**Assignment #2**

Nate Bitting

## Introduction

For this assignment, the objective is to leverage the Ames Housing dataset to develop several regression models using single and/or multiple predictor variables to predict SalePrice. The first step is to create two simple linear regression models using a single predictor variable and SalePrice as the response variable. We will then evaluate the goodness of fit and accuracy of the models in order to determine which model performs better. The second objective is to develop two multiple regression models using several predictor variables to predict SalePrice. We will then compare the results of each model to determine which is superior. Lastly, we will perform a log transformation on SalePrice to determine if it improves model performance.

In order to prepare the data, we performed several data cleansing steps to filter out data we did not want included in the sample used to train the regression models for this assignment. We first only included single family homes in residential zoning areas. We also only included homes with paved streets and those that had public utilities. Lastly, we performed a univariate procedure to identify the extreme observations so that we could exclude those outliers from our sample. In the next step before building our regression models, we created dummy variables for many of the categorical variables we were interested in using: Central Air, Fireplaces, Garage, Neighborhood, Kitchen Quality, Basement Quality, Exterior Quality, Fireplace Quality, Foundation type, and Month Sold.

## Results
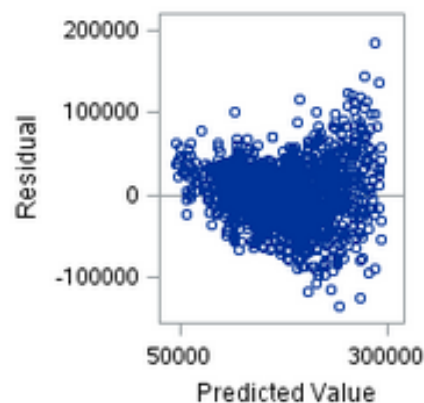
<p style="text-align:center;"><u>Simple Regression Model 1</u></p>

**Response Variable**: *SalePrice*     **Predictor Variable**: *total_SF*

**Equation**: *SalePrice = -19332 + 77.8244(total_SF)*

**Model 1 ANOVA Output**

**Model 1 Residual Plot**

The REG Procedure
Model: MODEL1
Dependent Variable: SalePrice

| Number of Observations Read | 1855 |
|---|---|
| Number of Observations Used | 1855 |

**Analysis of Variance**

| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
|---|---|---|---|---|---|
| Model | 1 | 4.392819E12 | 4.392819E12 | 4048.40 | <.0001 |
| Error | 1853 | 2.010645E12 | 1085075538 | | |
| Corrected Total | 1854 | 6.403464E12 | | | |

| Root MSE | 32940 | R-Square | 0.6860 |
|---|---|---|---|
| Dependent Mean | 170832 | Adj R-Sq | 0.6858 |
| Coeff Var | 19.28238 | | |

**Parameter Estimates**

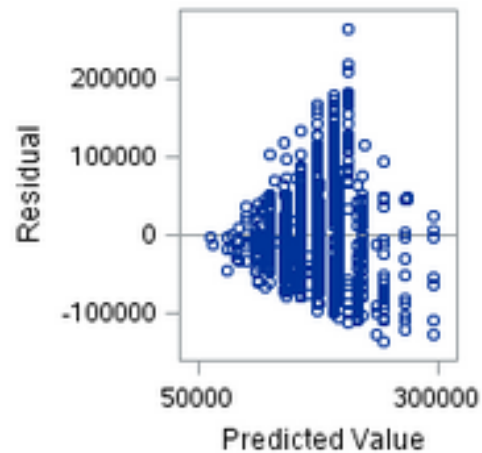| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > |t| |
|---|---|---|---|---|---|
| Intercept | 1 | -19332 | 3085.04498 | -6.27 | <.0001 |
| total_SF | 1 | 77.82440 | 1.22313 | 63.63 | <.0001 |

<u>Simple Regression Model 2</u>

**Response Variable**: *SalePrice*     **Predictor Variable**: *quality_index*
**Equation**: *SalePrice = -19332 + 77.8244(quality_index)*

**Model 2 ANOVA Output**

The REG Procedure
Model: MODEL1
Dependent Variable: SalePrice

| Number of Observations Read | 1855 |
|---|---|
| Number of Observations Used | 1855 |

Analysis of Variance

| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
|---|---|---|---|---|---|
| Model | 1 | 1.53233E12 | 1.53233E12 | 582.90 | <.0001 |
| Error | 1853 | 4.871134E12 | 2628782389 | | |
| Corrected Total | 1854 | 6.403464E12 | | | |

| | | | |
|---|---|---|---|
| Root MSE | 51272 | R-Square | 0.2393 |
| Dependent Mean | 170832 | Adj R-Sq | 0.2389 |
| Coeff Var | 30.01290 | | |

