

Applied Statistics Final Project

Jeffrey Dean, Ronald Tourtellot, Parth Patel

Report

For this portion of the project, we are tasked to review statistical data on the relationship between five predictor factors (floor, distance, view, end, and furnish) on price as a response variable for condominium offerings. Both significance testing and analysis of graphs will be utilized to determine whether models generated are good means of understanding regression and what conclusions to draw from this. For the remainder of the work the benchmark critical significance will be set at 5% as a standard practice.

Part 1- Scatterplots and Boxplots

For part 1, we need to provide graphs of the five models between each predictor and the response variable. Because three of the five predictors are Boolean figures as data (0,1), boxplots are required. For the continuous variables, a scatterplot is the preferred method of regression. These graphs are good starting position to understand the data provided in their relationship to a response variable, in this case, dollars. Below is the code and the output for the regression graphs.

```
model_1_price_floor <- lm(Applied_Statistics_Project_Condo_Data$PRICE100 ~
                          Applied_Statistics_Project_Condo_Data$FLOOR)
model_1_price_floor

model_1_plot <- plot(Applied_Statistics_Project_Condo_Data$FLOOR, Applied_Statistics_Project_Condo_Data$PRICE100,
                    main = "Model 1 Floor", xlab = "floor", ylab = "dollars")
abline(model_1_price_floor)
model_1_plot

model_2_price_distance <- lm(Applied_Statistics_Project_Condo_Data$PRICE100 ~
                             Applied_Statistics_Project_Condo_Data$DISTANCE)
model_2_price_distance

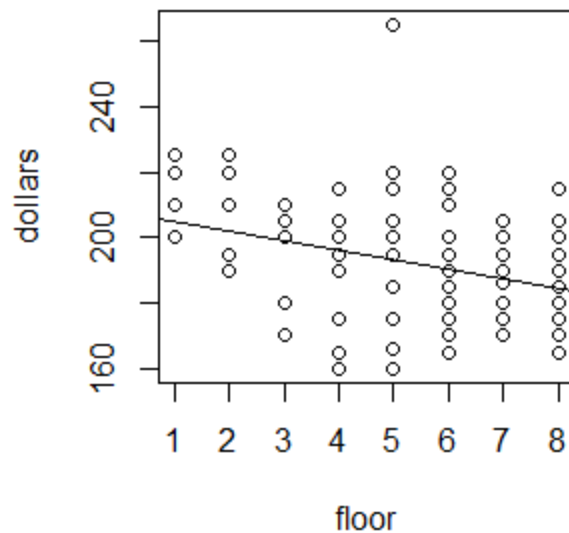
model_2_plot <- plot(Applied_Statistics_Project_Condo_Data$DISTANCE, Applied_Statistics_Project_Condo_Data$PRICE100,
                    main = "Model 2 Distance", xlab = "distance", ylab = "dollars")
abline(model_2_price_distance)
model_2_plot
```

```
lm(formula = Applied_Statistics_Project_Condo_Data$PRICE100 ~
    Applied_Statistics_Project_Condo_Data$FLOOR)
```

Coefficients:

(Intercept)	Applied_Statistics_Project_Condo_Data\$FLOOR
207.712	-2.881

Model 1 Floor

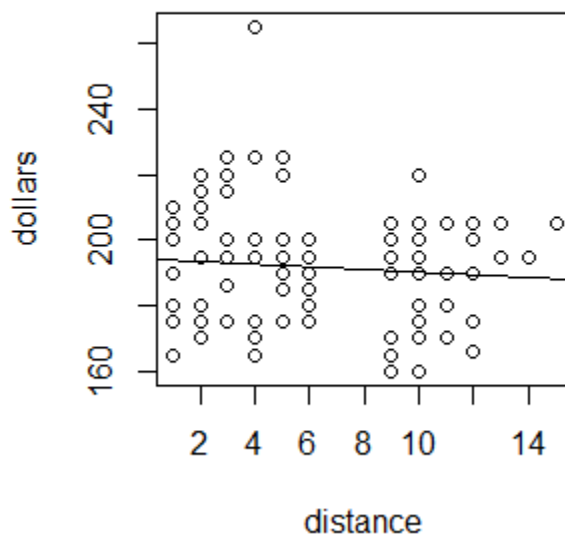


```
lm(formula = Applied_Statistics_Project_Condo_Data$PRICE100 ~  
    Applied_Statistics_Project_Condo_Data$DISTANCE)
```

