

Assignment #5

Nate Bitting

Introduction

The objective of this assignment is to perform automated variable selection techniques for identifying the “best” regression model for predicting sale price for homes in the Ames, Iowa area. The first phase includes the assessment of which predictor variables, based on common sense and business justification, made sense to include in a predictive model. After conducting some preliminary exploratory data analysis (EDA), it was concluded that the following predictor variable candidates would be considered in the model:

- **X1:** YearBuilt – year the home was built
- **X2:** total_SF – total square footage in the home
- **X3:** total_baths – total number of bathrooms in the home
- **X4:** good_kitchen – indicator variable to determine the condition of the kitchen
- **X5:** good_fireplace – indicator variable to determine the condition of the fireplace
- **X6:** good_exterior – indicator variable to determine the condition of the exterior of the home
- **X7:** quality_index – measure of Overall Condition * Overall Quality
- **X8:** central_air – indicator variable to determine if the home has central air or not
- **X9:** fireplace_ind – indicator variable to determine if there were any fireplaces or not
- **X10:** garage_ind – indicator variable to determine if there was a garage or not
- **X11:** good_basement_ind – indicator variable to determine the condition of the basement

After we assess which of the eleven predictor variables should be included in our model, we then will conduct an assessment of extreme observations, or outliers. The approach we will use for outlier detection is by outputting the studentized residuals for a model that includes all eleven variables. We will then remove any observation that has an absolute value of studentized residual exceeding 2. The next step will be to create a training and test dataset that could be used for building the model. The training set will be used to build the regression models and the test set will be used to test the predictive accuracy of the final selected model. Lastly, we will assign a prediction grade to each observation based on how accurate our predictions on the test set compare against the observed values. This will be the basis of our final assessment of how well the model will perform in a real-world setting using new observations.

Results

Model Identification by Automated Variable Selection

After performing our preliminary EDA and outlier removal process, the next step was to perform various automated variable selection methods in order to identify the “best” regression model. The automated variable selection methods used on this analysis included the following:

- Adjusted R-squared (Model_AdjR2)
- Maximum R-squared (Model_MaxR)
- Mallow’s Cp (Model_McP)
- Forward selection (Model_F)
- Backward selection (Model_B)
- Stepwise selection (Model_S)

In the next section we will go through each variable selection method and assess the final model and discuss each step in the variable selection process and provide an assessment of the results.

Adjusted R-Squared Model (Model_AdjR2)

The results of the adjusted R-squared variable selection method determined that all predictor variable candidates (X1-X11) should be included in the final model. The adjusted R-squared of the final model was 0.8887, which was far superior when a subset of the predictor variables were included in the model. In Table 1 below, you can see the summary output of the variable selection process with the final result in the first row.

Table 1 – Adjusted R-squared variable selection summary

Number in Model	Adjusted R-Square	R-Square	Variables in Model
11	0.8887	0.8898	YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior quality_index central_air fireplace_ind garage_ind good_basement_ind
10	0.8885	0.8894	YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior quality_index central_air fireplace_ind good_basement_ind
10	0.8885	0.8894	YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior quality_index central_air fireplace_ind garage_ind
9	0.8883	0.8891	YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior quality_index central_air fireplace_ind
10	0.8879	0.8888	YearBuilt total_SF total_baths good_fireplace good_exterior quality_index central_air fireplace_ind garage_ind good_basement_ind
10	0.8877	0.8887	YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior quality_index fireplace_ind garage_ind good_basement_ind
10	0.8877	0.8887	YearBuilt total_SF total_baths good_kitchen good_exterior quality_index central_air fireplace_ind garage_ind good_basement_ind
9	0.8877	0.8886	YearBuilt total_SF total_baths good_fireplace good_exterior quality_index central_air fireplace_ind good_basement_ind
9	0.8877	0.8885	YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior quality_index fireplace_ind good_basement_ind
9	0.8876	0.8884	YearBuilt total_SF total_baths good_fireplace good_exterior quality_index central_air fireplace_ind garage_ind
9	0.8875	0.8884	YearBuilt total_SF total_baths good_kitchen good_exterior quality_index central_air fireplace_ind garage_ind
9	0.8875	0.8883	YearBuilt total_SF total_baths good_kitchen good_exterior quality_index central_air fireplace_ind good_basement_ind
8	0.8875	0.8882	YearBuilt total_SF total_baths good_fireplace good_exterior quality_index central_air fireplace_ind

NOTE: Not all models included due to limited space

Based on these results, the adjusted R-squared method recommends we include all candidate predictor variables we originally considered for the model (X1-X11) to yield the highest possible adjusted R-squared value.

