

Project: Housing Prices Analysis

Albert Joseph, abj37

Marc Yarkony, mdy7

1.

a.

Quantitative	Qualitative
price (\$1000) bedrooms bathrooms sqft_living floors sqft_above yr_built yr_renovated latitude longitude sqft_living15	Id view condition grade

b. Price: The distribution is right-skewed. As shown in Figure 1.

c. As shown in Figure 2.

bedrooms: The distribution is right-skewed.

bathrooms: The distribution is left-skewed.

sqft_living: The distribution is right-skewed.

floors: The distribution is right-skewed.

sqft_above: The distribution is right-skewed.

yr_built: The distribution is left-skewed.

yr_renovated: The distribution is right-skewed.

latitude: The distribution is normal.

longitude: The distribution is normal.

sqft_living15: The distribution is right-skewed.

view: The distribution is right-skewed (zero has the highest frequency).

condition: The distribution is somewhat normal (three has the highest frequency).

grade: The distribution is normal.

d. There seems to be a positive correlation ($>.5$) between price and sqft_living, grade, sqft_living15. There seems to be a positive correlation between bedroom and sqft_living, There seems to be a positive correlation between bathrooms and sqft_living, floors, grade, sqft_above, yr_built, sqft_living15. Overall, as sqft_living increased, other variables associated with area increased as well (price, bathroom, etc.). As shown in Figure 3, Figure 4, and Figure 5.

- e. To predict price, we should use `sqft_living`, `grade`, and `sqft_living15` because there is a positive correlation ($>.5$) between these variables. Correlation between variables is shown in Figure 3.

2.

- a. At 95% confidence, the true average selling price of King County houses lies between [462.1161, 529.5443]. At 95% confidence, the true average square footage of King County houses lies between [1933.71, 2158.041]. As shown in Figure 6.
- b. The P-value for this test is .0054. Since $P < \alpha$, we reject the null hypothesis and conclude that houses with basements have a higher average selling price than houses without basements. As shown in Figure 7.
- c. CI: [21.419, 161.857] Since 0 is outside of CI, we reject the null hypothesis and conclude that houses with basements have a higher average selling price than houses without basements. As shown in Figure 7.
- d. Since there is overlap between the box plots, we cannot make a conclusion about the difference in means. As shown in Figure 8.

3.

- a.
 - i) Yes, there is a relationship between the Predictor and the Response. (See Figure 9)
 - ii) The strength of relationship between the predictor and the response can be measured by the correlation between the two. The correlation between the predictor and the response is 0.674729. (See Figure 9)
 - iii) The relationship between the predictor and the response is positive, evident from its positive correlation and its positive slope. (See Figure 9)
 - iv) Based on the linear fit line from our graph, `b1` returned a value of \$0.2028616 (In \$1000). Therefore, \$0.2028616 (In \$1000) is the cost for each additional square foot. (See Figure 9)
 - v) The predicted selling price for a 2000 square foot home is \$486.517691 (In \$1000). The 95% Confidence Interval for a selling price for a 2000 square foot house is [461.53, 511.5]. The 95% prediction Interval for selling price for a 2000 square foot home is [122.48, 859.55]. (See Figure 19)
- b. There are a few problems we noticed while analyzing the diagnostic plots of the least squares regression fit. The residual by predicted plot looked like a funnel plot, inferring that the variability of variables is different across points of the graph. For the actual by predicted plot, the plotted points were slightly lower than the preferred angle of 45 degrees. Lastly, the residual Normal Quantile plot had a handful of plots outside of the normal distribution lines.

c. **Quadratic fit - R^2 , R^2_{adj} = (0.457387, 0.452194)**

Price (\$1000) = $89.93 + 0.194 \cdot \text{sqft_living} + 1.28e^{-5} \cdot (\text{sqft_living})^2$. (See Figure 10)

Linear fit with logarithm transformation - R^2 , R^2_{adj} = (0.4361, 0.4334)

$\text{Log}(\text{price} (\$1000)) = 5.3135436 + 0.0003804 \cdot \text{sqft_living}$. (See Figure 11)

Linear Fit - R^2 , R^2_{adj} = (0.45526, 0.452666)

price (\$1000) = $80.794491 + 0.2028616 \cdot \text{sqft_living}$. (See Figure 12)

Although the Quadratic fit has a slightly better R^2 and R^2_{adj} value, it is not the only important criterion when picking models. Of the 3 models, the linear fit with logarithmic transformation had a significantly better RMSE. When looking over the linear models on the graph, it becomes even more evident that this model clearly plots the points better than the other two.

d. **Linear regression models to predict price for each predictor**

Bedrooms - price (\$1000) = $168.87516 + 95.474471 \cdot \text{bedrooms}$

Bathrooms - price (\$1000) = $119.01325 + 180.63356 \cdot \text{bathrooms}$

Sqft_living - price (\$1000) = $80.794491 + 0.2028616 \cdot \text{sqft_living}$

Floors - price (\$1000) = $296.20924 + 129.22028 \cdot \text{floors}$

View - price (\$1000) = $470.8305 + 120.45303 \cdot \text{view}$

Condition - price (\$1000) = $441.6712 + 16.263039 \cdot \text{condition}$

Grade - price (\$1000) = $-682.9709 + 154.16769 \cdot \text{grade}$

Sq_ft above - price (\$1000) = $169.16239 + 0.1851492 \cdot \text{sqft_above}$

Yr built - price (\$1000) = $-1162.812 + 0.8414649 \cdot \text{yr_built}$

renovated - price (\$1000) = $489.15857 + 0.0888655 \cdot \text{yr_renovated}$

Latitude - price (\$1000) = $-38511.59 + 820.36334 \cdot \text{latitude}$

Longitude - price (\$1000) = $15389.831 + 121.86585 \cdot \text{longitude}$

Sqft_living15 - price (\$1000) = $113.79225 + 0.1934628 \cdot \text{sqft_living15}$

With our models, the model with the statistically significant association between the response and the predictor was grade. Grade resulted in the highest correlation, R^2 value, R^2_{adj} value, and the lowest RMSE value.

(See Figure 13)

4.

- a. i) There is definitely a relationship between the predictors and the response with the R^2 value being higher than all other single predictor options and the RMSE lower than all other single predictor options. (See Figure 14)
- ii) The top predictors that seem to have a statistically significant relationship with the response are Latitude, Grade, view, year built, bathrooms and conditions (See Figure 14)
- iii) Because of the large coefficient value for “Latitude” the variable suggests that the response is significantly affected by each individual change in the predictor value. It is also a positive line based on the coefficient. (See Figure 14)

iv) The coefficient value for “grade” suggests a positive line as well at a significant slope value, causing a significant difference for each change in predictor value. (See Figure 14)

- b. Based on our plotted graph there are many more residual outliers plotted compared to when we plot the response to only one predictor. (See Figure 14)
- c. Our results from 4A resulted in a higher R^2 value and a lower RMSE than any plotted graph from 3D. When all predictors are taken into account for one plot it clearly makes a difference.
(Compare Rsquare of Figure 14 to every Rsquare in Figure 13)
- d. A few interactions that appear to be statistically significant are sqft_living*Grade, Grade*Latitude, Grade*View, and Sqft_Living*bedroom. The interactions are clear based on the non parallel lines. (See Figure 20)
- e. The transformations done were able to improve RMSE scores and R^2 values by a small margin. (See Figures 15-17)
- f. We narrowed down our Final Model predictors to the significant variables which were Latitude, grade, view, yr built, bathrooms and conditions. The log(x) transformation gave the best results compared to other transformations. Putting all of this together, our simplified model consisting of 6 predictors resulted in a R^2 value of .79, and a RMSE value of 0.2242. (See Figure 18)

5.