Parameter Estimates

| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > |t| |
|---|---|---|---|---|---|
| Intercept | 1 | 60185 | 4735.00274 | 12.71 | <.0001 |
| quality_index | 1 | 3261.05306 | 135.07005 | 24.14 | <.0001 |

**Model 2 Residual Plot**



By observing the ANOVA output of both Model 1 and Model 2, you can determine goodness of fit by reviewing the coefficient of determination (R-Square) of both models. Model 1 and Model 2 have an R-Square of 0.686 and 0.2303, respectively. Given Model 1 has a higher R-Square value and the MSE is much lower than Model 2, it is safe to say that Model 1 has a better fit than Model 2. The residual output also give us some insight to the model performance as well. By looking at the *Model 1 Residual Plot*, you will notice the residuals are fairly homoscedastic whereas the *Model 2 Residual Plot* has a fanning out pattern as the residuals increase as the quality_index gets larger. Thus, we can conclude that Model 1 performs better than Model 2.
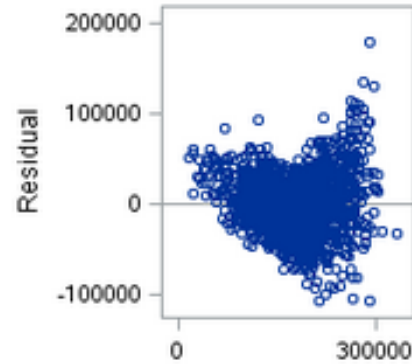
<p style="text-align:center"><u>Multiple Regression Model</u></p>

**Response Variable**: *SalePrice*     **Predictor Variables**: *total_SF, quality_index*
**Equation**: *SalePrice = -55909 + 70.31296(total_SF) + 1618.95131(quality_index)*

**Multiple Regression Model ANOVA Output**

The REG Procedure
Model: MODEL1
Dependent Variable: SalePrice

| Number of Observations Read | 1855 |
|---|---|
| Number of Observations Used | 1855 |

Analysis of Variance

| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
|---|---|---|---|---|---|
| Model | 2 | 4.729559E12 | 2.36478E12 | 2616.38 | <.0001 |
| Error | 1852 | 1.673904E12 | 903836072 | | |
| Corrected Total | 1854 | 6.403464E12 | | | |

| | | | |
|---|---|---|---|
| Root MSE | 30064 | R-Square | 0.7386 |
| Dependent Mean | 170832 | Adj R-Sq | 0.7383 |
| Coeff Var | 17.59849 | | |

Parameter Estimates

| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > |t| |
|---|---|---|---|---|---|
| Intercept | 1 | -55909 | 3393.91872 | -16.47 | <.0001 |
| total_SF | 1 | 70.31296 | 1.18221 | 59.48 | <.0001 |
| quality_index | 1 | 1618.95131 | 83.87469 | 19.30 | <.0001 |

**Multiple Regression Model Residual Plot**



Based on the results in the *Multiple Regression Model ANOVA Output* we can see that the MSE reduced significantly from Model 1 and Model 2 and the R-Square is much higher at 0.7386.  As a result of these findings, the Multiple Regression Model has a much better goodness of fit than both Model 1 and Model 2.

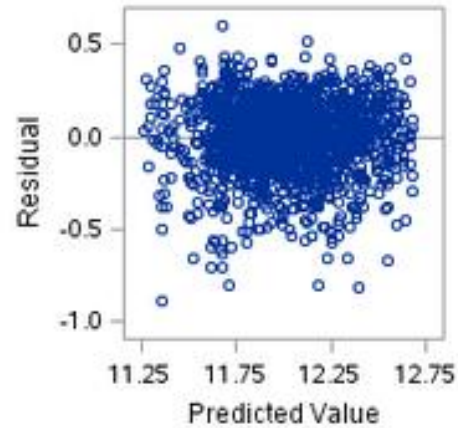## Regression Model 1 – Log Transformation on SalePrice

**Response Variable**: *log(SalePrice)*      **Predictor Variable**: *total_SF*

**Equation**: *log(SalePrice) = 10.91385 + 0.00044163(total_SF)*

### Model 1 Log Transformation ANOVA Output

The REG Procedure
Model: MODEL1
Dependent Variable: log_price

| Number of Observations Read | 1855 |
|---|---|
| Number of Observations Used | 1855 |