Maximum R-Squared Model (Model_MaxR)

Just like the adjusted R-squared variable selection method, the results of the maximum R-squared variable selection method also determined that all predictor variable candidates (X1-X11) should be included in the final model. The R-squared of the final model was 0.8898, which was far superior when a subset of the predictor variables were included in

the model. In Table 2 below, we can see the ANOVA and parameter estimates for the final suggested model from the maximum R-squared variable selection method.

Table 2 – Maximum R-squared variable selection summary

Maximum R-Square Improvement: Step 11					
Variable garage_ind Entered: R-Square = 0.8898 and C(p) = 12.0000					
Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	11	3.708008E12	3.370915E11	850.35	<.0001
Error	1159	4.594433E11	396413516		
Corrected Total	1170	4.167449E12			

Variable	Parameter Estimate	Standard Error	Type II SS	F Value	Pr > F
Intercept	-898145	79758	50267949163	126.81	<.0001
YearBuilt	448.06847	41.11992	47068764087	118.74	<.0001
total_SF	47.94454	1.24588	5.870452E11	1480.89	<.0001
total_baths	7137.00080	1098.01401	16748048368	42.25	<.0001
good_kitchen	5114.75146	1628.58768	3909986365	9.86	0.0017
good_fireplace	5587.92777	1662.48892	4478499220	11.30	0.0008
good_exterior	10620	1916.40496	12172904617	30.71	<.0001
quality_index	1441.45165	82.73133	1.203396E11	303.57	<.0001
central_air	-11873	3541.64889	4455031344	11.24	0.0008
fireplace_ind	8325.14250	1504.36520	12140183362	30.63	<.0001
garage_ind	7249.84653	3883.73858	1381357538	3.48	0.0622
good_basement_ind	3557.07044	1869.02846	1435825141	3.62	0.0573

Bounds on condition number: 3.1103, 228.45

The above model is the best 11-variable model found.

No further improvement in R-Square is possible.

Based on these results, the maximum R-squared method recommends we include all candidate predictor variables we originally considered for the model (X1-X11) that would yield the highest R-squared value.

Mallow's Cp Model (Model_MCp)

As with the previous two models, the results of the Mallow's Cp variable selection method also determined that all predictor variable candidates (X1-X11) should be included in the final model. The Mallow's Cp of the final model was 12.0, which was the smallest result from all other subset model alternatives. In Table 3 below, we can see some of the output from the Mallow's Cp variable selection process. The best model is shown in the first row of the table.

Table 3 – Mallow's Cp variable selection summary

Number in Model	C(p)	R-Square	Variables in Model
11	12.0000	0.8898	YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior quality_index central_air fireplace_ind garage_ind good_basement_ind
10	13.4846	0.8894	YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior quality_index central_air fireplace_ind good_basement_ind
10	13.6220	0.8894	YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior quality_index central_air fireplace_ind garage_ind
9	14.7703	0.8891	YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior quality_index central_air fireplace_ind
10	19.8634	0.8888	YearBuilt total_SF total_baths good_fireplace good_exterior quality_index central_air fireplace_ind garage_ind good_basement_ind
9	20.5344	0.8886	YearBuilt total_SF total_baths good_fireplace good_exterior quality_index central_air fireplace_ind good_basement_ind
9	20.9530	0.8885	YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior quality_index fireplace_ind good_basement_ind
10	21.2383	0.8887	YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior quality_index fireplace_ind garage_ind good_basement_ind
10	21.2975	0.8887	YearBuilt total_SF total_baths good_kitchen good_exterior quality_index central_air fireplace_ind garage_ind good_basement_ind
9	21.7200	0.8884	YearBuilt total_SF total_baths good_fireplace good_exterior quality_index central_air fireplace_ind garage_ind
8	22.0781	0.8882	YearBuilt total_SF total_baths good_fireplace good_exterior quality_index central_air fireplace_ind
9	22.6109	0.8884	YearBuilt total_SF total_baths good_kitchen good_exterior quality_index central_air fireplace_ind garage_ind
9	22.7771	0.8883	YearBuilt total_SF total_baths good_kitchen good_exterior quality_index central_air fireplace_ind good_basement_ind
8	22.9991	0.8881	YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior quality_index fireplace_ind

