with data. Fig 3.1 shows a plot of SalePrice vs. GrLivArea and by removing the outliers the linearity assumption is reasonably met.

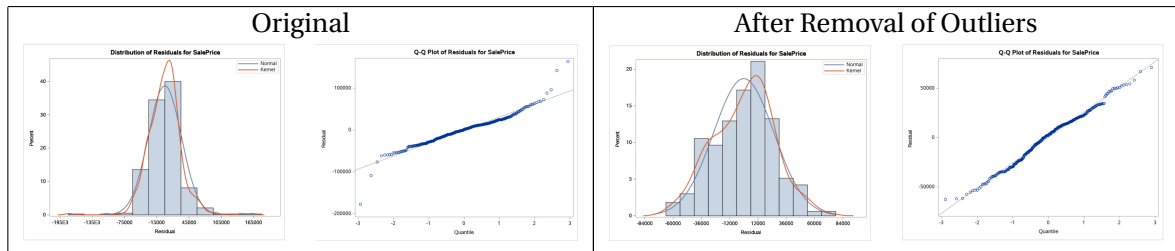


Figure 3.2: Normality Plots.

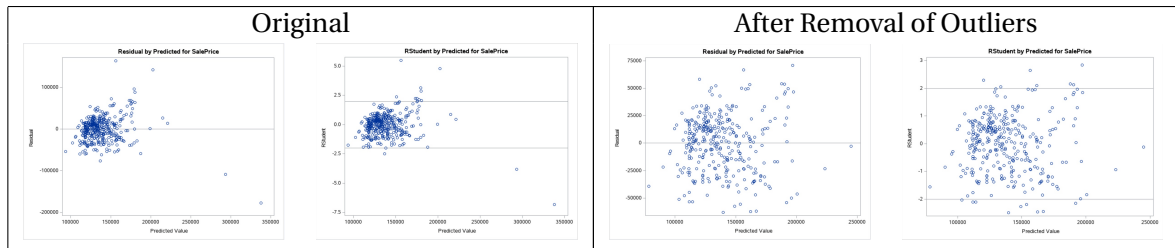


Figure 3.3: Residual Plots.

CONSTANT VARIANCE The residual plot, Fig 3.3 resembles somewhat of a random scatter of points around the 0 line, although there is a slight suspicion of non-constant variance judging from the dense cloud around the predicted value of \$130,000. Also shown is the Studentized Residual Plot which is very similar to the residual plot, although this plot identifies potential outlying observations.

NORMALITY Based upon the histograms and q-q plots in Fig 3.2 there is no evidence to suggest that normality of the data. Additionally the random scatter associated with the residual plots in Fig 3.3 also support the normality assumption.

INDEPENDENCE The independence assumption can be assumed to be maintained since these are all unique sales in a free housing market.

3.3.2. Outliers and Influential Points

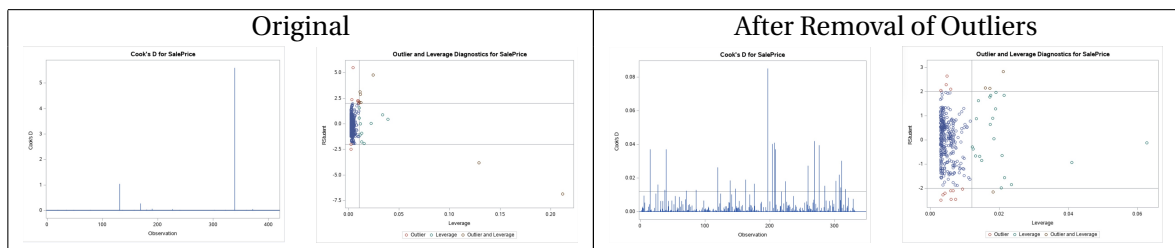


Figure 3.4: Influential Point Plots.

OUTLIERS / INFLUENTIAL OBSERVATIONS There are distinct outliers that can be seen both in the Cook's D and Outlier and Leverage Diagnostics seen in Fig 3.4. By removing the observations that resulted from non-normal sales conditions such as foreclosures and sale prices that were significantly outside the population, the Cook's D and Outlier and Leverage Diagnostics

significantly improved. We are confident that the removal of these observations was appropriate since they do not represent the population as a whole. Additionally we are confident that the remaining data does not contain significant outliers or leverage points that need to be addressed further.

The model is a reasonable fit without transformations. The removal of observations not reflective of the population, such as non-normal sale conditions and home sales greater than \$300,000, seems appropriate and enables the required model assumptions to be met.