Coefficients:

(Intercept)	Applied_Statistics_Project_Condo_Data\$DISTANCE
194.3488	-0.3915

Model 2 Distance



From the graphs, we can deduce a mildly linear relationship around the regression line. Without further analysis there is not much substance to understand between the data of the relationship. A few outliers are present, but otherwise this is a rough outline of what is occurring.

The next set of three regressors are better explained using boxplots as the x variable is denominated in 0 or 1. 0 in this case denotes an absence of a condition (view, end condo, or furnish) and 1 means the condition is present.

```
model_3_price_view <- lm(Applied_Statistics_Project_Condo_Data$PRICE100 ~ Applied_Statistics_Project_Condo_Data$VIEW)
model_3_price_view

model_3_plot <- boxplot(Applied_Statistics_Project_Condo_Data$PRICE100 ~ Applied_Statistics_Project_Condo_Data$VIEW)
model_3_plot

model_4_price_end <- lm(Applied_Statistics_Project_Condo_Data$PRICE100 ~ Applied_Statistics_Project_Condo_Data$END)
model_4_price_end

model_4_plot <- boxplot(Applied_Statistics_Project_Condo_Data$PRICE100 ~ Applied_Statistics_Project_Condo_Data$END)
model_4_plot
```

```
model_5_price_furnish<- lm(Applied_Statistics_Project_Condo_Data$PRICE100 ~
                           Applied_Statistics_Project_Condo_Data$FURNISH)
model_5_price_furnish

model_5_plot <- boxplot(Applied_Statistics_Project_Condo_Data$PRICE100 ~
                        Applied_Statistics_Project_Condo_Data$FURNISH)
model_5_plot
```

```
lm(formula = Applied_Statistics_Project_Condo_Data$PRICE100 ~
    Applied_Statistics_Project_Condo_Data$VIEW)
```

Coefficients:

(Intercept)	Applied_Statistics_Project_Condo_Data\$VIEW
173.94	27.06

```
lm(formula = Applied_Statistics_Project_Condo_Data$PRICE100 ~
    Applied_Statistics_Project_Condo_Data$END)
```

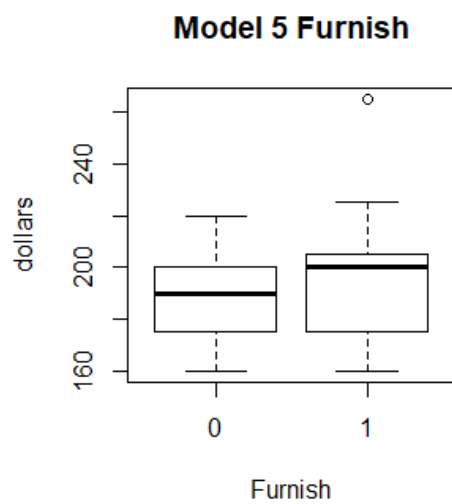
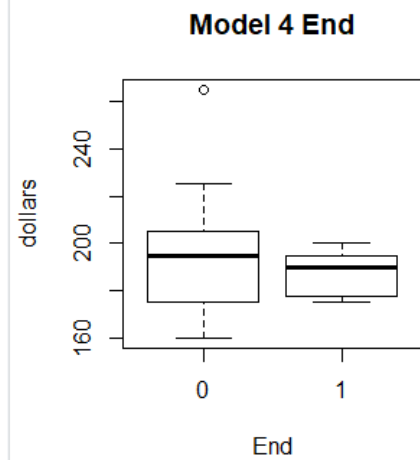
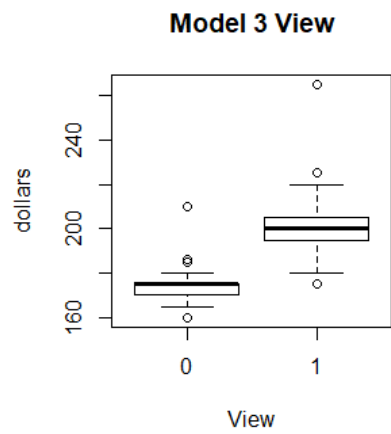
Coefficients:

(Intercept)	Applied_Statistics_Project_Condo_Data\$END
192.141	-4.999

```
lm(formula = Applied_Statistics_Project_Condo_Data$PRICE100 ~
    Applied_Statistics_Project_Condo_Data$FURNISH)
```

Coefficients:

(Intercept)	Applied_Statistics_Project_Condo_Data\$FURNISH
189.508	5.425



Each boxplot explains data in the clearest manner given their parameters. Model 3 represents if a condo has a view or not, showing a loose representation that a view creates more dollar value to a condo, then if did not have a view. With one outlier exception a view has an average value of 200 (in '000), while a non-view condo has a value of 170. There is a tighter distribution for no view condos than one with view, supporting there may be more extraneous factors for an apartment with a view. The distribution of an end condo (or not) tells a different distribution. Both boxplots are similar with "an end condo" distribution being tighter, with a slightly lower average dollar value. Finally, furnished (or not) distributions are very similar in both spacing and size. Furnishing has a slightly higher command of dollar value, accounting for the built-in cost of paying to furnish the rooms that normally is borne by the consumer. The cost is marginal, so it would be unlikely for the data to be too different in scope.

Linear Model

For part 2, we are tasked with developing a base model with each of the input predictor variables relative to dollar in price. We developed this model and reviewed the key indicators that were shown from the output.

```
#full model
Condo_linear_model_5pred<- lm(Applied_Statistics_Project_Condo_Data$PRICE100 ~
                             Applied_Statistics_Project_Condo_Data$FLOOR +
                             Applied_Statistics_Project_Condo_Data$DISTANCE +
                             Applied_Statistics_Project_Condo_Data$VIEW +
                             Applied_Statistics_Project_Condo_Data$END +
                             Applied_Statistics_Project_Condo_Data$FURNISH)
summary(Condo_linear_model_5pred)
plot(Condo_linear_model_5pred)

#term by term model
Condo_linear_model_5pred_plot <- plot(Applied_Statistics_Project_Condo_Data$PRICE100 ~
                                     Applied_Statistics_Project_Condo_Data$FLOOR +
                                     Applied_Statistics_Project_Condo_Data$DISTANCE +
                                     Applied_Statistics_Project_Condo_Data$VIEW +
                                     Applied_Statistics_Project_Condo_Data$END +
                                     Applied_Statistics_Project_Condo_Data$FURNISH)

abline(Condo_linear_model_5pred)
Condo_linear_model_5pred_plot

#residuals
Condo_linear_5pred_resid <- resid(Condo_linear_model_5pred)
Condo_linear_5pred_resid
plot(Condo_linear_5pred_resid)
```

Below is output found for the model using a linear model function in R. Looking at the range of t statistics, we can see that all the variables are significant on their own as part of the model with the notable except of the "floor" predictor. The model supports all of the significant factors should be included for the final model, except floor at this point. The R^2 is 70.58%, which means the regression line supports a majority of the regression output. F statistic is 47.98, which is significant at any level needed. Therefore, we should support using these predictors and this model further for our analysis.