- a. Validation Column created in JMP file. (See JMP File, Last column)
- b. A few variables that help understand the association between houses built before 1980 and the variables are Floors, Grade, and Bathroom. Floors being the best indicator of the three. (See Figure 21)
- c. Completed, See JMP file. Split up the proportion of Training and Test as .8 to .2 (80% to 20%).
- d. The test error (Misclassification) of the model obtained is 0.1887 (See Figure 23)
- e. The value of K that seemed to perform the best data was when there was a range of k from 1-75 and the K value was k=75. This resulted in a misclassification of .16327 (See Figure 22)

- f. Our models returned slightly better classification results for the Logistic Regression Method instead of the K-Nearest neighbor method. Therefore, we would recommend using the Logistic Regression classification model to predict whether or not a house was built before 1980.

Bonus:

For the extra credit, we decided to put ourselves in the shoes of perspective home buyer in the King County area, Washington State, between May 2014 and May 2015. As a home buyer, there are a few factors that are important: bedrooms, bathrooms, sqft_living, floors, sqft_above, yr_built, yr_renovated, and sqft_living15. Since there are hundreds of houses to choose from, we performed cluster analysis to identify which houses are similar to each other, so that home buyers could fit themselves into a cluster, and then decide which houses to buy from there.

We decided to use Kmeans clustering because that was the easiest method for us to understand and apply. In order to determine how many clusters to use, we ran a factor analysis on the previously mentioned variables in JMP. From the elbow plot, we learned that two clusters would be sufficient. The factor analysis is shown in Bonus Figure 1.

Next, we ran the Kmeans cluster JMP, and found that in cluster one there would be 75 homes, and a cluster two there would be 137 homes. The average price for cluster one is \$689.893 and the average price for cluster two is \$389.591. Cluster one would be best for a home buyer with a higher budget who wants a newer home in larger living space. Cluster two would be best for a home buyer with a lower budget who wants a smaller living space, and older home. The Kmeans cluster results are shown in Bonus Figure 2.

Appendix:

Figure 1:

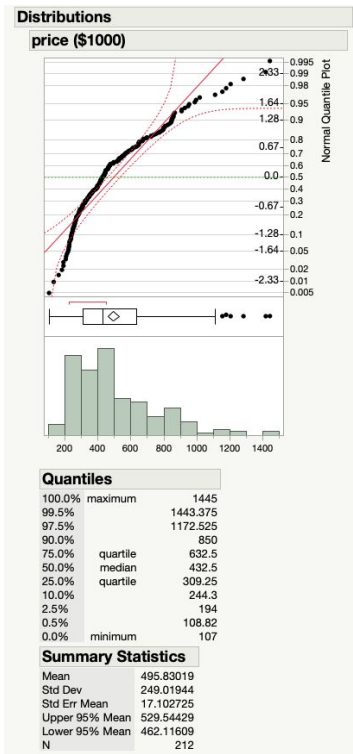


Figure 2:

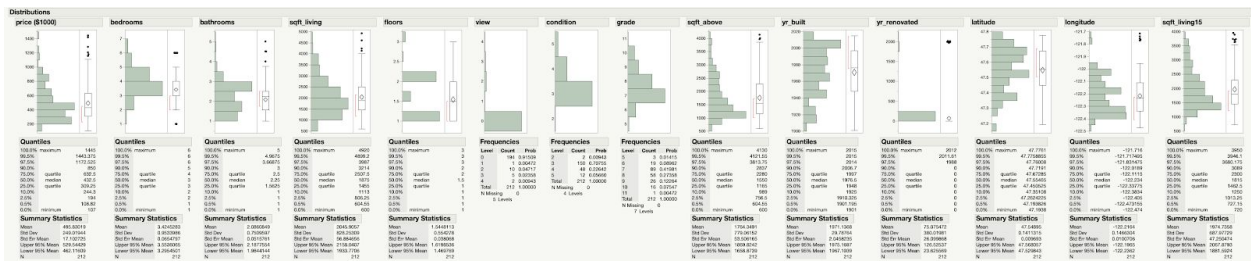


Figure 3:

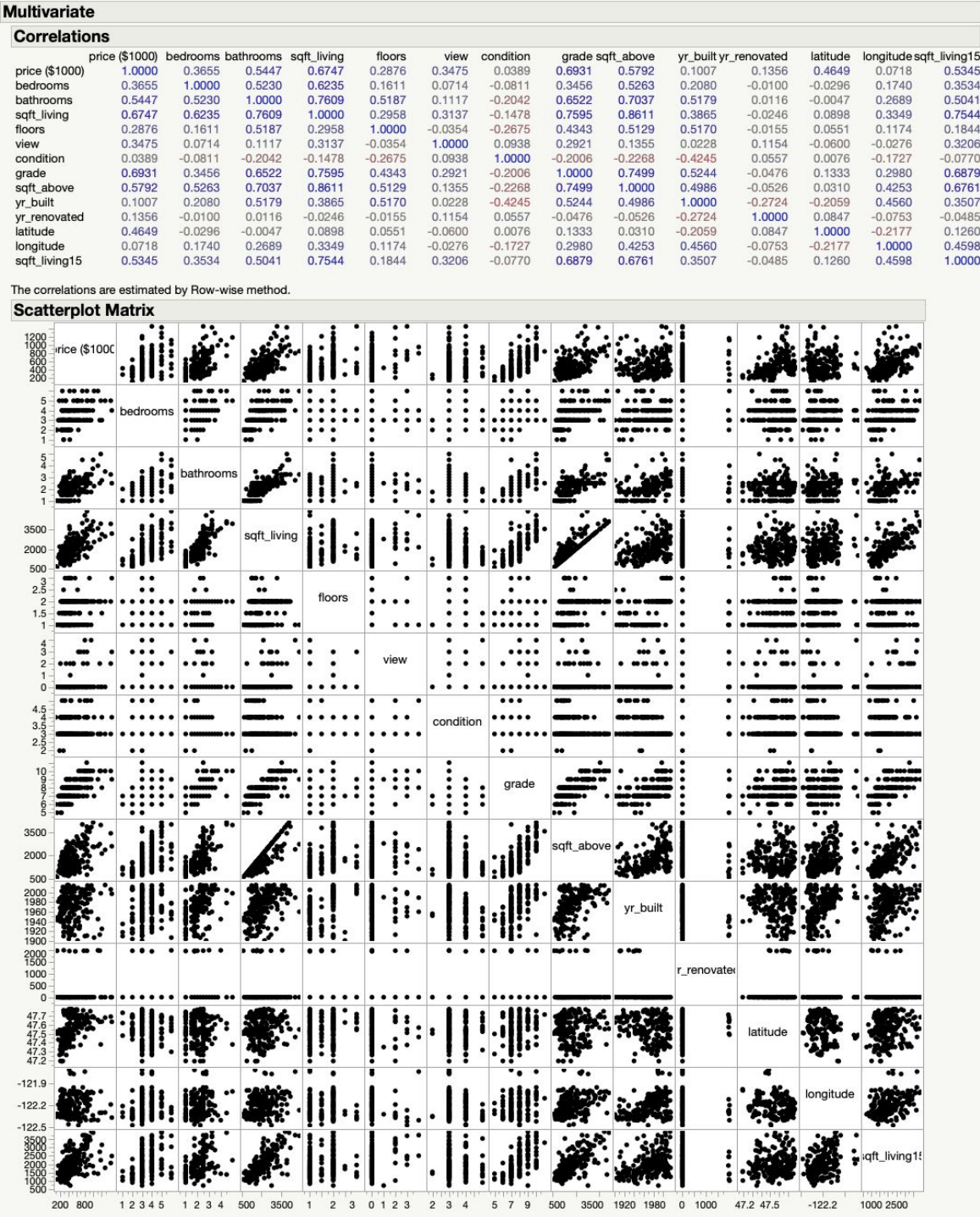


Figure 4:

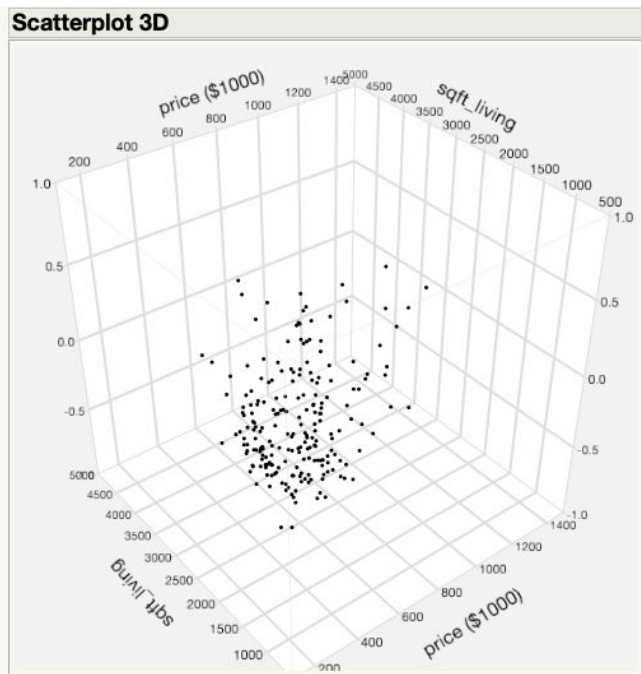


Figure 5:

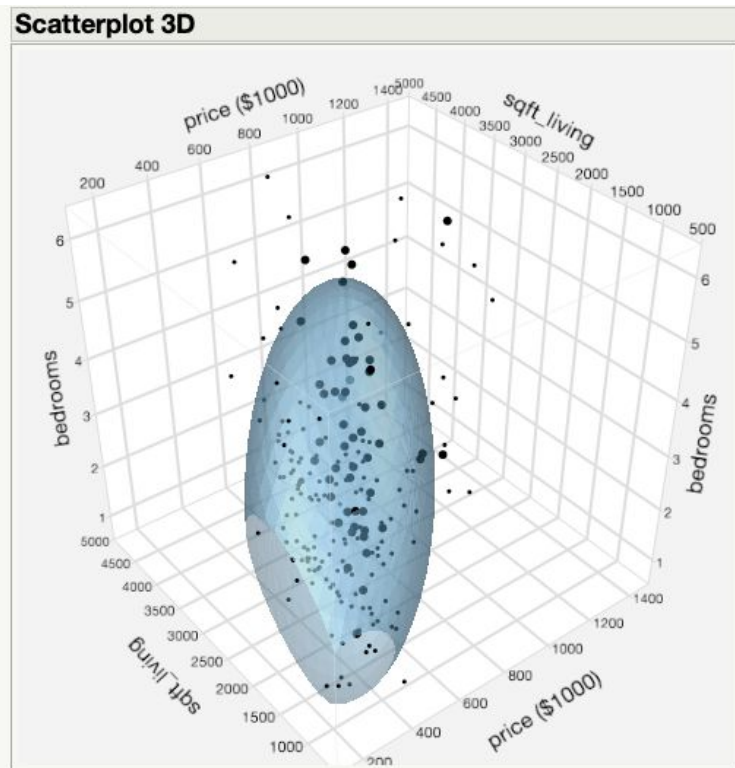


Figure 6:

Distributions									
price (\$1000)					sqft_living				
Confidence Intervals					Confidence Intervals				
Parameter	Estimate	Lower CI	Upper CI	1-Alpha	Parameter	Estimate	Lower CI	Upper CI	1-Alpha
Mean	495.8302	462.1161	529.5443	0.950	Mean	2045.906	1933.771	2158.041	0.950
Std Dev	249.0194	227.3574	275.2796	0.950	Std Dev	828.2531	756.204	915.596	0.950

Figure 7:

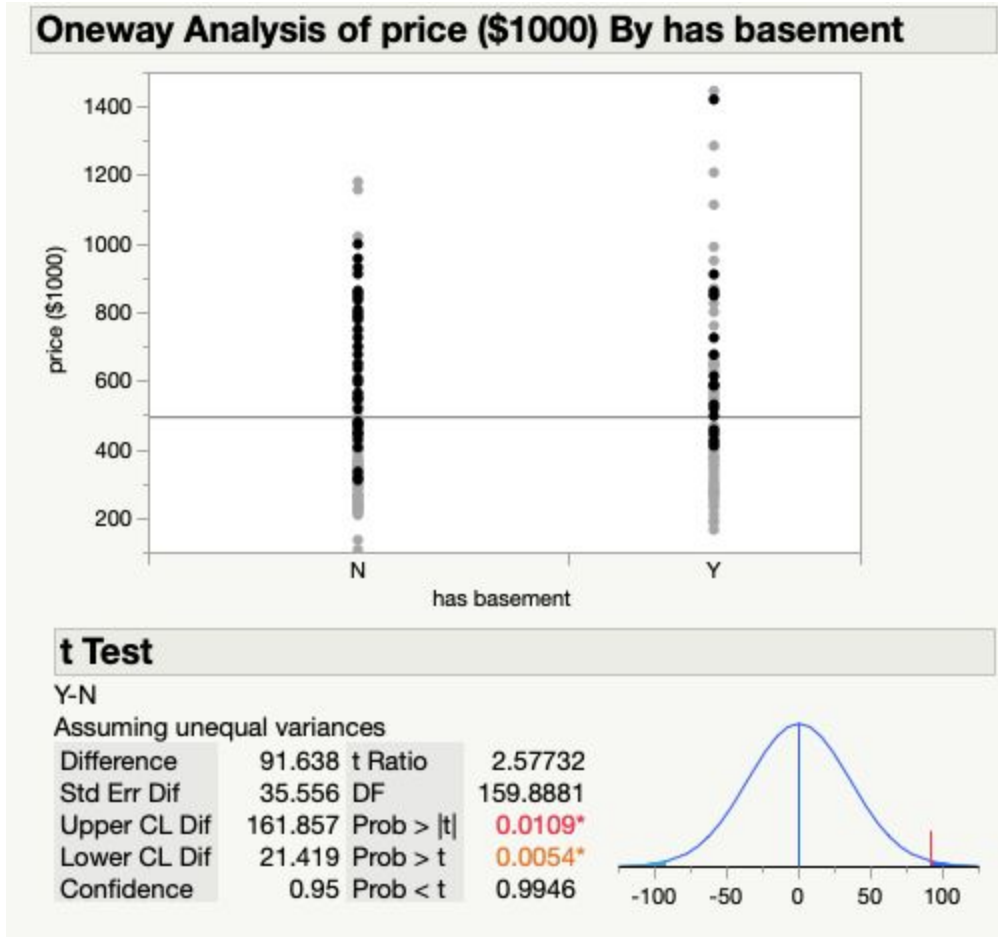


Figure 8:

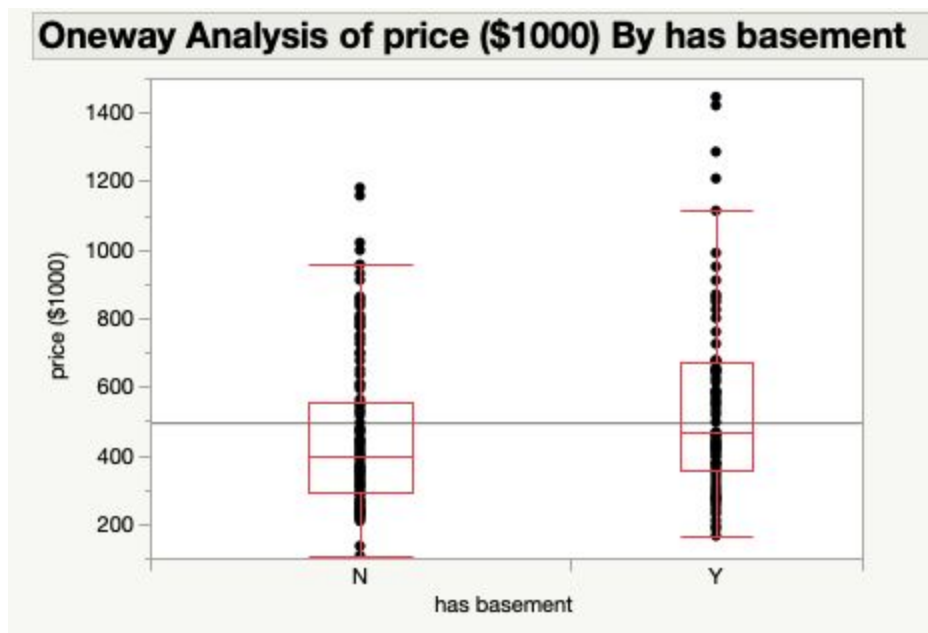


Figure 9

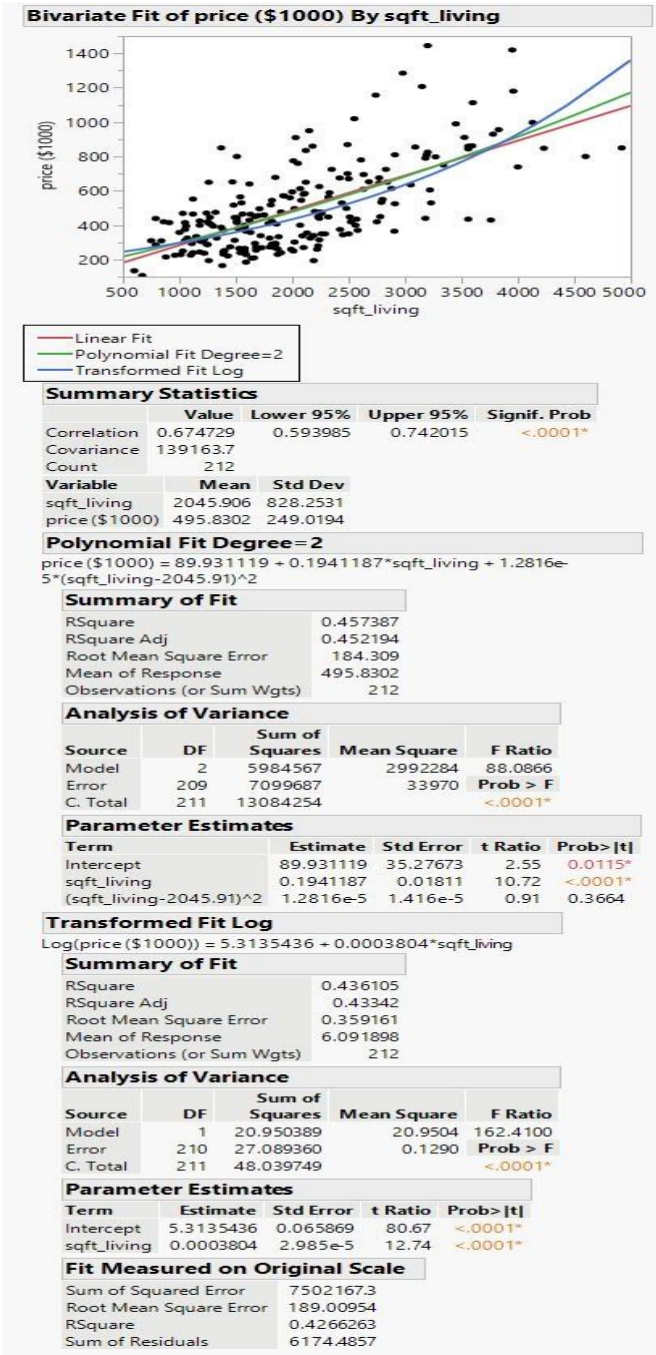


Figure 10

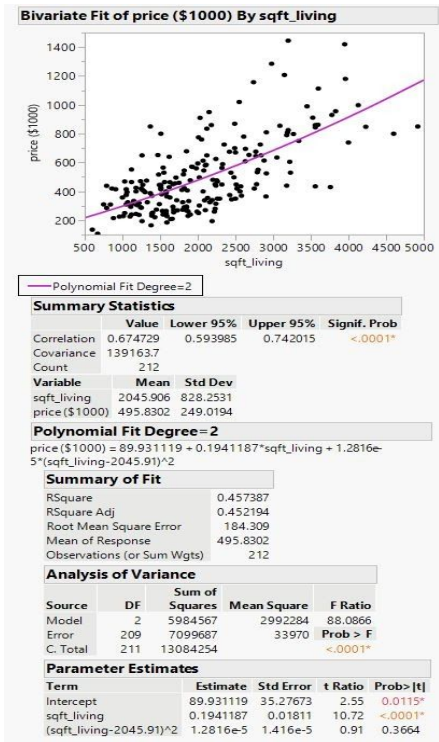


Figure 11

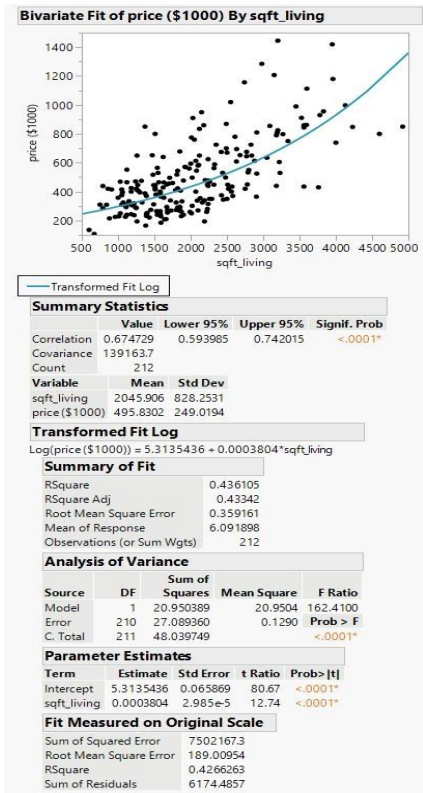


Figure 12

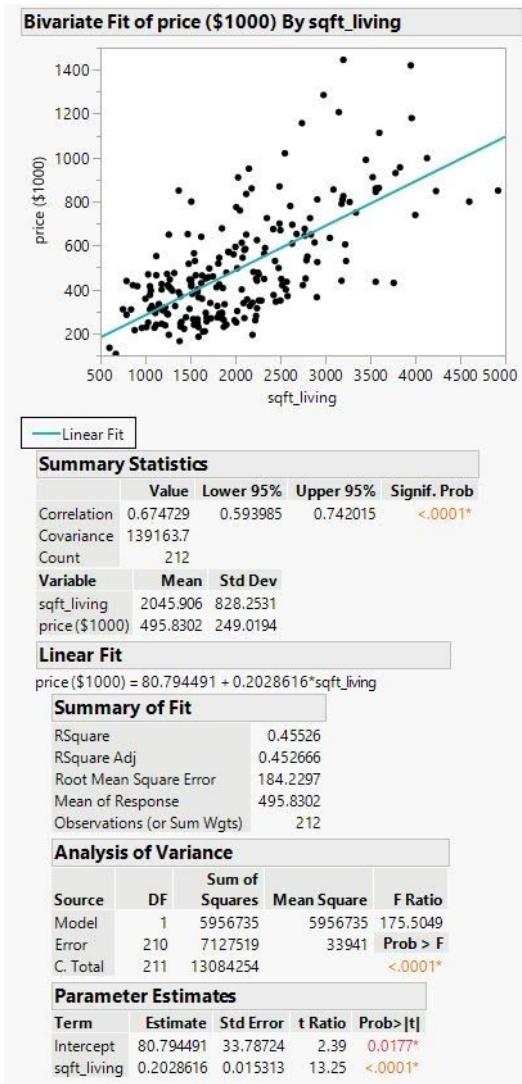


Figure 13

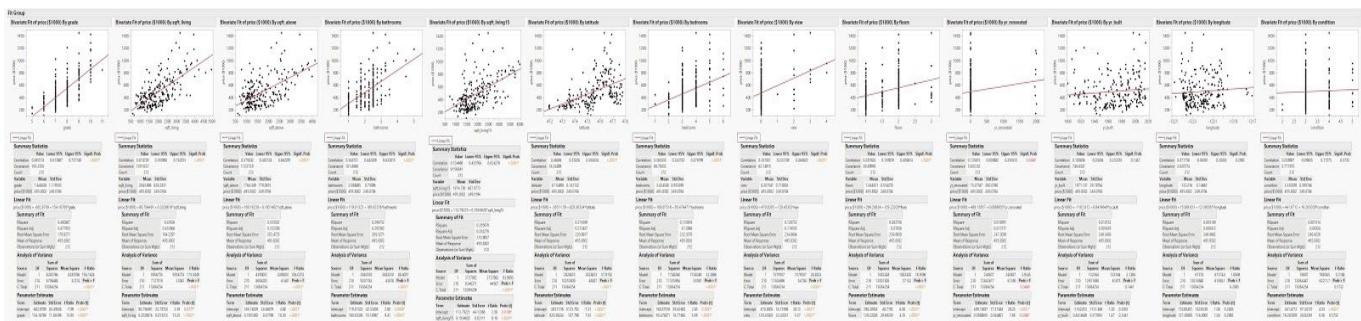


Figure 14)