3.3.3. Effect by Neighborhood

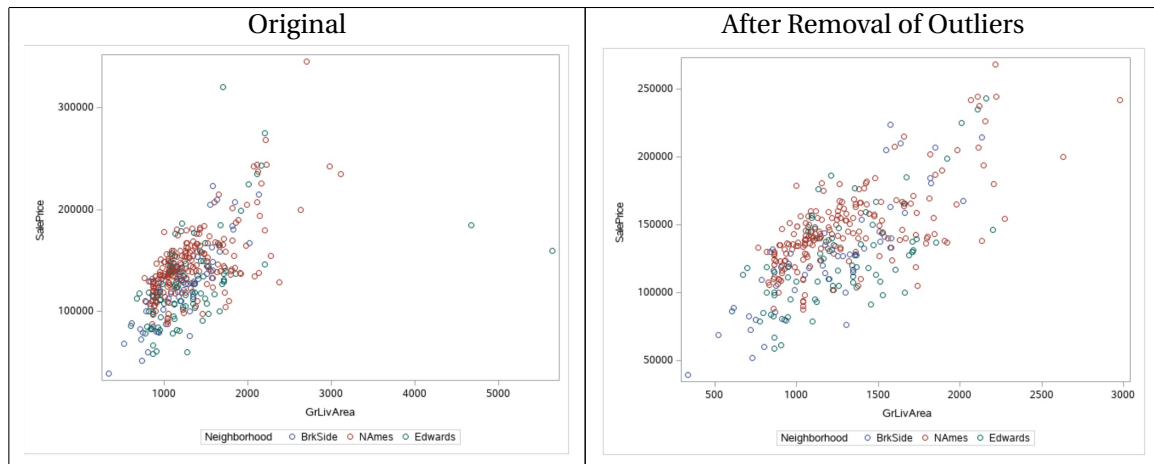


Figure 3.5: Scatterplot of Sale Prices vs Living Area by Neighborhood.

MODEL A model was developed to individually look at the neighborhoods of interest to determine if the sale price is impacted by neighborhood. Fig 3.5 shows the sales price vs living area by neighborhood. The model shows (Table 3.1) that the intercept and slope for Brookside and Edwards do differ from North Ames with statistical significance (p-values: < .0001 and .0077). As such we have determined to choose this model which allows for different intercept and slopes based upon the neighborhood. The resultant model can be seen in Fig 3.6.

Parameter	Estimate		Standard Error	t Value	Pr > t	95% Confidence Limits	
Intercept	74982.03784	B	5925.64914	12.65	<.0001	63324.70074	86639.37493
GrLivArea100	54.93056	B	4.36399	12.59	<.0001	46.34542	63.51570
Neighborhood BrkSide	-55776.16653	B	12012.18226	-4.64	<.0001	-79407.34256	-32144.99051
Neighborhood Edwards	-30273.99848	B	11296.53114	-2.68	0.0077	-52497.29731	-8050.69965
Neighborhood NAmes	0.00000	B
GrLivArea*Neighborhood BrkSide	33.31303	B	9.33444	3.57	0.0004	14.94970	51.67637
GrLivArea*Neighborhood Edwards	8.24334	B	8.45713	0.97	0.3304	-8.39409	24.88077
GrLivArea*Neighborhood NAmes	0.00000	B

Table 3.1: Results of Neighborhood Impact on Sales Price

3.4. Model Metrics

Model metrics can be seen in Table 3.2.

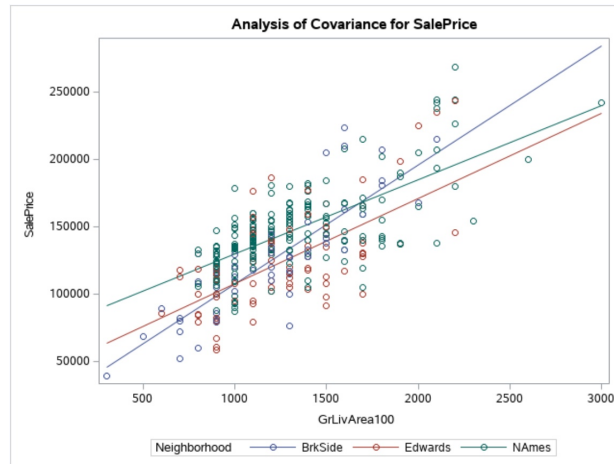


Figure 3.6: Linear Regression Model by Neighborhood.

Root MSE	24094
Dependent Mean	137977
R-Square	0.5333
Adj R-Sq	0.5291
AIC	7037.53791
AICC	7037.72196
BIC	6705.63517
C(p)	4.00000
PRESS	1.961216E11
SBC	6718.75845
ASE	573508240

Table 3.2: Results of Neighborhood Impact on Sales Price

3.4.1. $Adj R^2$

ADJUSTED R^2 obtained for this model is 0.53.

3.4.2. Internal Press

PRESS obtained by this model is $1.96E11$.

3.5. Parameters

3.5.1. Estimates

The parameter estimates can be seen in Table 3.1. With these estimates a separate equation can be written for each neighborhood to predict sale price based on living area.

$$SalePrice = 74982 + 54.93(GrLivArea) - 55776(BrkSide) - 30274(Edwards) + 33.31(GrLivArea)(BrkSide) + 8.24(GrLivArea)(Edwards)$$

NORTH AMES

$$SalePrice = 74982 + 54.93(GrLivArea)$$

EDWARDS

$$\begin{aligned} \text{SalePrice} &= 74982 + 54.93(\text{GrLivArea}) - 30274 + 8.24(\text{GrLivArea}) \\ &= 44708 + 63.17(\text{GrLivArea}) \end{aligned}$$

BROOKSIDE

$$\begin{aligned} \text{SalePrice} &= 74982 + 54.93(\text{GrLivArea}) - 55776 + 33.31(\text{GrLivArea}) \\ &= 19206 + 88.24(\text{GrLivArea}) \end{aligned}$$

3.5.2. Interpretation

β_0 The intercept in this model provides an estimate (74982) of the sale price of a home in North Ames (reference neighborhood) with a living area of zero. Of course, this is extrapolation and does not have a clear, practical meaning.