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	177.7035	4.1684	42.631	< 2e-16 ***
Applied_Statistics_Project_Condo_Data\$FLOOR	-0.7151	0.5308	-1.347	0.18090
Applied_Statistics_Project_Condo_Data\$DISTANCE	-0.8733	0.2449	-3.565	0.00056 ***
Applied_Statistics_Project_Condo_Data\$VIEW	31.2728	2.2312	14.016	< 2e-16 ***
Applied_Statistics_Project_Condo_Data\$END	-17.8078	3.9820	-4.472	2.05e-05 ***
Applied_Statistics_Project_Condo_Data\$FURNISH	9.9838	2.0515	4.867	4.24e-06 ***

 signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 9.905 on 100 degrees of freedom
 Multiple R-squared: 0.7058, Adjusted R-squared: 0.6911
 F-statistic: 47.98 on 5 and 100 DF, p-value: < 2.2e-16

Quadratic Models

For part three it was important to assess whether any quadratic coefficients included in the model were significant to assess if the model had a quadratic component rather than a purely linear one. Adding a coefficient for the square of each continuous variable was done and a model for both terms being quadratic was also used as well. Significance in this portion would lead us to use the predictor variable in the final model.

```
#Quadratic 1 Floor
Condo_quadratic_model_floor<- lm(Applied_Statistics_Project_Condo_Data$PRICE100 ~
  Applied_Statistics_Project_Condo_Data$FLOOR +
  Applied_Statistics_Project_Condo_Data$DISTANCE +
  Applied_Statistics_Project_Condo_Data$VIEW
  + Applied_Statistics_Project_Condo_Data$END +
  Applied_Statistics_Project_Condo_Data$FURNISH +
  Applied_Statistics_Project_Condo_Data$FLOOR.SQ)

summary(Condo_quadratic_model_floor)
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	188.9570	6.2137	30.410	< 2e-16 ***
Applied_Statistics_Project_Condo_Data\$FLOOR	-6.3512	2.4068	-2.639	0.00966 **
Applied_Statistics_Project_Condo_Data\$DISTANCE	-0.8059	0.2410	-3.344	0.00117 **
Applied_Statistics_Project_Condo_Data\$VIEW	30.8095	2.1886	14.077	< 2e-16 ***
Applied_Statistics_Project_Condo_Data\$END	-18.8011	3.9126	-4.805	5.51e-06 ***
Applied_Statistics_Project_Condo_Data\$FURNISH	10.9033	2.0408	5.343	5.86e-07 ***
Applied_Statistics_Project_Condo_Data\$FLOOR.SQ	0.5650	0.2356	2.398	0.01836 *

 signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 9.677 on 99 degrees of freedom
 Multiple R-squared: 0.7219, Adjusted R-squared: 0.7051
 F-statistic: 42.84 on 6 and 99 DF, p-value: < 2.2e-16

Above is data for model 1- quadratic for the “floor” variable being used as a quadratic. In R we ran a model and found that at a 5% level, floor squared was significant, even though in the initial model, the “floor” variable was not significant. The F statistic was significant and the R² was in a reasonable position so the addition of this variable seemed appropriate.

```
#Quadratic 2 Distance
Condo_quadratic_model_distance<- lm(Applied_Statistics_Project_Condo_Data$PRICE100 ~
    Applied_Statistics_Project_Condo_Data$FLOOR +
    Applied_Statistics_Project_Condo_Data$DISTANCE +
    Applied_Statistics_Project_Condo_Data$VIEW
    + Applied_Statistics_Project_Condo_Data$END +
    Applied_Statistics_Project_Condo_Data$FURNISH +
    Applied_Statistics_Project_Condo_Data$DISTANCE.SQ)

summary(Condo_quadratic_model_distance)
```

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	181.13669	4.87889	37.127	< 2e-16	***
Applied_Statistics_Project_Condo_Data\$FLOOR	-0.75244	0.52940	-1.421	0.158368	
Applied_Statistics_Project_Condo_Data\$DISTANCE	-2.31088	1.10030	-2.100	0.038251	*
Applied_Statistics_Project_Condo_Data\$VIEW	30.98001	2.23311	13.873	< 2e-16	***
Applied_Statistics_Project_Condo_Data\$END	-16.02423	4.18362	-3.830	0.000225	***
Applied_Statistics_Project_Condo_Data\$FURNISH	10.29314	2.05639	5.005	2.43e-06	***
Applied_Statistics_Project_Condo_Data\$DISTANCE.SQ	0.10359	0.07731	1.340	0.183335	

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 9.865 on 99 degrees of freedom
Multiple R-squared: 0.711, Adjusted R-squared: 0.6935
F-statistic: 40.6 on 6 and 99 DF, p-value: < 2.2e-16

For output 2, “distance squared” was used for the quadratic factor added to the initial model. From the output above, we see the t statistic for “distance squared” is not significant at any level, even though distance by itself initially did show to be significant. This supports the variable had more value as a simple linear relationship and this is reflected in our final model to fit best.