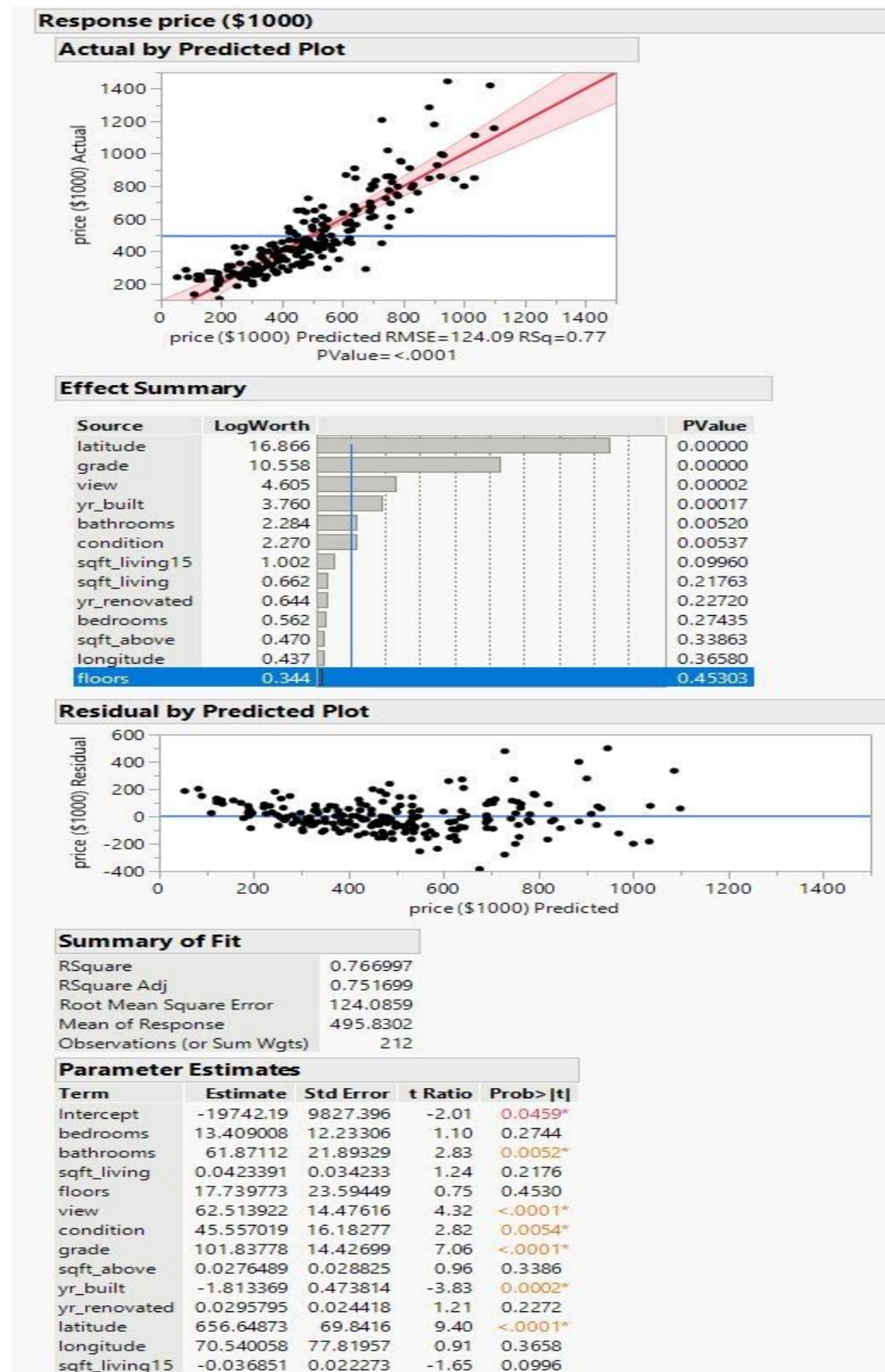


Figure 15

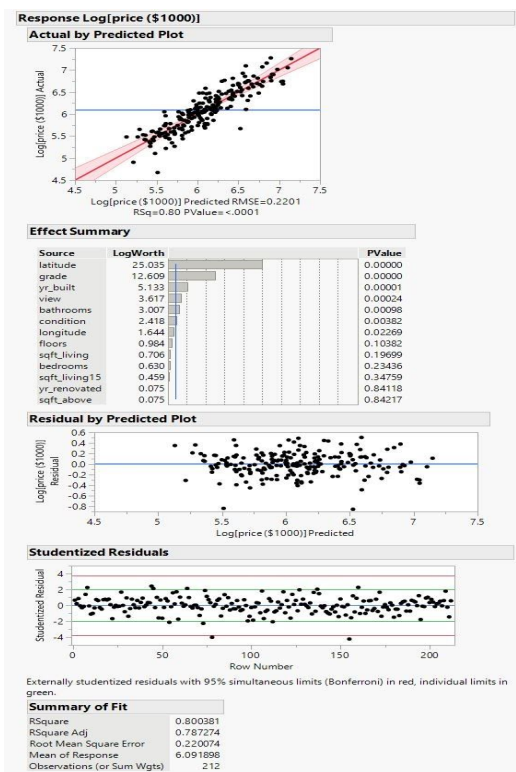


Figure 16

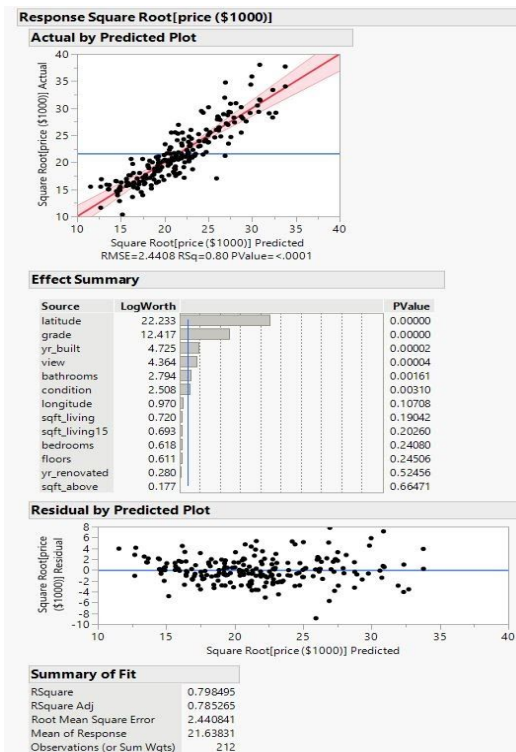


Figure 17

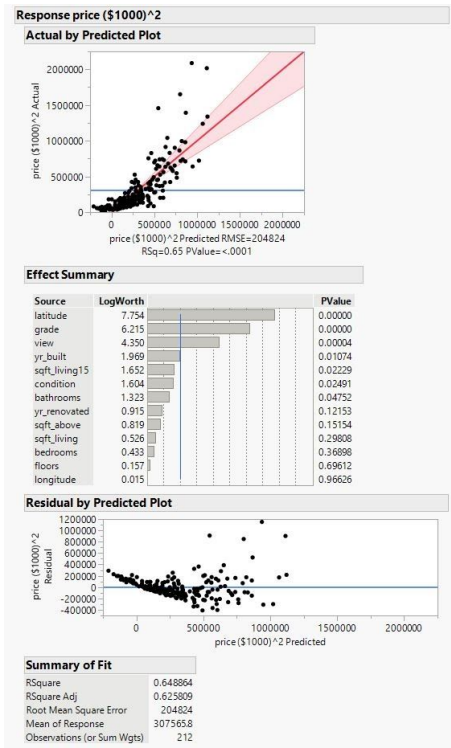


Figure 18

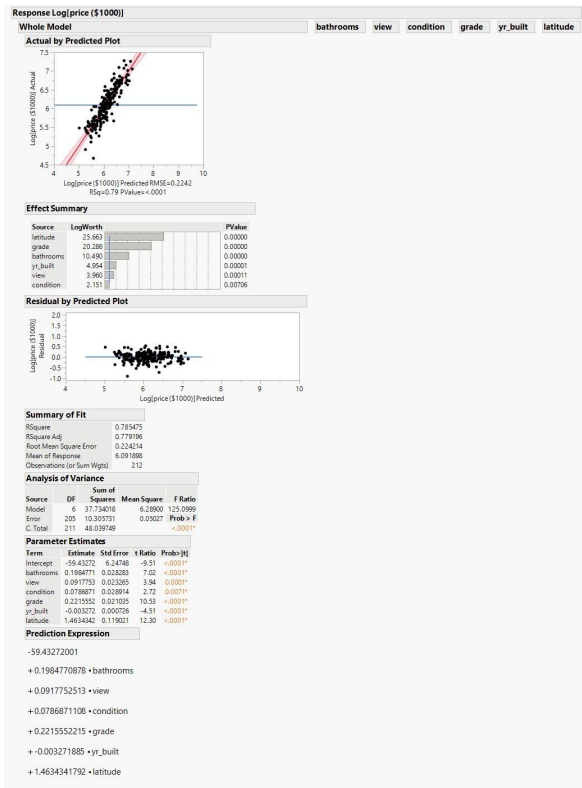


Figure 19

Lower 95% Indiv price (\$1000)	Upper 95% Indiv price (\$1000)	Lower 95% Mean price (\$1000)	Upper 95% Mean price (\$1000)
71.660693359	799.94389215	409.30807354	462.29651197
126.54205733	854.60779256	465.61960595	515.53024395
-50.94988683	679.12055147	277.29470221	350.87596244
-75.57951088	655.06339135	250.21218068	329.27169979
158.97824912	887.087313	497.763403	548.30215913
161.00338362	889.11941052	499.74191051	550.38088363
81.837670676	810.0530749	419.9211921	471.96955348
245.83320475	974.693334	580.06208983	640.46444891
-75.57951088	655.06339135	250.21218068	329.27169979
-110.5325065	621.0434427	211.6242683	298.88666792
81.837670676	810.0530749	419.9211921	471.96955348
22.724211935	751.50680522	357.38578754	416.84522962
-42.74795233	687.14754503	286.28742366	358.11216904
368.25559273	1099.7620989	690.669702	777.34798962
306.1564177	1036.0870815	635.03113138	707.2123678
414.17583404	1147.1581939	731.48750386	829.84652409
223.65847994	952.23850664	559.48706205	616.40992453
85.906710583	814.09849902	424.14415491	475.8610547
370.25483428	1101.8200894	692.45111244	779.62381119
87.940855267	816.12158635	426.25073353	477.81170809
-124.9456575	607.05596967	195.67415998	286.43615215
436.09176605	1169.8718141	750.9132875	855.05029261
49.249338429	777.72569492	385.69741531	441.27761804
24.766104459	753.52214472	359.577576	418.71067318
63.514609822	791.86104763	400.76251335	454.6131441
6.3800894275	735.39307162	339.77675421	401.99640684
122.48302778	850.52535809	461.53616625	511.49921962
280.04386999	1009.455613	611.38491041	678.11457258
26.80774738	755.53773381	361.767171	420.57831019
34.971822506	763.60258674	370.50288789	428.07152136
-1.797956925	727.34218991	330.92596012	394.61827287
-57.10394456	673.10291316	270.5381308	345.4608378