β_1 For each 100 square foot increase in the living area of a home in North Ames, the estimated sale price increases \$54.93.

β_2 This is the adjustment of the intercept for a home in Brookside with respect to a home in North Ames. For a living area of zero, the home in Brookside has an estimated sale price of \$55,776 less than a home in North Ames.

β_3 This is the adjustment of the intercept for a home in Edwards with respect to a home in North Ames. For a living area of zero, the home in Edwards has an estimated sale price of \$30,274 less than a home in North Ames.

β_4 For each 100 square foot increase in the living area of a home in Brookside, the estimated sale price increases \$33 from the change with the home in North Ames.

β_5 For each 100 square foot increase in the living area of a home in Edwards, the estimated sale price increases \$8 from the change with the home in North Ames.

3.5.3. Confidence Intervals

The confidence intervals for the estimates can be seen in Table 3.1.

3.6. Conclusion

The square feet of above ground living area is a statistically significant feature to use to predict the sale price of homes in the North Ames, Edwards, and Brookside neighborhoods of Ames, Iowa. The existing sale prices of homes that have sold in these neighborhoods that are under \$300,000 and underwent a normal sales condition meet all the assumptions required to generate an appropriate linear regression model. We also examined the differences in predicted prices in these three neighborhoods and determined that the sale prices of homes in each of these neighborhoods do differ from each other. Homes in North Ames under 1600 square feet are predicted to sell for the highest price as compared to the other two neighborhoods. The price per square foot in Brookside increases at the highest rate of the three neighborhoods.

Homes greater than approximately 1000 square feet in Brookside are predicted to sell for higher prices than comparable homes in Edwards. Homes greater than approximately 1600 square feet in Brookside are predicted to sell for higher prices than comparable homes in North Ames. Homes in the Edwards neighborhood are predicted to sell for the lowest price of all three neighborhoods with the exception of homes that are smaller than 1000 square feet.

4. ANALYSIS QUESTION 2

4.1. Restatement of Problem

Build the most predictive model for sales prices of homes in all of Ames Iowa. This includes all neighborhoods. Your group is limited to only the techniques we have learned in 6371 (no random forests or other methods we have not yet covered). Specifically, you should produce 4 models: one from forward selection, one from backwards elimination, one from stepwise selection, and one that you build custom. The custom model could be one of the three preceding models or one that you build by adding or subtracting variables at your will. Generate an adjusted R², CV Press and Kaggle Score for each of these models and clearly describe which model you feel is the best in terms of being able to predict future sale prices of homes in Ames, Iowa.

4.2. Model Selection

4.2.1. *Stepwise*

4.2.2. *Forward*

4.2.3. *Backward*

4.2.4. *Custom*

4.3. Checking Assumptions

4.3.1. *Residual Plots*

4.3.2. *Influential point analysis (Cook's D and Leverage)*

4.3.3. *Make sure to address each assumption*

4.4. Comparing Competing Models

Predictive Models	Adjusted R ²	CV Press	Kaggle Score
Forward	XX	XX	XX
Backward	XX	XX	XX
Stepwise	XX	XX	XX
CUSTOM	XX	XX	XX

Table 4.1: Analysis Results

4.4.1. *Adj R²*

4.4.2. *Internal CV Press*

4.4.3. *Kaggle Score*

4.5. Conclusion

A. SOURCE CODE FOR ANALYSIS 1

Listing 1: Analysis 1 SAS Code.

```
1  *-----*
2  | Import train.csv |
3  | Set REFFILE for train.csv |
4  *-----* ;
5
6  * FILENAME REFFILE '/home/mwolfe0/train.csv';
7  FILENAME REFFILE
8  '/folders/myfolders/MSDS6371/GroupProject/Datasets/train.csv';
9
10 PROC IMPORT DATAFILE=REFFILE DBMS=CSV REPLACE OUT=TRAIN;
11     GETNAMES=YES;
12 RUN;
13
14 *-----*
15 | Subset the data to only include homes sold in the |
16 | neighborhoods of interest – Names, BrkSide, and Edwards |
17 | Round the gross living area to the nearest 100 SF |
18 | Keep only the variables of Neighborhood, GrLivArea, |
19 | and SalePrice in the dataset |
20 *-----* ;
21
22 DATA HOMES1;
23 SET TRAIN (KEEP=Neighborhood GrLivArea SalePrice);
24 IF Neighborhood EQ "Names" |
25     Neighborhood EQ "BrkSide" |
26     Neighborhood EQ "Edwards";
27 GrLivArea100 = ROUND(GrLivArea, 100); /*FLOOR(GrLivArea);*/
28 RUN;
29
30 *-----*
31 | Descriptive statistics on the HOMES1 dataset for |
32 | GrLivArea and SalePrice |
33 *-----* ;
34
35 PROC UNIVARIATE DATA=HOMES1;
36     CLASS Neighborhood;
37     VAR GrLivArea SalePrice;
38 RUN;
39
40 *-----*
41 | Scatter plot of sale prices in the three neighborhoods |
42 | vs Gross Living Area |
43 *-----* ;
44
45 PROC SGPLOT DATA=HOMES1;
46     SCATTER X=GrLivArea Y=SalePrice;
47     REG X=GrLivArea Y=SalePrice;
```

```

48 RUN;
49
50 *-----*
51 | Regression model of homes in the three neighborhoods |
52 | combined for Sale Price based on Gross Living Area |
53 | to check assumptions on the data in these three |
54 | neighborhoods |
55 *-----*;
56
57 PROC REG DATA=HOMES1;
58     MODEL SalePrice=GrLivArea / CLB;
59     RUN;
60
61 *-----*
62 | Regression model of homes in the three neighborhoods |
63 | using an equal slope model |
64 *-----*;
65
66 PROC GLM DATA=HOMES1;
67     CLASS Neighborhood;
68     MODEL SalePrice=GrLivArea Neighborhood / CLPARM;
69     RUN;
70
71 *-----*
72 | Regression model of homes in the three neighborhoods |
73 | using an equal intercept model (slopes differ) |
74 *-----*;
75
76 PROC GLM DATA=HOMES1;
77     CLASS Neighborhood;
78     MODEL SalePrice=GrLivArea*Neighborhood / CLPARM;
79     RUN;
80
81 *-----*
82 | Regression model of homes in the three neighborhoods |
83 | using a model that allows slopes and intercepts to |
84 | vary |
85 *-----*;
86
87 PROC GLM DATA=HOMES1;
88     CLASS Neighborhood;
89     MODEL SalePrice=Neighborhood GrLivArea*Neighborhood / CLPARM;
90     RUN;
91
92 *-----*
93 | Remove Outliers: |
94 | SaleCondition is not normal (confounding effect on |
95 | prices) |
96 | SalePrice is greater than 300,000 since they are not |
97 | representative of the overall population in these |