Based on these results, the Mallow's Cp method recommends we include all candidate predictor variables we originally considered for the model (X1-X11) that would yield the lowest Mallow's Cp value.

Forward Selection Model (Model_F)

As with all proceeding variable selection methods, the results of the forward variable selection method also determined that all predictor variable candidates (X1-X11) should be included in the final model. In Table 4 below, we can see the summary of the forward selection method and you'll notice that the p-value from the nested F-tests did not increase until the 7th variable was entered into the model. For this method, we chose a *s/entry* value of 0.15 as our threshold for variables to be allowed to enter into the model.

Table 4 – Forward variable selection summary

Summary of Forward Selection							
Step	Variable Entered	Number Vars In	Partial R-Square	Model R-Square	C(p)	F Value	Pr > F
1	total_SF	1	0.7805	0.7805	1350.89	3711.90	<.0001
2	YearBuilt	2	0.0571	0.8176	752.324	365.85	<.0001
3	quality_index	3	0.0498	0.8675	230.301	438.91	<.0001
4	good_exterior	4	0.0073	0.8748	155.436	68.08	<.0001
5	fireplace_ind	5	0.0064	0.8812	89.7125	63.18	<.0001
6	total_baths	6	0.0048	0.8880	41.7044	48.56	<.0001
7	good_fireplace	7	0.0011	0.8871	31.6481	11.82	0.0006
8	central_air	8	0.0011	0.8882	22.0781	11.44	0.0007
9	good_kitchen	9	0.0009	0.8891	14.7703	9.27	0.0024
10	good_basement_ind	10	0.0003	0.8894	13.4846	3.28	0.0704
11	garage_ind	11	0.0003	0.8898	12.0000	3.48	0.0622

Based on these results, the forward selection method recommends we include all candidate predictor variables we originally considered for the model (X1-X11).

Backward Selection Model (Model_F)

The backward variable selection method actually did not iterate through any removal steps as all predictor variable were kept in the model. This result is exactly the same as all previous selection methods. A p-value of 0.15 was used for the *s/stay* option, which resulted in none of the variables being removed from the model.

Stepwise Selection Model (Model_S)

The stepwise selection method was the final option used for model selection. As with all previous methods, the stepwise selection method indicated that all predictor variable candidates should remain in the model. The stepwise variable selection summary is shown in Table 5 below. You can see that all variables were entered into the model at each step with the Model R-Square increasing and Cp decreasing with the addition of each new variable. The final "best" model from this method indicates that all variables should remain in the model.

Table 5 – Forward variable selection summary

Summary of Stepwise Selection								
Step	Variable Entered	Variable Removed	Number Vars In	Partial R-Square	Model R-Square	C(p)	F Value	Pr > F
1	total_SF		1	0.7805	0.7805	1350.89	3711.90	<.0001
2	YearBuilt		2	0.0571	0.8176	752.324	385.85	<.0001
3	quality_index		3	0.0498	0.8675	230.301	438.91	<.0001
4	good_exterior		4	0.0073	0.8748	155.436	68.08	<.0001
5	fireplace_ind		5	0.0084	0.8812	89.7125	63.18	<.0001
6	total_baths		6	0.0048	0.8860	41.7044	48.56	<.0001
7	good_fireplace		7	0.0011	0.8871	31.6481	11.82	0.0006
8	central_air		8	0.0011	0.8882	22.0781	11.44	0.0007
9	good_kitchen		9	0.0009	0.8891	14.7703	9.27	0.0024
10	good_basement_ind		10	0.0003	0.8894	13.4846	3.28	0.0704
11	garage_ind		11	0.0003	0.8898	12.0000	3.48	0.0622

Model Comparison

Given all variable selection methods resulted in the same model that included all eleven predictor variable candidates, the model fit criteria is exactly the same for all models. Table 6 below shows the model fit criteria from the models created using the training sample. In the next couple sections we will explore if any multicollinearity exists in the model as well as the operational accuracy of the final model that includes all predictor variables using the test dataset (out-of-sample).