```
#Quadratic 3 Floor and Distance
Condo_quadratic_model_floor_distance<- lm(Applied_Statistics_Project_Condo_Data$PRICE100 ~
    Applied_Statistics_Project_Condo_Data$FLOOR +
    Applied_Statistics_Project_Condo_Data$DISTANCE +
    Applied_Statistics_Project_Condo_Data$VIEW
    + Applied_Statistics_Project_Condo_Data$END +
    Applied_Statistics_Project_Condo_Data$FURNISH +
    Applied_Statistics_Project_Condo_Data$FLOOR.SQ +
    Applied_Statistics_Project_Condo_Data$DISTANCE.SQ)

summary(Condo_quadratic_model_floor_distance)
```

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	193.47806	6.80467	28.433	< 2e-16	***
Applied_Statistics_Project_Condo_Data\$FLOOR	-6.68650	2.39856	-2.788	0.00638	**
Applied_Statistics_Project_Condo_Data\$DISTANCE	-2.45099	1.07281	-2.285	0.02449	*
Applied_Statistics_Project_Condo_Data\$VIEW	30.44960	2.18447	13.939	< 2e-16	***
Applied_Statistics_Project_Condo_Data\$END	-16.80737	4.08536	-4.114	8.10e-05	***
Applied_Statistics_Project_Condo_Data\$FURNISH	11.30581	2.04185	5.537	2.58e-07	***
Applied_Statistics_Project_Condo_Data\$FLOOR.SQ	0.59430	0.23461	2.533	0.01289	*
Applied_Statistics_Project_Condo_Data\$DISTANCE.SQ	0.11879	0.07551	1.573	0.11893	

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 9.606 on 98 degrees of freedom
Multiple R-squared: 0.7288, Adjusted R-squared: 0.7094
F-statistic: 37.62 on 7 and 98 DF, p-value: < 2.2e-16

For the final portion of part 3 we wanted to assess the model where both “distance” and “floor” variables both had a quadratic component included to see how the output may change. The above

model is similar to the prior two in terms of both the significance of the F statistic and the calculated R^2 . It is also similar that the “Floor Squared” variable is significant in this case as well at the 5% level. This supports that is a good idea for us to add this portion to our final model for condo pricing.

Interaction Term Models

For the final portion we needed to examine the interaction terms and support what a good model should include to reflect our findings. It would be too pedantic to show every interaction term model we will provide an example of how one was conducted and the results for all of the test as an output. The remainder of the models not shown were conducted in the same way with different predictor inputs. Finally, we took the significant portions of the prior models to construct the best option for our final model of predictors vs. response variable.

For the portion of interaction terms, we cataloged the five predictors 1 to 5 and built a two-way interaction model output for each pairing. 5 variables in a two-way test yield 10 models, each testing the t statistic of the interaction term and reviewing the results.