```

```

98 | three neighborhoods. |
99 | Keep only the variables of Neighborhood, GrLivArea, |
100 | and SalePrice in the dataset |
101 *-----* ;
102
103 DATA HOMES2;
104 SET TRAIN (KEEP=Neighborhood GrLivArea SalePrice SaleCondition);
105 IF Neighborhood EQ "NAmes" |
106     Neighborhood EQ "BrkSide" |
107     Neighborhood EQ "Edwards";
108 IF SalePrice LT 300000;
109 IF SaleCondition EQ "Normal";
110 GrLivArea100 = ROUND(GrLivArea, 100);
111 RUN;
112
113 *-----*
114 | Descriptive statistics on the HOMES1 dataset for |
115 | GrLivArea and SalePrice |
116 *-----* ;
117
118 PROC UNIVARIATE DATA=HOMES2;
119     CLASS Neighborhood;
120     VAR GrLivArea SalePrice;
121 RUN;
122
123 *-----*
124 | Scatter plot of sale prices in the three neighborhoods |
125 | vs Gross Living Area |
126 *-----* ;
127
128 PROC SGPLOT DATA=HOMES2;
129     SCATTER X=GrLivArea Y=SalePrice;
130     REG X=GrLivArea Y=SalePrice;
131 RUN;
132
133 *-----*
134 | Regression model of homes in the three neighborhoods |
135 | using a model that allows slopes and intercepts to |
136 | vary |
137 | Output 95% confidence limit for parameter estimates |
138 *-----* ;
139
140 PROC REG DATA=HOMES2 PLOTS=ALL;
141     MODEL SalePrice=GrLivArea / CLB;
142     RUN;
143
144 *-----*
145 | Regression model of homes in the three neighborhoods |
146 | using an equal slope model |
147 *-----* ;

```

```

148
149 PROC GLM DATA=HOMES2;
150     CLASS Neighborhood;
151     MODEL SalePrice=GrLivArea Neighborhood / CLPARM;
152     RUN;
153
154 *-----*
155 | Regression model of homes in the three neighborhoods |
156 | using an equal intercept model (slopes differ)       |
157 *-----*;
158
159 PROC GLM DATA=HOMES2;
160     CLASS Neighborhood;
161     MODEL SalePrice=GrLivArea*Neighborhood / CLPARM;
162     RUN;
163
164 *-----*
165 | Regression model of homes in the three neighborhoods |
166 | using a model that allows slopes and intercepts to   |
167 | vary                                                  |
168 *-----*;
169
170 PROC GLM DATA=HOMES2;
171     CLASS Neighborhood;
172     MODEL SalePrice=Neighborhood GrLivArea*Neighborhood / CLPARM;
173     RUN;
174
175
176 *-----*
177 | Alternate Method with interaction terms               |
178 |                                                       |
179 | Keep only the variables of Neighborhood, GrLivArea,  |
180 | and SalePrice in the dataset                         |
181 | d1 = NAmes, d2 = BrkSide, Control = Edwards         |
182 *-----*;
183
184 DATA HOMES3;
185 SET TRAIN (KEEP=Neighborhood GrLivArea SalePrice SaleCondition);
186 IF Neighborhood EQ "NAmes" |
187     Neighborhood EQ "BrkSide" |
188     Neighborhood EQ "Edwards";
189 IF SalePrice LT 300000;
190 IF SaleCondition EQ "Normal";
191 GrLivArea100 = ROUND(GrLivArea, 100);
192 IF Neighborhood = 'NAmes' THEN d1 = 1; ELSE d1=0;
193 IF Neighborhood = 'BrkSide' THEN d2 = 1; ELSE d2=0;
194     int1 = d1*GrLivArea100; int2 = d2*GrLivArea100;
195 RUN;
196
197 *-----*