**Analysis of Variance**

| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
|---|---|---|---|---|---|
| Model | 1 | 141.45747 | 141.45747 | 4105.56 | <.0001 |
| Error | 1853 | 63.84525 | 0.03446 | | |
| Corrected Total | 1854 | 205.30271 | | | |

| | | | |
|---|---|---|---|
| Root MSE | 0.18562 | R-Square | 0.6890 |
| Dependent Mean | 11.99297 | Adj R-Sq | 0.6889 |
| Coeff Var | 1.54775 | | |

**Parameter Estimates**

| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > |t| |
|---|---|---|---|---|---|
| Intercept | 1 | 10.91385 | 0.01738 | 627.80 | <.0001 |
| total_SF | 1 | 0.00044163 | 0.00000689 | 64.07 | <.0001 |

### Model 1 Log Transformation Residual Plot
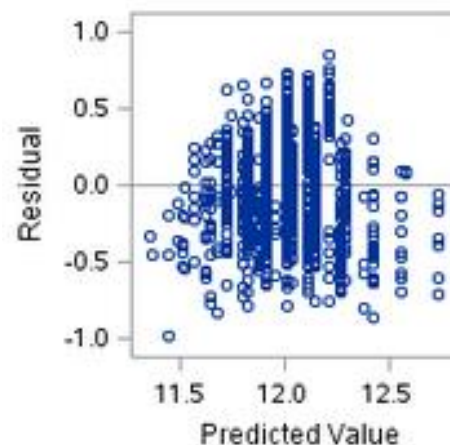


## Regression Model 2 – Log Transformation on SalePrice

**Response Variable**: *log(SalePrice)*      **Predictor Variable**: *quality_index*

**Equation**: *log(SalePrice) = -11.32999 + 0.01954(quality_index)*

### Model 2 Log Transformation ANOVA Output

The REG Procedure
Model: MODEL1
Dependent Variable: log_price

| Number of Observations Read | 1855 |
|---|---|
| Number of Observations Used | 1855 |

**Analysis of Variance**

| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
|---|---|---|---|---|---|
| Model | 1 | 55.01426 | 55.01426 | 678.31 | <.0001 |
| Error | 1853 | 150.28846 | 0.08111 | | |
| Corrected Total | 1854 | 205.30271 | | | |

| | | | |
|---|---|---|---|
| Root MSE | 0.28479 | R-Square | 0.2680 |
| Dependent Mean | 11.99297 | Adj R-Sq | 0.2676 |
| Coeff Var | 2.37464 | | |

**Parameter Estimates**

| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > |t| |
|---|---|---|---|---|---|
| Intercept | 1 | 11.32999 | 0.02630 | 430.79 | <.0001 |
| quality_index | 1 | 0.01954 | 0.00075025 | 26.04 | <.0001 |

### Model 2 Log Transformation Residual Plot

## Multiple Regression Model – Log Transformation on SalePrice
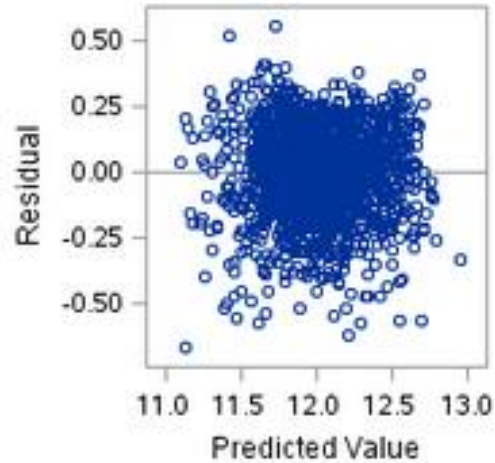
**Response Variable**: *log(SalePrice)*        **Predictor Variables**: *total_SF, quality_index*
**Equation**: *log(SalePrice) = 10.68008 + 0.00039362(total_SF) + 0.01035(quality_index)*

**Multiple Regression Model ANOVA Output**

The REG Procedure
Model: MODEL1
Dependent Variable: log_price

| Number of Observations Read | 1855 |
|---|---|
| Number of Observations Used | 1855 |

**Analysis of Variance**

| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
|---|---|---|---|---|---|
| Model | 2 | 155.21251 | 77.60626 | 2869.36 | <.0001 |
| Error | 1852 | 50.09020 | 0.02705 | | |
| Corrected Total | 1854 | 205.30271 | | | |

| Root MSE | 0.16446 | R-Square | 0.7560 |
|---|---|---|---|
| Dependent Mean | 11.99297 | Adj R-Sq | 0.7558 |
| Coeff Var | 1.37129 | | |