```
Condo_Interaction_model_1_2 <- lm(Applied_Statistics_Project_Condo_Data$PRICE100 ~
                                Applied_Statistics_Project_Condo_Data$FLOOR +
                                Applied_Statistics_Project_Condo_Data$DISTANCE +
                                Applied_Statistics_Project_Condo_Data$VIEW
                                + Applied_Statistics_Project_Condo_Data$END + Applied_Statistics_Project_Condo_Data$FURNISH +
                                Applied_Statistics_Project_Condo_Data$FLOOR:Applied_Statistics_Project_Condo_Data$DISTANCE)
summary(Condo_Interaction_model_1_2)
```

Coefficients:

	Estimate	Std. Error	t value
(Intercept)	181.6401	6.1940	29.325
Applied_Statistics_Project_Condo_Data\$FLOOR	-1.4010	0.9582	-1.462
Applied_Statistics_Project_Condo_Data\$DISTANCE	-1.5457	0.8193	-1.887
Applied_Statistics_Project_Condo_Data\$VIEW	31.1402	2.2394	13.905
Applied_Statistics_Project_Condo_Data\$END	-17.6442	3.9917	-4.420
Applied_Statistics_Project_Condo_Data\$FURNISH	9.9749	2.0542	4.856
Applied_Statistics_Project_Condo_Data\$FLOOR:Applied_Statistics_Project_Condo_Data\$DISTANCE	0.1171	0.1361	0.860

	Pr(> t)
(Intercept)	< 2e-16 ***
Applied_Statistics_Project_Condo_Data\$FLOOR	0.1469
Applied_Statistics_Project_Condo_Data\$DISTANCE	0.0622 .
Applied_Statistics_Project_Condo_Data\$VIEW	< 2e-16 ***
Applied_Statistics_Project_Condo_Data\$END	2.53e-05 ***
Applied_Statistics_Project_Condo_Data\$FURNISH	4.49e-06 ***
Applied_Statistics_Project_Condo_Data\$FLOOR:Applied_Statistics_Project_Condo_Data\$DISTANCE	0.3918

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 9.917 on 99 degrees of freedom
Multiple R-squared: 0.708, Adjusted R-squared: 0.6903
F-statistic: 40 on 6 and 99 DF, p-value: < 2.2e-16

The output above includes each of the initial model predictors and one interaction term that modifies the condition of the model. In this case this was variable 1 and variable 2, “floor” and “distance”. The variable proved to be not significant at any level for the interaction, while the F statistic for the model used was significant. This supports we do not need to further investigate this interaction. Below is the

Excel spreadsheet of the findings if the process was iterated 9 more times in the same conditions for the remaining variable interactions. The findings found significance between variables 1 and 3 and 3 and 5. These are interactions for floor and view and view and furnishing. This makes sense as floor matters a lot with a good view as a low story condo doesn't matter if it has a good view since its low to the ground. Comparatively, a good view condo is probably furnished as condos with a good view usually are priced higher and the upsell to provided furnishing luxury has high margins for that range of clientele. Both variables 2 and 4 and 3 and 4 were inconclusive for the test run. The significance was very strong compared to the remainder of the test output so I suspect some degree of multicollinearity between these variables, so dropping their involvement should provide a better final model overall. This is a judgement call based on what was given in those outputs.

Final Project- Condo Interactions			5% Level Comparison					
Interaction Terms	Significance	Conclusion				Terms		
1:2	>0.1	Not significant				1 Floor		
1:3	<0.01	Very significant				2 Distance		
1:4	>0.1	Not significant				3 View		
1:5	>0.1	Not significant				4 End		
2:3	>0.1	Not significant				5 Furnish		
2:4	NA	Inconclusive						
2:5	>0.1	Not significant						
3:4	NA	Inconclusive						
3:5	<0.05	Significant						
4:5	>0.1	Not significant						
Good to Examine	1:3, 3:5							

Final Model

Given all of the work prior for the linear model, experimenting with the quadratic data, and data on the interaction terms it is possible to build a concise final model given what inputs are considered good factors of the predictors relative to the dollar value of the condos as a response variable. Below is the output run and the results developed from the model and our inference.

```
#4 - Final Model
Condo_Final_Model <- lm(Applied_Statistics_Project_Condo_Data$PRICE100 ~
  Applied_Statistics_Project_Condo_Data$FLOOR +
  Applied_Statistics_Project_Condo_Data$DISTANCE +
  Applied_Statistics_Project_Condo_Data$VIEW + Applied_Statistics_Project_Condo_Data$END
+ Applied_Statistics_Project_Condo_Data$FURNISH +
  Applied_Statistics_Project_Condo_Data$FLOOR.SQ +
  Applied_Statistics_Project_Condo_Data$FLOOR:Applied_Statistics_Project_Condo_Data$VIEW
+ Applied_Statistics_Project_Condo_Data$VIEW:Applied_Statistics_Project_Condo_Data$FURNISH)

summary(Condo_Final_Model)