```

```

198 | Plots to check assumptions |
199 | d1 = NAmes, d2 = BrkSide, Control = Edwards |
200 *-----* ;
201
202 PROC SGPLOT DATA=HOMES3;
203 HISTOGRAM GrLivArea100;
204 DENSITY GrLivArea100/TYPE=NORMAL;
205 TITLE "Histogram of Gross Living Area in NAmes, BrkSide, and Edwards";
206 RUN;
207
208 PROC SGPLOT DATA=HOMES3;
209 SCATTER X=GrLivArea100 Y=SalePrice;
210 TITLE "Gross Living Area vs Sale Price in NAmes, BrkSide, and Edwards";
211 RUN;
212
213 PROC REG DATA=HOMES3;
214 model SalePrice = GrLivArea100/CLB;
215 RUN;
216
217 *-----*
218 | Run regression model with interaction terms using dummy|
219 | variables |
220 | d1 = NAmes, d2 = BrkSide, Control = Edwards |
221 | Output 95% confidence limit for parameter estimates |
222 *-----* ;
223 PROC REG DATA=HOMES3;
224     model SalePrice = GrLivArea100 d1 d2 int1 int2/VIF CLB;
225     title
226     'Regression of Sale Price on Gross Living Area
227     with Interaction Terms';
228     RUN;
229
230 *-----*
231 | center the interaction terms based on the means of |
232 | GrLivArea100 and d1 and d2 to correct for the |
233 | inflated VIF |
234 *-----* ;
235
236 PROC MEANS DATA=HOMES3;
237 var GrLivArea100 d1 d2;
238 run;
239
240 DATA center;
241 set Homesp1b;
242 cent1 = (GrLivArea100 - 1280.72)*(d1-0.588);
243 cent2 = (GrLivArea100 - 1280.72)*(d2-0.151);
244 RUN;
245
246 DATA center;
247 set HOMES3;

```

```

248 cent1 = (GrLivArea100 - 1283.2)*(d1-0.593);
249 cent2 = (GrLivArea100 - 1283.2)*(d2-0.164);
250 RUN;
251
252 PROC REG DATA=center PLOTS=ALL;
253 model SalePrice = GrLivArea100 d1 d2 cent1 cent2/VIF CLB;
254 title
255         'Regression of Sale Price on Gross Living Area
256         with Interaction Terms';
257 RUN;
258
259 PROC GLM DATA=HOMES3 PLOT=ALL;
260 CLASS Neighborhood;
261 model SalePrice=GrLivArea100|Neighborhood/solution CLPARM;
262 RUN;

```

B. SOURCE CODE FOR ANALYSIS 2

B.1. Forward Selection

Listing 2: Forward Selection SAS Code.

```
1  *-----*
2  | Import train.csv |
3  | Import test.csv |
4  | Set REFFILE for train.csv |
5  | Set REFFILE2 for test.csv |
6  *-----* ;
7
8  *FILENAME REFFILE '/home/mwolfe0/train.csv';
9  *FILENAME REFFILE2 '/home/mwolfe0/test.csv';
10 FILENAME REFFILE
11 '/folders/myfolders/MSDS6371/GroupProject/Datasets/train.csv';
12 FILENAME REFFILE2
13 '/folders/myfolders/MSDS6371/GroupProject/Datasets/test.csv';
14
15 PROC IMPORT DATAFILE=REFFILE DBMS=CSV REPLACE OUT=TRAIN;
16     GETNAMES=YES;
17 RUN;
18
19 PROC IMPORT DATAFILE=REFFILE2 DBMS=CSV REPLACE OUT=TEST;
20     GETNAMES=YES;
21 RUN;
22
23 *-----*
24 | Combine train and test into one datafile HOMES |
25 *-----* ;
26
27 DATA HOMES;
28     SET TRAIN TEST;
29     IF LotFrontage EQ "NA" THEN LotFrontage = 0;
30     LotFront = input(LotFrontage, 8.);
31     drop LotFrontage;
32     RENAME LotFront=LotFrontage;
33 RUN;
34
35 *-----*
36 | Code for forward selection |
37 | Set seed to a constant for model comparison |
38 | Class variable input with split option to allow |
39 |     classification variable to be able to enter or |
40 |     leave the model independently |
41 | Stop=CV specifies the model will stop when the |
42 |     predicted residual sum of square is reached with |
43 |     k-fold cross validation |
44 | CVMMethod specifies how subsets are formed for |
45 |     cross validation |
```



```

46 | OUTPUT Dataset to RESULTS with the predicted variable |
47 | based on the final model |
48 *-----*
49
50 PROC GLMSELECT DATA=HOMES SEED=71669132;
51     CLASS MSSubClass MSZoning Street Alley LotShape LandContour
52         Utilities LotConfig LandSlope Neighborhood Condition1
53         Condition2 BldgType HouseStyle OverallQual OverallCond
54         RoofMatl Exterior1st Exterior2nd MasVnrArea ExterQual
55         ExterCond Foundation BsmtQual BsmtExposure BsmtFinType1
56         BsmtFinType2 Heating HeatingQC CentralAir Electrical
57         KitchenQual Functional FireplaceQu GarageType
58         GarageFinish GarageQual GarageCond PavedDrive PoolQC
59         Fence MiscFeature SaleType SaleCondition RoofStyle
60         BsmtCond MasVnrType
61         / split;
62     MODEL SalePrice= LotArea YearBuilt YearRemodAdd BsmtFinSF1
63         BsmtFinSF2 BsmtUnfSF TotalBsmtSF _1stFlrSF _2ndFlrSF
64         LowQualFinSF GrLivArea BsmtFullBath BsmtHalfBath
65         FullBath HalfBath BedroomAbvGr KitchenAbvGr
66         TotRmsAbvGrd Fireplaces GarageYrBlt GarageCars
67         GarageArea WoodDeckSF OpenPorchSF EnclosedPorch
68         _3SsnPorch ScreenPorch PoolArea MiscVal MoSold YrSold
69         MSSubClass MSZoning Street Alley LotShape LandContour
70         Utilities LotConfig LandSlope Neighborhood Condition1
71         Condition2 BldgType HouseStyle OverallQual OverallCond
72         RoofMatl Exterior1st Exterior2nd MasVnrArea ExterQual
73         ExterCond Foundation BsmtQual BsmtExposure BsmtFinType1
74         BsmtFinType2 Heating HeatingQC CentralAir Electrical
75         KitchenQual Functional FireplaceQu GarageType
76         GarageFinish GarageQual GarageCond PavedDrive PoolQC
77         Fence MiscFeature SaleType SaleCondition LotFrontage
78         RoofStyle BsmtCond MasVnrType
79     / selection =forward(stop=CV) cvmethod=random(5) stats=all;
80     OUTPUT OUT=RESULTS P=PREDICT;
81 RUN;
82
83 *-----*
84 | Create a datafile RESULTS_FW of predicted values for |
85 | SalePrice for house id greater than 1460 which |
86 | is where the Kaggle test set data begins. |
87 *-----*
88
89 DATA RESULTS_FW;
90     SET RESULTS;
91
92     IF PREDICT > 0 THEN
93         SalePrice=Predict;
94
95     IF PREDICT < 0 THEN