Figure 20

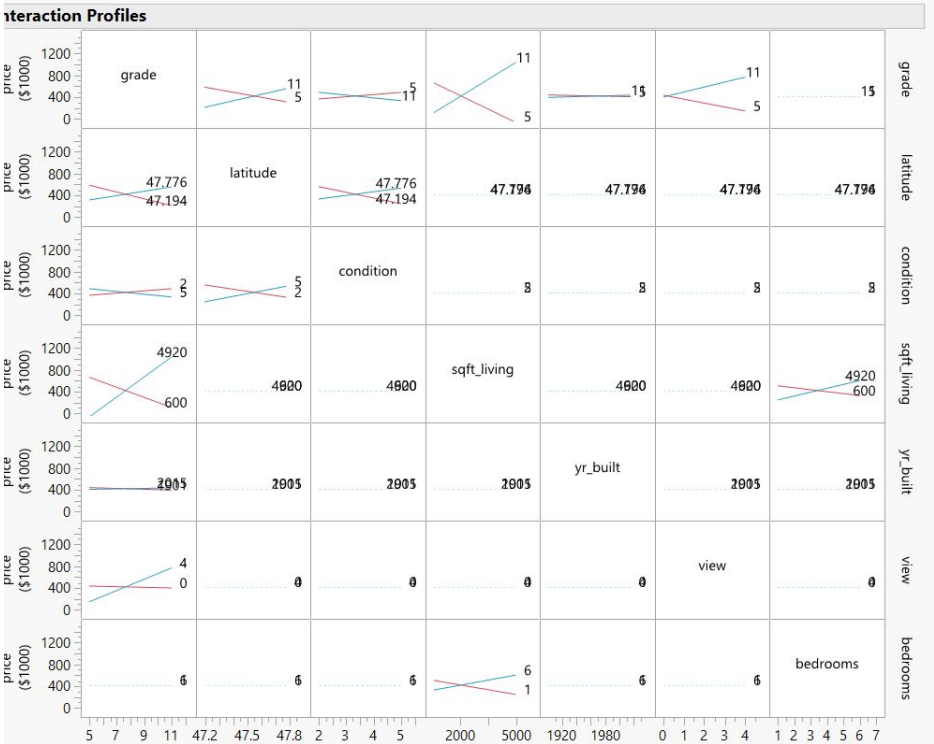


Figure 21

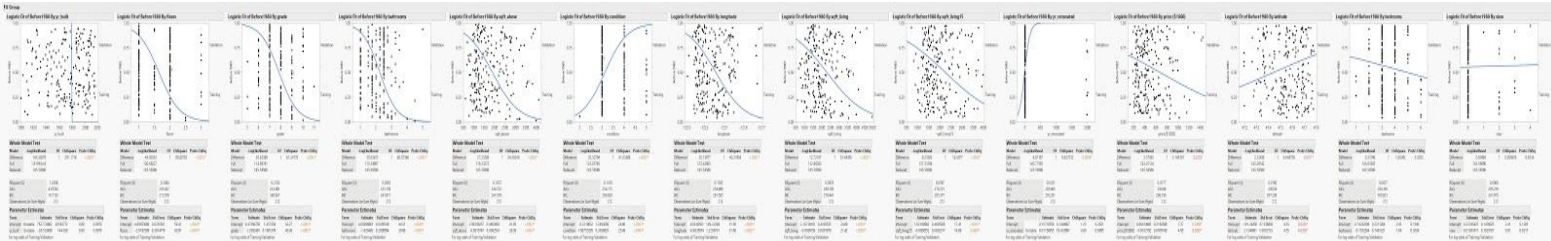
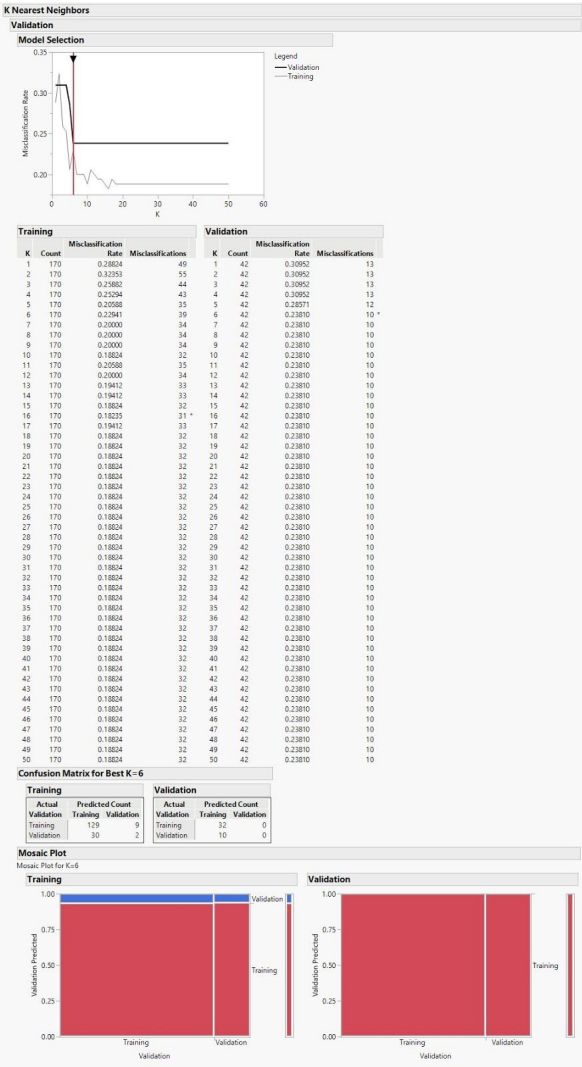


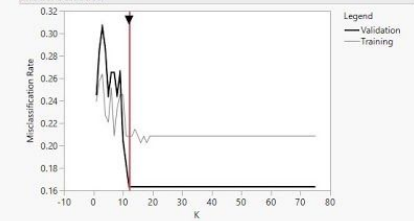
Figure 22



K Nearest Neighbors

Validation

Model Selection



Training				Validation			
K	Count	Misclassification Rate	Misclassifications	K	Count	Misclassification Rate	Misclassifications
1	163	0.2926	39	1	49	0.2490	12
2	163	0.2576	42	2	49	0.2871	14
3	163	0.2638	43	3	49	0.3061	15
4	163	0.2099	37	4	49	0.2871	14
5	163	0.2286	36	5	49	0.2490	12
6	163	0.2513	41	6	49	0.2631	13
7	163	0.2089	34	7	49	0.2631	13
8	163	0.2313	38	8	49	0.2490	12
9	163	0.2454	40	9	49	0.2631	13
10	163	0.2454	40	10	49	0.2040	10
11	163	0.2089	34	11	49	0.1836	9
12	163	0.2089	34	12	49	0.1632	8
13	163	0.2089	34	13	49	0.1632	8
14	163	0.2142	35	14	49	0.1632	8
15	163	0.2089	34	15	49	0.1632	8
16	163	0.2045	33	16	49	0.1632	8
17	163	0.2089	34	17	49	0.1632	8
18	163	0.2045	33	18	49	0.1632	8
19	163	0.2089	34	19	49	0.1632	8
20	163	0.2089	34	20	49	0.1632	8
21	163	0.2089	34	21	49	0.1632	8
22	163	0.2089	34	22	49	0.1632	8
23	163	0.2089	34	23	49	0.1632	8
24	163	0.2089	34	24	49	0.1632	8
25	163	0.2089	34	25	49	0.1632	8
26	163	0.2089	34	26	49	0.1632	8
27	163	0.2089	34	27	49	0.1632	8
28	163	0.2089	34	28	49	0.1632	8
29	163	0.2089	34	29	49	0.1632	8
30	163	0.2089	34	30	49	0.1632	8
31	163	0.2089	34	31	49	0.1632	8
32	163	0.2089	34	32	49	0.1632	8
33	163	0.2089	34	33	49	0.1632	8
34	163	0.2089	34	34	49	0.1632	8
35	163	0.2089	34	35	49	0.1632	8
36	163	0.2089	34	36	49	0.1632	8
37	163	0.2089	34	37	49	0.1632	8
38	163	0.2089	34	38	49	0.1632	8
39	163	0.2089	34	39	49	0.1632	8
40	163	0.2089	34	40	49	0.1632	8
41	163	0.2089	34	41	49	0.1632	8
42	163	0.2089	34	42	49	0.1632	8
43	163	0.2089	34	43	49	0.1632	8
44	163	0.2089	34	44	49	0.1632	8
45	163	0.2089	34	45	49	0.1632	8
46	163	0.2089	34	46	49	0.1632	8
47	163	0.2089	34	47	49	0.1632	8
48	163	0.2089	34	48	49	0.1632	8
49	163	0.2089	34	49	49	0.1632	8
50	163	0.2089	34	50	49	0.1632	8
51	163	0.2089	34	51	49	0.1632	8
52	163	0.2089	34	52	49	0.1632	8
53	163	0.2089	34	53	49	0.1632	8
54	163	0.2089	34	54	49	0.1632	8
55	163	0.2089	34	55	49	0.1632	8
56	163	0.2089	34	56	49	0.1632	8
57	163	0.2089	34	57	49	0.1632	8
58	163	0.2089	34	58	49	0.1632	8
59	163	0.2089	34	59	49	0.1632	8
60	163	0.2089	34	60	49	0.1632	8
61	163	0.2089	34	61	49	0.1632	8
62	163	0.2089	34	62	49	0.1632	8
63	163	0.2089	34	63	49	0.1632	8
64	163	0.2089	34	64	49	0.1632	8
65	163	0.2089	34	65	49	0.1632	8
66	163	0.2089	34	66	49	0.1632	8
67	163	0.2089	34	67	49	0.1632	8
68	163	0.2089	34	68	49	0.1632	8
69	163	0.2089	34	69	49	0.1632	8
70	163	0.2089	34	70	49	0.1632	8
71	163	0.2089	34	71	49	0.1632	8
72	163	0.2089	34	72	49	0.1632	8
73	163	0.2089	34	73	49	0.1632	8
74	163	0.2089	34	74	49	0.1632	8
75	163	0.2089	34	75	49	0.1632	8