**Parameter Estimates**

| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > \|t\| |
|---|---|---|---|---|---|
| Intercept | 1 | 10.68008 | 0.01857 | 575.26 | <.0001 |
| total_SF | 1 | 0.00039362 | 0.00000647 | 60.87 | <.0001 |
| quality_index | 1 | 0.01035 | 0.00045882 | 22.55 | <.0001 |

**Multiple Regression Model Residual Plot**



After rebuilding the models leveraging the log transformation of the SalePrice for all three models we can leverage the ANOVA output and residuals to determine if these new models perform better than the original models.  For Model 1 and Model 2, we can see there was no improvement in their R-Square values as they remained exactly the same before the transformation on SalePrice.  However, after the log transformation on SalePrice, the Multiple Regression Model showed some improvement its R-Square value going from 0.7386 up to 0.7560.

## Conclusion

In conclusion, performing the log transformation on SalePrice in the Multiple Regression Model had the best goodness of fit than all other models.  Performing the log transformations on the simple linear regression models had no impact on model performance or fit than the original models.  By leveraging the ANOVA output from REG procedure, comparison of multiple models is a very straightforward endeavor given the limited number of regressors used in this assignment.  As you begin to add more predictor variables to your model, you can continue to tweak and fine tune the models for optimal performance.

**Code:**

```sas
1  * Code used to get the data into my library;
2  ods graphics on;
3  libname mydata '/courses/d6fc9ae5ba27fe300/c_3505/SAS_Data/' access=readonly;
4  proc datasets library=mydata; run; quit;
5  proc print data=mydata.anscombe; run; quit;
6
7
8  Data building;
9      SET mydata.ames_housing_data;
10
11     * filter on only single family homes that meet specific criteria;
12     if (SaleCondition = 'Normal');
13     if (BldgType = '1Fam'); * single family homes only;
14     if (Zoning in ('RH','RL','RP','RM')); * residential zones only;
15     if (Street='Pave'); * paved streets only;
16     if (Utilities='AllPub'); * only homes with public utilities;
17
18     log_price = log(SalePrice); *create a variable for the natural log of SalePrice;
19
20     * create new variables by combining multiple variables in the housing dataset;
21     total_SF = max(GrLivArea,0) + max(TotalBsmtSF,0);
22     total_baths = max(FullBath,0) + max(BsmtFullBath,0);
23     total_halfbaths = max(HalfBath,0) + max(BsmtHalfBath,0);
24     total_baths_calc = total_baths + total_halfbaths;
25
26     * filter on typical homes only by removing outliers;
27     if (total_SF < 4000 and total_SF > 800);
28
29     * Neighborhood dummy variables;
30     if (Neighborhood = 'Blmngtn') then Blmngtn=1; else Blmngtn=0;
31     if (Neighborhood = 'Blueste') then Blueste=1; else Blueste=0;
32     if (Neighborhood = 'BrDale') then BrDale=1; else BrDale=0;
33     if (Neighborhood = 'BrkSide') then BrkSide=1; else BrkSide=0;
34     if (Neighborhood = 'ClearCr') then ClearCr=1; else ClearCr=0;
35     if (Neighborhood = 'CollgCr') then CollgCr=1; else CollgCr=0;
36     if (Neighborhood = 'Crawfor') then Crawfor=1; else Crawfor=0;
37     if (Neighborhood = 'Edwards') then Edwards=1; else Edwards=0;
38     if (Neighborhood = 'Gilbert') then Gilbert=1; else Gilbert=0;
39     if (Neighborhood = 'Greens') then Greens=1; else Greens=0;
40     if (Neighborhood = 'GrnHill') then GrnHill=1; else GrnHill=0;
41     if (Neighborhood = 'IDOTRR') then IDOTRR=1; else IDOTRR=0;
42     if (Neighborhood = 'Landmrk') then Landmrk=1; else Landmrk=0;
43     if (Neighborhood = 'MeadowV') then MeadowV=1; else MeadowV=0;
44     if (Neighborhood = 'Mitchel') then Mitchel=1; else Mitchel=0;
45     if (Neighborhood = 'NAmes') then NAmes=1; else NAmes=0;
46     if (Neighborhood = 'NPkVill') then NPkVill=1; else NPkVill=0;
47     if (Neighborhood = 'NWAmes') then NWAmes=1; else NWAmes=0;
48     if (Neighborhood = 'NoRidge') then NoRidge=1; else NoRidge=0;
49     if (Neighborhood = 'NridgHt') then NridgHt=1; else NridgHt=0;
50     if (Neighborhood = 'OldTown') then OldTown=1; else OldTown=0;
51     if (Neighborhood = 'SWISU') then SWISU=1; else SWISU=0;
52     if (Neighborhood = 'Sawyer') then Sawyer=1; else Sawyer=0;
53     if (Neighborhood = 'SawyerW') then SawyerW=1; else SawyerW=0;
54     if (Neighborhood = 'Somerst') then Somerst=1; else Somerst=0;
55     if (Neighborhood = 'StoneBr') then StoneBr=1; else StoneBr=0;
56     if (Neighborhood = 'Timber') then Timber=1; else Timber=0;
57     if (Neighborhood = 'Veenker') then Veenker=1; else Veenker=0;
58
```