```

```

96         SalePrice=10000;
97     KEEP id SalePrice;
98     WHERE id > 1460;
99 RUN;
100
101 *-----*
102 | Export a datafile for predicted values for      |
103 | SalePrice for house id greater than 1460 which  |
104 | is where the Kaggle test set data begins.      |
105 *-----*;
106
107 *FILENAME REFFILE3 '/home/mwolfe0/results_fw.csv';
108 FILENAME REFFILE3
109 '/folders/myfolders/MSDS6371/GroupProject/Datasets/results_fw.csv';
110
111 PROC EXPORT DATA=RESULTS_FW FILE=REFFILE3 DBMS=CSV REPLACE;
112 RUN;

```

B.2. Backward Selection

Listing 3: Backward Selection SAS Code.

```

1  *-----*
2  | Import train.csv                               |
3  | Import test.csv                               |
4  | Set REFFILE for train.csv                     |
5  | Set REFFILE2 for test.csv                     |
6  *-----*;
7
8  *FILENAME REFFILE '/home/mwolfe0/train.csv';
9  *FILENAME REFFILE2 '/home/mwolfe0/test.csv';
10 FILENAME REFFILE
11 '/folders/myfolders/MSDS6371/GroupProject/Datasets/train.csv';
12 FILENAME REFFILE2
13 '/folders/myfolders/MSDS6371/GroupProject/Datasets/test.csv';
14
15 PROC IMPORT DATAFILE=REFFILE DBMS=CSV REPLACE OUT=TRAIN;
16     GETNAMES=YES;
17 RUN;
18
19 PROC IMPORT DATAFILE=REFFILE2 DBMS=CSV REPLACE OUT=TEST;
20     GETNAMES=YES;
21 RUN;
22
23 *-----*
24 | Combine train and test into one datafile HOMES |
25 *-----*;
26
27 DATA HOMES;
28     SET TRAIN TEST;

```

```

29      IF LotFrontage EQ "NA" THEN LotFrontage = 0;
30      LotFront = input(LotFrontage, 8.);
31      drop LotFrontage;
32      RENAME LotFront=LotFrontage;
33  RUN;
34
35  *-----*
36  | Code for backward selection |
37  | Set seed to a constant for model comparison |
38  | Class variable input with split option to allow |
39  |   classification variable to be able to enter or |
40  |   leave the model independently |
41  | Stop=10 specifies the model will stop selection at the |
42  |   first step for which the selected model has 10 |
43  |   effects |
44  | CVMMethod specifies how subsets are formed for |
45  |   cross validation |
46  | OUTPUT Dataset to RESULTS with the predicted variable |
47  |   based on the final model |
48  *-----*;
49
50  PROC GLMSELECT DATA=HOMES SEED=71669132;
51      CLASS MSSubClass MSZoning Street Alley LotShape LandContour
52          Utilities LotConfig LandSlope Neighborhood Condition1
53          Condition2 BldgType HouseStyle OverallQual OverallCond
54          RoofMatl Exterior1st Exterior2nd MasVnrArea ExterQual
55          ExterCond Foundation BsmtQual BsmtExposure BsmtFinType1
56          BsmtFinType2 Heating HeatingQC CentralAir Electrical
57          KitchenQual Functional FireplaceQu GarageType
58          GarageFinish GarageQual GarageCond PavedDrive PoolQC
59          Fence MiscFeature SaleType SaleCondition RoofStyle
60          BsmtCond MasVnrType
61          / split;
62      MODEL SalePrice= LotArea YearBuilt YearRemodAdd BsmtFinSF1
63          BsmtFinSF2 BsmtUnfSF TotalBsmtSF _1stFlrSF _2ndFlrSF
64          LowQualFinSF GrLivArea BsmtFullBath BsmtHalfBath
65          FullBath HalfBath BedroomAbvGr KitchenAbvGr
66          TotRmsAbvGrd Fireplaces GarageYrBlt GarageCars
67          GarageArea WoodDeckSF OpenPorchSF EnclosedPorch
68          _3SsnPorch ScreenPorch PoolArea MiscVal MoSold YrSold
69          MSSubClass MSZoning Street Alley LotShape LandContour
70          Utilities LotConfig LandSlope Neighborhood Condition1
71          Condition2 BldgType HouseStyle OverallQual OverallCond
72          RoofMatl Exterior1st Exterior2nd MasVnrArea ExterQual
73          ExterCond Foundation BsmtQual BsmtExposure BsmtFinType1
74          BsmtFinType2 Heating HeatingQC CentralAir Electrical
75          KitchenQual Functional FireplaceQu GarageType
76          GarageFinish GarageQual GarageCond PavedDrive PoolQC
77          Fence MiscFeature SaleType SaleCondition LotFrontage
78          RoofStyle BsmtCond MasVnrType