Condo_Final_Model_resid <- resid(Condo_Final_Model)
Condo_Final_Model_resid
plot(Condo_Final_Model_resid)
```

Coefficients:

	Estimate	Std. Error	t value
(Intercept)	164.7126	12.1036	13.609
Applied_Statistics_Project_Condo_Data\$FLOOR	-0.4857	3.1599	-0.154
Applied_Statistics_Project_Condo_Data\$DISTANCE	-0.8154	0.2331	-3.498
Applied_Statistics_Project_Condo_Data\$VIEW	52.1034	10.3183	5.050
Applied_Statistics_Project_Condo_Data\$END	-19.0903	3.7718	-5.061
Applied_Statistics_Project_Condo_Data\$FURNISH	9.5861	3.5806	2.677
Applied_Statistics_Project_Condo_Data\$FLOOR.SQ	0.2845	0.2457	1.158
Applied_Statistics_Project_Condo_Data\$FLOOR:Applied_Statistics_Project_Condo_Data\$VIEW	-3.7059	1.4840	-2.497
Applied_Statistics_Project_Condo_Data\$VIEW:Applied_Statistics_Project_Condo_Data\$FURNISH	3.1196	4.2943	0.726

	Pr(> t)
(Intercept)	< 2e-16 ***
Applied_Statistics_Project_Condo_Data\$FLOOR	0.87817
Applied_Statistics_Project_Condo_Data\$DISTANCE	0.00071 ***
Applied_Statistics_Project_Condo_Data\$VIEW	2.07e-06 ***
Applied_Statistics_Project_Condo_Data\$END	1.98e-06 ***
Applied_Statistics_Project_Condo_Data\$FURNISH	0.00872 **
Applied_Statistics_Project_Condo_Data\$FLOOR.SQ	0.24964
Applied_Statistics_Project_Condo_Data\$FLOOR:Applied_Statistics_Project_Condo_Data\$VIEW	0.01420 *
Applied_Statistics_Project_Condo_Data\$VIEW:Applied_Statistics_Project_Condo_Data\$FURNISH	0.46932

Residual standard error: 9.325 on 97 degrees of freedom
Multiple R-squared: 0.747, Adjusted R-squared: 0.7262
F-statistic: 35.8 on 8 and 97 DF, p-value: < 2.2e-16

Provided is both the F table data and residual plot for the inputs in order to control for the risk of heteroskedasticity. The table uses each of the linear inputs because all were considered significant in their own right except the “floor” variable. “Floor” is included because the “Floor.Squared” variable was a significant component found earlier and it is good statistical practice to have the original variable included if the quadratic component is allowed to be involved. Otherwise since we found that the “Floor:View” and “View:Furnish” interactions were significant, then it would be a good idea to include it here as well. The results show an interesting depiction for the output. Four of the five predictor variables were significant (without floor) and the interaction between “View” and “Floor” was significant. This supports the distance, view, end, and furnish provide a valid explanation on their own of the price of a condo independently, each portion’s inclusion adding to the price of the real estate as piecemeal. This makes sense as the price of the condo goes up incrementally when features are provided to a tenant and appeal to potential customers in different ways. Floor is only significant when it interacts with view so the logical conclusion is that floor only matters to a potential tenant if the view is better or that a view only matters if the floor is high enough. It is a two-way interaction test. Tenants can’t enjoy and consume a good view on the lower floors and the higher floor patrons only volunteer to spend more if the view provides a value.

In looking at the model parameters we see an F-Statistic at around a 36, which is highly significant at any level, R^2 is about 75%, which useful as a gauge of the model. Below is the output for the residuals to reaffirm the distribution is normal and does not have a change in variance. The residuals plotted are standardized around 0, and with few outlier exceptions, provides a lower and upper bound that is parallel. This supports there is no or minimal shift in variance and further strengthens the model delivered.

Final Model Residual Plot