Table 6 – Model Comparison from Training and Test samples

		Model_AdjR2	Model_MaxR	Model_MCp	Model_F	Model_B	Model_S
Training Sample	Predictor(s) Selected	X1 - X11	X1 - X11	X1 - X11	X1 - X11	X1 - X11	X1 - X11
	Adjusted R2	0.8887	0.8887	0.8887	0.8887	0.8887	0.8887
	AIC	23195.3592	23195.3592	23195.3592	23195.3592	23195.3592	23195.3592
	BIC	23197.6074	23197.6074	23197.6074	23197.6074	23197.6074	23197.6074
	Mallow's Cp	12.0000	12.0000	12.0000	12.0000	12.0000	12.0000
	MSE	396413516	396413516	396413516	396413516	396413516	396413516
	MAE	15372.37	15372.37	15372.37	15372.37	15372.37	15372.37
Test Sample	MSE	402854448	402854448	402854448	402854448	402854448	402854448
	MAE	15488.47	15488.47	15488.47	15488.47	15488.47	15488.47

Based on the results in Table 6, we can see that the predictive ability of the final model built from the training sample performed very well with the test sample data.

Multicollinearity Assessment

Given the end result for all variable selection methods determined all predictor variable candidates should be included in the final model, we performed an assessment for multicollinearity on the final model only. To assess multicollinearity we leverage the variable inflation factor (VIF) statistic as shown in Table 7 below.

Table 7 – Parameter Estimates with Variable Inflation Factors (VIF)

Parameter Estimates						
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Variance Inflation
Intercept	1	-898145	79758	-11.26	<.0001	0
YearBuilt	1	448.06847	41.11992	10.90	<.0001	3.11033
total_SF	1	47.94454	1.24588	38.48	<.0001	2.18016
total_baths	1	7137.00080	1098.01401	6.50	<.0001	1.83676
good_kitchen	1	5114.75146	1628.58766	3.14	0.0017	1.92785
good_fireplace	1	5587.92777	1662.48892	3.36	0.0008	1.46337
good_exterior	1	10620	1916.40496	5.54	<.0001	2.29912
quality_index	1	1441.45185	82.73133	17.42	<.0001	1.38250
central_air	1	-11873	3541.64889	-3.35	0.0008	1.22244
fireplace_ind	1	8325.14250	1504.36520	5.53	<.0001	1.66208
garage_ind	1	7249.84653	3883.73858	1.87	0.0622	1.11225
good_basement_ind	1	3557.07044	1869.02846	1.90	0.0573	2.57132

Common practical experience suggests that if any VIFs exceeds 5 or 10, it is an indication that multicollinearity exists in the model. Based on the VIF values for each parameter in Table 7, we conclude that no multicollinearity exists among the predictor variables given the fact that none of the VIFs exceeds 5.

Operational Validation

By leveraging the eleven predictor variable candidates suggested from the variable selection and performing model adequacy checking, the final model resulted strong predictive accuracy. To assess the operational accuracy of the final model, we placed the predictive scores, absolute value for each observation's actual vs predicted value for the response variable SalePrice, into three categories: Grade 1 (within 10% of the observed value), Grade 2 (between 10-15% of the observed value), Grade 3 (everything else).

Table 8 – Frequency Table for Prediction Grade Categories

Prediction_Grade	Frequency	Percent	Cumulative Frequency	Cumulative Percent
01: Grade 1	344	64.66	344	64.66
02: Grade 2	101	18.98	445	83.65
03: Grade 3	87	16.35	532	100.00

Based on the results of the prediction scores using the test sample dataset, we can see in Table 8 that 83.7% of the predicted values were within 15% of the observed value. Depending on the success criteria from management, this could be considered a usable model for predicting the sale price for homes in the Ames, Iowa area.