```

```

79      / selection =backward(stop=10) cvmethod=random(5)
80      stats=ADJRSQ stats=PRESS;
81      OUTPUT OUT=RESULTS P=PREDICT;
82  RUN;
83
84  *-----*
85  | Create a datafile RESULTS_BW of predicted values for |
86  | SalePrice for house id greater than 1460 which      |
87  | is where the Kaggle test set data begins.          |
88  *-----*
89
90  DATA RESULTS_BW;
91      SET RESULTS;
92
93      IF PREDICT > 0 THEN
94          SalePrice=Predict;
95
96      IF PREDICT < 0 THEN
97          SalePrice=10000;
98      KEEP id SalePrice;
99      WHERE id > 1460;
100 RUN;
101
102 *-----*
103 | Export a datafile for predicted values for           |
104 | SalePrice for house id greater than 1460 which      |
105 | is where the Kaggle test set data begins.          |
106 *-----*
107
108 *FILENAME REFFILE3 '/home/mwolfe0/results_bw.csv';
109 FILENAME REFFILE3
110 '/folders/myfolders/MSDS6371/GroupProject/Datasets/results_bw.csv';
111
112 PROC EXPORT DATA=RESULTS_BW FILE=REFFILE3 DBMS=CSV REPLACE;
113 RUN;

```

B.3. Stepwise Selection

Listing 4: Stepwise Selection SAS Code.

```

1  *-----*
2  | Import train.csv |
3  | Import test.csv  |
4  | Set REFFILE for train.csv |
5  | Set REFFILE2 for test.csv |
6  *-----*
7
8  *FILENAME REFFILE '/home/mwolfe0/train.csv';
9  *FILENAME REFFILE2 '/home/mwolfe0/test.csv';
10 FILENAME REFFILE

```

```

11  '/folders/myfolders/MSDS6371/GroupProject/Datasets/train.csv';
12  FILENAME REFFILE2
13  '/folders/myfolders/MSDS6371/GroupProject/Datasets/test.csv';
14
15  PROC IMPORT DATAFILE=REFFILE DBMS=CSV REPLACE OUT=TRAIN;
16      GETNAMES=YES;
17  RUN;
18
19  PROC IMPORT DATAFILE=REFFILE2 DBMS=CSV REPLACE OUT=TEST;
20      GETNAMES=YES;
21  RUN;
22
23  *-----*
24  | Combine train and test into one datafile HOMES |
25  *-----*
26
27  DATA HOMES;
28      SET TRAIN TEST;
29      IF LotFrontage EQ "NA" THEN LotFrontage = 0;
30      LotFront = input(LotFrontage, 8.);
31      drop LotFrontage;
32      RENAME LotFront=LotFrontage;
33  RUN;
34
35  *-----*
36  | Code for stepwise selection |
37  | Set seed to a constant for model comparison |
38  | Class variable input with split option to allow |
39  | classification variable to be able to enter or |
40  | leave the model independently |
41  | Stop=CV specifies the model will stop when the |
42  | predicted residual sum of square is reached with |
43  | k-fold cross validation |
44  | CVMMethod specifies how subsets are formed for |
45  | cross validation |
46  | OUTPUT Dataset to RESULTS with the predicted variable |
47  | based on the final model |
48  *-----*
49
50  PROC GLMSELECT DATA=HOMES SEED=71669132;
51      CLASS MSSubClass MSZoning Street Alley LotShape LandContour
52          Utilities LotConfig LandSlope Neighborhood Condition1
53          Condition2 BldgType HouseStyle OverallQual OverallCond
54          RoofMatl Exterior1st Exterior2nd MasVnrArea ExterQual
55          ExterCond Foundation BsmtQual BsmtExposure BsmtFinType1
56          BsmtFinType2 Heating HeatingQC CentralAir Electrical
57          KitchenQual Functional FireplaceQu GarageType
58          GarageFinish GarageQual GarageCond PavedDrive PoolQC
59          Fence MiscFeature SaleType SaleCondition RoofStyle
60          BsmtCond MasVnrType