**Code Continued:**

```sas
58
59      * KitchenQual dummy variable;
60      if (KitchenQual in ('Ex', 'Gd')) then good_kitchen=1; else good_kitchen=0;
61
62      * FireplaceQu dummy variable;
63      if (FireplaceQu in ('Ex', 'Gd')) then good_fireplace=1; else good_fireplace=0;
64
65      * ExterQual dummy variable;
66      if (ExterQual in ('Ex', 'Gd')) then good_exterior=1; else good_exterior=0;
67
68      * Foundation dummy variables;
69      if (Foundation = 'BrkTil') then Foundation_BrkTil=1; else Foundation_BrkTil=0;
70      if (Foundation = 'CBlock') then Foundation_CBlock=1; else Foundation_CBlock=0;
71      if (Foundation = 'PConc') then Foundation_PConc=1; else Foundation_PConc=0;
72      if (Foundation = 'Slab') then Foundation_Slab=1; else Foundation_Slab=0;
73      if (Foundation = 'Stone') then Foundation_Stone=1; else Foundation_Stone=0;
74      if (Foundation = 'Wood') then Foundation_Wood=1; else Foundation_Wood=0;
75
76      * MoSold dummy variables;
77      if (MoSold = 1) then jan_sold=1; else jan_sold=0;
78      if (MoSold = 2) then feb_sold=1; else feb_sold=0;
79      if (MoSold = 3) then mar_sold=1; else mar_sold=0;
80      if (MoSold = 4) then apr_sold=1; else apr_sold=0;
81      if (MoSold = 5) then may_sold=1; else may_sold=0;
82      if (MoSold = 6) then jun_sold=1; else jun_sold=0;
83      if (MoSold = 7) then jul_sold=1; else jul_sold=0;
84      if (MoSold = 8) then aug_sold=1; else aug_sold=0;
85      if (MoSold = 9) then sep_sold=1; else sep_sold=0;
86      if (MoSold = 10) then oct_sold=1; else oct_sold=0;
87      if (MoSold = 11) then nov_sold=1; else nov_sold=0;
88      if (MoSold = 12) then dec_sold=1; else dec_sold=0;
89
90      * Construct a composite quality index;
91      quality_index = OverallCond*OverallQual;
92
93      * Central Air Indicator;
94      if (CentralAir='Y') then central_air=1; else central_air=0;
95      * Fireplace Indicator;
96      if (Fireplaces>0) then fireplace_ind=1; else fireplace_ind=0;
97      * Garage Indicator;
98      if (GarageCars>0) then garage_ind=1; else garage_ind=0;
99      * Good Basement Indicator;
100     if (BsmtQual in ('Ex','Gd')) or (BsmtCond in ('Ex','Gd')) then good_basement_ind=1; else good_basement_ind=0;
101
102 run; quit;
103
104 * Review new dataset to ensure no missing values based on filtering;
105 proc contents data=building;
106 run; quit;
107
108 * figure out what a typical home's square footage is in Ames;
109 proc univariate data=building;
110     var total_SF;
111 run; quit;
112
113 * 1st simple regression model;
114 proc reg data=building;
115     model SalePrice = total_SF;
116 run; quit;
```

**Code Continued:**

```sas
118  * 2nd simple regression model;
119  proc reg data=building;
120      model SalePrice = quality_index;
121  run; quit;
122
123  * Multiple Regression Model;
124  proc reg data=building;
125      model SalePrice = total_SF quality_index;
126  run; quit;
127
128  * 1st simple regression model with log transformation on SalePrice;
129  proc reg data=building;
130      model log_price = total_SF;
131  run; quit;
132
133  * 2nd simple regression model with log transformation on SalePrice;
134  proc reg data=building;
135      model log_price = quality_index;
136  run; quit;
137
138  * Multiple Regression Model with log transformation on SalePrice;
139  proc reg data=building;
140      model log_price = total_SF quality_index;
141  run; quit;
```