Conclusion

One can understand how automated variable selection methods are extremely useful in the model selection process, however, in this particular exercise all methods suggested the same “best” model. Given the fact that we had a biased selection of the original candidate predictor variables, we could further explore the incorporation of additional variables from the dataset. After working with the Ames housing dataset for the past several weeks, we have learned the importance of proper exploratory data analysis leveraging both statistics and visual aids, such as scatter plots, histograms, or correlation matrices. The final model selection process of actually calculating the predictive accuracy by using out-of-sample observations on the model was the most beneficial takeaway from all assignments to date. This was the final step in determining how well the models would perform in the final intended environment.

Our final model yielded a fairly high Adjusted R-squared of 0.8887, which shows extremely good fit to the data. However, I think if we wanted to explore ways of further improving predictive accuracy of our model, we should partner with subject matter experts in real estate in the Ames, Iowa area. By partnering with a realtor in the area, they could further inform us which predictor variables could be important when determining sale price for homes in the area. This would of course introduce bias to the model, but we need to ensure we include variables that are not only statistically significant, but also those that are practical from a business perspective.

SAS Code Output

```

1  *-----
2  * Nate Bitting
3  * Assignment 5
4  *-----;
5
6  * Code used to get the data into my library;
7  ods graphics on;
8  libname mydata '/courses/d6fc9ae5ba27fe300/c_3505/SAS_Data/' access=readonly;
9  proc datasets library=mydata; run; quit;
10
11  *-----
12  * Create original dataset by filtering out unnecessary data and adding new categorical
13  * features to the dataset
14  *-----;
15
16  * Smart data formats for sale price and square footage;
17  proc format;
18  value price_sfmt
19  . = '10: Missing'
20  1 -< 100000 = '01: [1; 100,000)'
21  100000 -< 150000 = '02: [100,000; 150,000)'
22  150000 -< 200000 = '03: [150,000; 200,000)'
23  200000 -< 250000 = '04: [200,000; 250,000)'
24  250000 -< 300000 = '05: [250,000; 300,000)'
25  300000 -< 350000 = '06: [300,000; 350,000)'
26  350000 -< 400000 = '07: [350,000; 400,000)'
27  400000 - high = '08: [400,000+]'
28  other = '09: Invalid Value'
29  ; * use a semi-colon to end each format in the proc format statement;
30  * Note on how we use the < to create open intervals.
31  * The dash - will create closed intervals, and
32  * hence the dash should only be used with discrete values;
33  value sqft_sfmt
34  . = '08: Missing'
35  1 - 1000 = '01: [1; 1,000]'
36  1000 -< 1500 = '02: (1,000; 1,500]'
37  1500 -< 2000 = '03: (1,500; 2,000]'
38  2000 -< 2500 = '04: (2,000; 2,500]'
39  2500 -< 3000 = '05: (2,500; 3,000]'
40  3000 - high = '06: (3,000+]'
41  other = '07: Invalid Value'
42  ;
43  run;
44
45  * Dataset before removing outliers;
46  Data building;
47  SET mydata.ames_housing_data;
48
49  * filter on only single family homes that meet specific criteria;
50  if (SaleCondition = 'Normal');
51  if (BldgType = '1Fam'); * single family homes only;
52  if (Zoning in ('RH', 'RL', 'RP', 'RM')); * residential zones only;
53  if (Street='Pave'); * paved streets only;
54  if (Utilities='AllPub'); * only homes with public utilities;
55
56  log_price = log(SalePrice); *create a variable for the natural log of SalePrice;
57
58  * create new variables by combining multiple variables in the housing dataset;
59  total_SF = max(GrLivArea,0) + max(TotalBsmtSF,0);
60  total_baths = max(FullBath,0) + max(BsmtFullBath,0);
61  total_halfbaths = max(HalfBath,0) + max(BsmtHalfBath,0);
62  total_baths_calc = total_baths + total_halfbaths;
63
64  * Neighborhood dummy variables;
65  if (Neighborhood = 'Blmngtn') then Blmngtn=1; else Blmngtn=0;
66  if (Neighborhood = 'Blueste') then Blueste=1; else Blueste=0;
67  if (Neighborhood = 'BrDale') then BrDale=1; else BrDale=0;
68  if (Neighborhood = 'BrkSide') then BrkSide=1; else BrkSide=0;
69  if (Neighborhood = 'ClearCr') then ClearCr=1; else ClearCr=0;
70  if (Neighborhood = 'CollgCr') then CollgCr=1; else CollgCr=0;
71  if (Neighborhood = 'Crawfor') then Crawfor=1; else Crawfor=0;
72  if (Neighborhood = 'Edwards') then Edwards=1; else Edwards=0;
73  if (Neighborhood = 'Gilbert') then Gilbert=1; else Gilbert=0;
74  if (Neighborhood = 'Greens') then Greens=1; else Greens=0;
75  if (Neighborhood = 'GrnHill') then GrnHill=1; else GrnHill=0;
76  if (Neighborhood = 'IDOTRR') then IDOTRR=1; else IDOTRR=0;
77  if (Neighborhood = 'Landmrk') then Landmrk=1; else Landmrk=0;
78  if (Neighborhood = 'MeadowV') then MeadowV=1; else MeadowV=0;
79  if (Neighborhood = 'Mitchel') then Mitchel=1; else Mitchel=0;
80  if (Neighborhood = 'NAmes') then NAmes=1; else NAmes=0;