Confusion Matrix for Best K = 12

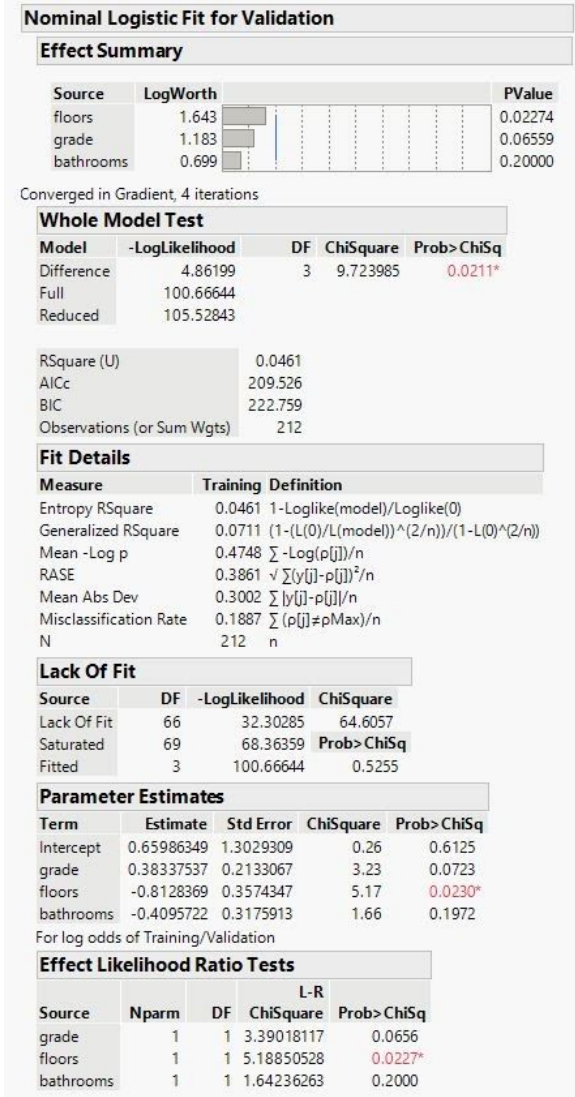
Training			Validation		
Actual	Predicted Count		Actual	Predicted Count	
Validation	Training	Validation	Validation	Training	Validation
Training	128	1	Training	41	0
Validation	33	1	Validation	8	0

Mosaic Plot

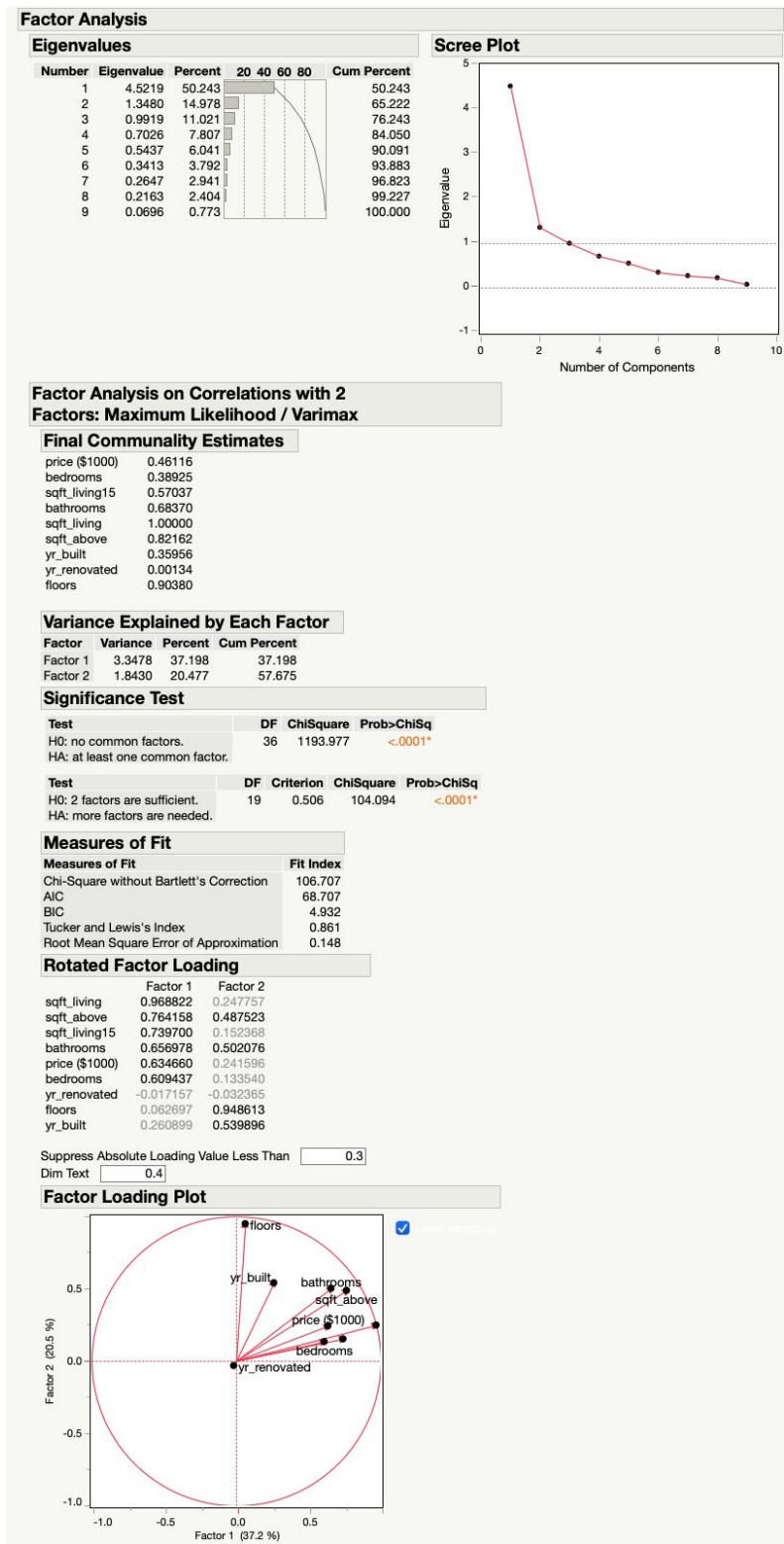
Mosaic Plot for K=12



Figure 23



Bonus Figure 1:



Bonus Figure 2:

Iterative Clustering

Cluster Comparison

Method	NCluster	CCC	Best
K Means Cluster	2	-2.267	Optimal CCC

Columns Scaled Individually

K Means NCluster=2

Columns Scaled Individually

Cluster Summary

Cluster	Count	Step	Criterion
1	75	9	0
2	137		

Cluster Means

Cluster	price (\$1000)	bedrooms	bathrooms	sqft_living15	sqft_living	sqft_above	yr_built	yr_renovated	floors
1	689.893333	3.96	2.73	2602.50667	2912.82667	2621.76	1988.61333	79.8133333	1.87333333
2	389.591241	3.13138686	1.73357664	1631.06569	1571.31387	1294.9635	1961.56934	72.4817518	1.3649635