```

```

61         / split;
62     MODEL SalePrice= LotArea YearBuilt YearRemodAdd BsmtFinSF1
63         BsmtFinSF2 BsmtUnfSF TotalBsmtSF _1stFlrSF _2ndFlrSF
64         LowQualFinSF GrLivArea BsmtFullBath BsmtHalfBath
65         FullBath HalfBath BedroomAbvGr KitchenAbvGr
66         TotRmsAbvGrd Fireplaces GarageYrBlt GarageCars
67         GarageArea WoodDeckSF OpenPorchSF EnclosedPorch
68         _3SsnPorch ScreenPorch PoolArea MiscVal MoSold YrSold
69         MSSubClass MSZoning Street Alley LotShape LandContour
70         Utilities LotConfig LandSlope Neighborhood Condition1
71         Condition2 BldgType HouseStyle OverallQual OverallCond
72         RoofMatl Exterior1st Exterior2nd MasVnrArea ExterQual
73         ExterCond Foundation BsmtQual BsmtExposure BsmtFinType1
74         BsmtFinType2 Heating HeatingQC CentralAir Electrical
75         KitchenQual Functional FireplaceQu GarageType
76         GarageFinish GarageQual GarageCond PavedDrive PoolQC
77         Fence MiscFeature SaleType SaleCondition LotFrontage
78         RoofStyle BsmtCond MasVnrType
79     / selection=stepwise(stop=CV) cvmethod=random(5) stats=all;
80     OUTPUT OUT=RESULTS P=PREDICT;
81 RUN;
82
83     *----------------------------------------------------------------*
84     | Create a datafile RESULTS_SW of predicted values for      |
85     | SalePrice for house id greater than 1460 which            |
86     | is where the Kaggle test set data begins.                |
87     *----------------------------------------------------------------*;
88
89 DATA RESULTS_SW;
90     SET RESULTS;
91
92     IF PREDICT > 0 THEN
93         SalePrice=Predict;
94
95     IF PREDICT < 0 THEN
96         SalePrice=10000;
97     KEEP id SalePrice;
98     WHERE id > 1460;
99 RUN;
100
101     *----------------------------------------------------------------*
102     | Export a datafile for predicted values for                |
103     | SalePrice for house id greater than 1460 which            |
104     | is where the Kaggle test set data begins.                |
105     *----------------------------------------------------------------*;
106
107 *FILENAME REFFILE3 '/home/mwolfe0/results_sw.csv';
108 FILENAME REFFILE3
109 '/folders/myfolders/MSDS6371/GroupProject/Datasets/results_sw.csv';
110

```

```
111 PROC EXPORT DATA=RESULTS_SW FILE=REFFILE3 DBMS=CSV REPLACE;  
112 RUN;
```

C. HIGH LEVEL SUMMARY OF DATA

Attribute Type		Description	Features
Categorical (Qualitative)	Nominal	Only provide enough information to distinguish one object from another	MSSubClass, MSZoning, Street, Alley, LotShape, LandContour, Utilities, LotConfig, LandSlope, Neighborhood, Condition1, Condition2, BldgType, HouseStyle, RoofMatl, Exterior1st, Exterior2nd, MasVnrType, Foundation, BsmtExposure, HeatingCentralAir, Electrical, GarageType, GarageFinish, PavedDrive, MiscFeature, SaleType, SaleCondition, RoofStyle (30 variables)
	Ordinal	Provide enough information to order objects	OverallQual, OverallCond, ExterQual, ExterCond, BsmtQual, BsmtCond, BsmtFinType1, BsmtFinType2, HeatingQC, KitchenQual, Functional, FireplaceQu, GarageQual, GarageCond, PoolQC, Fence (16 variables)
Numeric (Quantitative)	Interval	Interval attributes difference between values are meaningful	YearBuilt, YearRemodAdd, GarageYrBlt, MoSold, YrSold (5 variables)
	Ratio	Differences and ratios are meaningful	LotFrontage, LotArea, BsmtFinSF1, BsmtFinSF2, BsmtUnfSF, TotalBsmtSF, 1stFlrSF, 2ndFlrSF, LowQualFinSF, GrLivArea, BsmtFullBath, MasVnrArea, BsmtHalfBath, FullBath, HalfBath, Bedroom, Kitchen, TotRmsAbvGrd, Fireplaces, GarageCars, GarageArea, WoodDeckSF, OpenPorchSF, EnclosedPorch, 3SsnPorch, ScreenPorch, PoolArea, MiscVal (28 variables)

Table C.1: Summary of Dataset

D. DETAILED DATA DESCRIPTION

Classification Variables	
MSSubClass	Identifies the type of dwelling involved in the sale
20	1-STORY 1946 & NEWER ALL STYLES
30	1-STORY 1945 & OLDER
40	1-STORY W/FINISHED ATTIC ALL AGES
45	1-1/2 STORY - UNFINISHED ALL AGES
50	1-1/2 STORY FINISHED ALL AGES
60	2-STORY 1946 & NEWER
70	2-STORY 1945 & OLDER
75	2-1/2 STORY ALL AGES
80	SPLIT OR MULTI-LEVEL
85	SPLIT FOYER
90	DUPLEX - ALL STYLES AND AGES
120	1-STORY PUD (Planned Unit Development) - 1946 & NEWER
150	1-1/2 STORY PUD - ALL AGES
160	2-STORY PUD - 1946 & NEWER
180	PUD - MULTILEVEL - INCL SPLIT LEV/FOYER
190	2 FAMILY CONVERSION - ALL STYLES AND AGES