```



```

81 if (Neighborhood = 'NridgHt') then NridgHt=1; else NridgHt=0;
82 if (Neighborhood = 'OldTown') then OldTown=1; else OldTown=0;
83 if (Neighborhood = 'SWISU') then SWISU=1; else SWISU=0;
84 if (Neighborhood = 'Sawyer') then Sawyer=1; else Sawyer=0;
85 if (Neighborhood = 'SawyerW') then SawyerW=1; else SawyerW=0;
86 if (Neighborhood = 'Somerst') then Somerst=1; else Somerst=0;
87 if (Neighborhood = 'StoneBr') then StoneBr=1; else StoneBr=0;
88 if (Neighborhood = 'Timber') then Timber=1; else Timber=0;
89 if (Neighborhood = 'Veenker') then Veenker=1; else Veenker=0;
90
91 * KitchenQual dummy variable;
92 if (KitchenQual in ('Ex', 'Gd')) then good_kitchen=1; else good_kitchen=0;
93
94 * FireplaceQu dummy variable;
95 if (FireplaceQu in ('Ex', 'Gd')) then good_fireplace=1; else good_fireplace=0;
96
97 * ExterQual dummy variable;
98 if (ExterQual in ('Ex', 'Gd')) then good_exterior=1; else good_exterior=0;
99
100 * Foundation dummy variables;
101 if (Foundation = 'BrkTil') then Foundation_BrkTil=1; else Foundation_BrkTil=0;
102 if (Foundation = 'CBlock') then Foundation_CBlock=1; else Foundation_CBlock=0;
103 if (Foundation = 'PConc') then Foundation_PConc=1; else Foundation_PConc=0;
104 if (Foundation = 'Slab') then Foundation_Slab=1; else Foundation_Slab=0;
105 if (Foundation = 'Stone') then Foundation_Stone=1; else Foundation_Stone=0;
106 if (Foundation = 'Wood') then Foundation_Wood=1; else Foundation_Wood=0;
107
108 * MoSold dummy variables;
109 if (MoSold = 1) then jan_sold=1; else jan_sold=0;
110 if (MoSold = 2) then feb_sold=1; else feb_sold=0;
111 if (MoSold = 3) then mar_sold=1; else mar_sold=0;
112 if (MoSold = 4) then apr_sold=1; else apr_sold=0;
113 if (MoSold = 5) then may_sold=1; else may_sold=0;
114 if (MoSold = 6) then jun_sold=1; else jun_sold=0;
115 if (MoSold = 7) then jul_sold=1; else jul_sold=0;
116 if (MoSold = 8) then aug_sold=1; else aug_sold=0;
117 if (MoSold = 9) then sep_sold=1; else sep_sold=0;
118 if (MoSold = 10) then oct_sold=1; else oct_sold=0;
119 if (MoSold = 11) then nov_sold=1; else nov_sold=0;
120 if (MoSold = 12) then dec_sold=1; else dec_sold=0;
121
122 * Construct a composite quality index;
123 quality_index = OverallCond*OverallQual;
124
125 * Central Air Indicator;
126 if (CentralAir='Y') then central_air=1; else central_air=0;
127 * Fireplace Indicator;
128 if (Fireplaces>0) then fireplace_ind=1; else fireplace_ind=0;
129 * Garage Indicator;
130 if (GarageCars>0) then garage_ind=1; else garage_ind=0;
131 * Good Basement Indicator;
132 if (BsmtQual in ('Ex','Gd')) or (BsmtCond in ('Ex','Gd'))
133 then good_basement_ind=1;
134 else good_basement_ind=0;
135
136 *apply the put function to create the categorical variables for the various scales of price and sqft;
137 price_cat = put(SalePrice,price_sfmt.);
138 sqft_cat = put(total_SF,sqft_sfmt.);
139
140 if (YearBuilt>1920);
141 run; quit;
142
143
144 * Review new dataset to ensure no missing values based on filtering;
145 proc contents data=building;
146 run; quit;
147
148
149 -----
150 * Include all predictor variable candidates and identify outliers using the Studentized
151 * Residuals
152 -----;
153
154 proc reg data=building;
155 model SalePrice = YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior
156 quality_index central_air fireplace_ind garage_ind good_basement_ind;
157 output out=outdata (keep = SalePrice YearBuilt total_SF total_baths good_kitchen
158 good_fireplace good_exterior quality_index central_air fireplace_ind
159 garage_ind good_basement_ind studentr) rstudent=studentr;
160 run;

```

```

162 * Assess the plot of the studentized residuals to see how many exceed an absolute value of 2;
163 proc univariate data=outdata plot;
164     var studentr;
165 run;
166
167 * Remove any observations with an absolute studentized residual greater than 2;
168 data clean_data;
169     SET outdata;
170     if abs(studentr) > 2 then delete;
171 run;
172
173 data clean_data_2;
174     set clean_data;
175     * generate a uniform(0,1) random variable with seed set to 123;
176     u = uniform(123);
177     if (u < 0.70) then train = 1;
178     else train = 0;
179     if (train=1) then train_response=SalePrice;
180     else train_response=.;
181     if (train=0) then test_response=SalePrice;
182     else test_response=.;
183 run;
184
185 -----
186 * Build a model using the adjusted r-squared variable selection: Model_AdjR2
187 -----;
188
189 proc reg data=clean_data_2;
190     model train_response = YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior
191         quality_index central_air fireplace_ind garage_ind good_basement_ind /
192     selection = adjrsq vif;
193 run;
194 * The Adjusted R2 selection method suggested to use all 11 predictor variables;
195 -----;
196
197
198 -----
199 * Build a model using the adjusted r-squared variable selection: Model_MaxR
200 -----;
201
202 proc reg data=clean_data_2;
203     model train_response = YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior
204         quality_index central_air fireplace_ind garage_ind good_basement_ind /
205     selection = MAXR vif;
206 run;
207 * The MaxR selection method suggested to use all 11 predictor variables;
208 -----;
209
210
211 -----
212 * Build a model using the adjusted r-squared variable selection: Model_MCP
213 -----;
214
215 proc reg data=clean_data_2;
216     model train_response = YearBuilt total_SF total_baths good_kitchen good_fireplace
217         good_exterior quality_index central_air fireplace_ind garage_ind good_basement_ind /
218     selection = cp;
219 run;
220 * The Mallows Cp selection method suggested to use all 11 predictor variables;
221
222
223 -----
224 * Build a model using the adjusted r-squared variable selection: Model_F
225 -----;
226
227 proc reg data=clean_data_2;
228     model train_response = YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior
229         quality_index central_air fireplace_ind garage_ind good_basement_ind /
230     selection = FORWARD slentry=.15;
231 run;
232 * The Forward selection method suggested to use all 11 predictor variables;
233 -----;
234
235
236 -----
237 * Build a model using the adjusted r-squared variable selection: Model_B
238 -----;
239
240 proc reg data=clean_data_2;
241     model train_response = YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior
242         quality_index central_air fireplace_ind garage_ind good_basement_ind /

```

```

243 selection = BACKWARD slstay=.15;
244 run;
245 * The Backward selection method suggested to use all 11 predictor variables;
246 -----;
247
248
249 -----
250 * Build a model using the adjusted r-squared variable selection: Model S
251 -----;
252
253 proc reg data=clean_data_2;
254     model train_response = YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior
255         quality_index central_air fireplace_ind garage_ind good_basement_ind /
256         selection = STEPWISE aic bic mse adjrsq slentry=.15 slstay=.15;
257 run;
258 * The Stepwise selection method suggested to use all 11 predictor variables;
259 -----;
260
261 -----
262 * Calculate the AIC, BIC, Adjusted R-squared, Mallows Cp, MSE for the training sample
263 -----;
264 proc reg data=clean_data_2;
265     model train_response = YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior
266         quality_index central_air fireplace_ind garage_ind good_basement_ind /
267         selection = adjrsq aic bic cp mse;
268 run;
269
270 -----
271 * Calculate the Mean Absolute Error for the training sample
272 -----;
273 proc reg data=clean_data_2;
274     model train_response = YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior
275         quality_index central_air fireplace_ind garage_ind good_basement_ind;
276     output out=residuals_final (keep = resid_final) r=resid_final;
277 run;
278
279 data abs_resid;
280     set residuals_final;
281     abs_resid = abs(resid_final);
282 run;
283
284 proc means data=abs_resid mean;
285     var abs_resid;
286 run;
287
288 -----
289 * Use the suggested model that includes all 11 predictor candidates and
290 * perform cross validation by outputting the estimated values from the test set
291 -----;
292 proc reg data=clean_data_2 outest=RegOut;
293     SalePrice_Hat: model train_response = YearBuilt total_SF total_baths good_kitchen
294         good_fireplace good_exterior quality_index central_air fireplace_ind
295         garage_ind good_basement_ind;
296     title "Model with all 11 Predictor Candidates";
297 run;
298
299 proc print data=RegOut;
300     title2 'OUTEST= Data Set from PROC REG';
301 run;
302
303 * calculate the estimated values using the test dataset;
304 proc score data=clean_data_2 score=RegOut out=RScoreP type=parms;
305     var YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior quality_index
306         central_air fireplace_ind garage_ind good_basement_ind;
307 run;
308
309 proc score data=clean_data_2 score=RegOut out=RScoreR type=parms;
310     var train_response YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior
311         quality_index central_air fireplace_ind garage_ind good_basement_ir
312 run;
313
314 * Output the scores using the model built with the training dataset against the test dataset;
315 proc score data=test_dataset score=RegOut out=NewPred type=parms
316     nostd predict;
317     var train_response YearBuilt total_SF total_baths good_kitchen good_fireplace good_exterior
318         quality_index central_air fireplace_ind garage_ind good_basement_ir
319 run;
320
321 * Smart data formats for sale price and square footage;
322 proc format;
323 value pred_acc_sfmt

```

```

324 0 -< .1 = '01: Grade 1'
325 .1 -< .15 = '02: Grade 2'
326 other = '03: Grade 3'
327 ;
328 run;
329
330 *create a new dataset that contains the test set data and the predictive scores;
331 data prediction_Data;
332 set NEWPRED;
333 if (train_response = null);
334 pred_score = abs(test_response / SalePrice_Hat - 1);
335 abs_resid = abs(test_response - SalePrice_Hat);
336 error_term = test_response - SalePrice_Hat;
337 sq_error = error_term**2;
338 Prediction_Grade = put(pred_score,pred_acc_sfmt.);
339 run;
340
341 * calculate the MAE and MSE from the test sample;
342 proc means data=prediction_Data mean;
343 var abs_resid sq_error;
344 run;
345
346 * create a frequency table to show the operational accuracy of the model;
347 proc freq data=prediction_Data;
348 TABLES Prediction_Grade;
349 run;
350

```