

Predicting a Movie's Revenue  
Melissa Marsh and Kai Hartley  
April 17, 2020

**Introduction:**

The oldest surviving film is *Roundhay Garden Scene* from 1888; a 2.11 second long film showing the director's family walking in a garden. Film has come a long way from where it started over 100 years ago. When it first started, it was mostly for the elite who could afford to attend. Nowadays, going to the movies is a popular pastime across the world and, with increased interest comes increased money.

Movies can reach larger audiences than ever before, and producers are continually trying to create new movies that will generate a large amount of revenue. The question remains, can you predict how well a movie will be received before it is released? This compilation of work will attempt to answer the question a hand and use linear regression to predict a movie's revenue.

**Initial Model:**

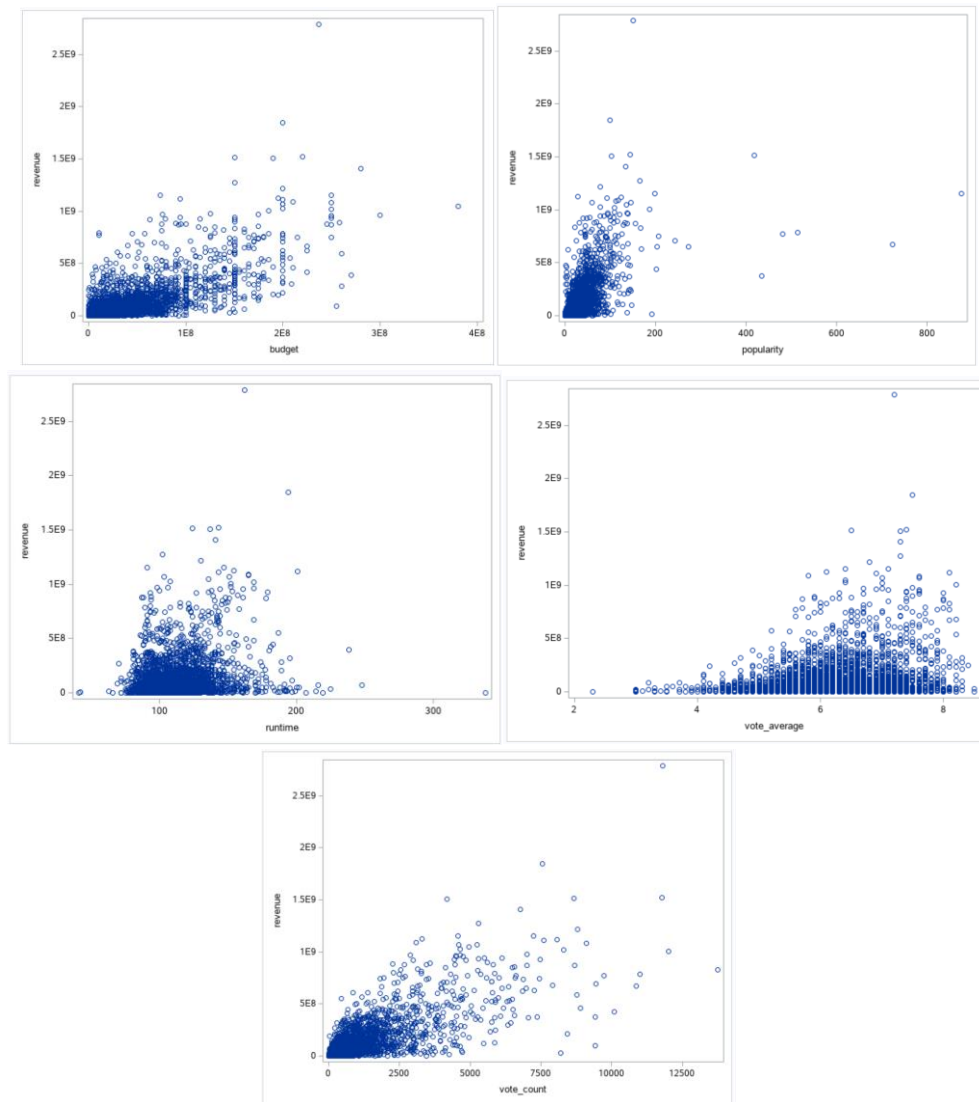
The dataset for our modelling was collected from TMDB. In this dataset, we are trying to determine the effects of multiple prediction variables on a movie's revenue. The full model includes eight prediction variables. All variables are described in *Table 1* below.

Variable	Description
Revenue	Amount of revenue generated by movie in US dollars adjusted to 2010 inflation
Budget	Budget of movie in US dollars adjusted to 2010 inflation
Popularity	Popularity score determined by TMDB
Runtime	Length of movie in minutes
Vote Average	Average voting score given by TMDB users on a scale from 1-10
Vote Count	Number of votes received by TMDB users
Genre	Primary genre of the movie (1=Action, 2=Adventure, 3=Animation, 4=Comedy, 5=Crime, 6=Documentary, 7=Drama, 8=Family, 9=Fantasy, 10=Foreign, 11=History, 12=Horror, 13=Music, 14=Mystery, 15=Romance, 16=Science Fiction, 17=Thriller, 18=War, 19=Western)
English	Whether or not the original language for the movie was English (1=English 0=not English)
US Production	Whether or not the country of the primary production company was the United States (1=US production, 0=not US production)

*Table 1: Description of variables in TMDB dataset.*

Before fitting a model using OLS regression, we first examined the distribution of the data by creating the scatterplots of revenue vs. the quantitative variables budget, popularity, runtime, vote average, and vote count. These scatterplots are shown in *Figure 1*. The scatterplots show a potential linear relationship between the quantitative variables and revenue; however, the distribution of these data points is not even. We have some points that appear to look like potential influential points or outliers. Although these points exist, their influence can be reduced through a variety of techniques.

We also examined histograms and boxplots to determine the distribution of the data. Not surprisingly, the histograms for the quantitative variables showed non-normality for all variables except for vote average which was distributed normally. Additionally, the variables budget, popularity and vote count had prominent right skew. The boxplots showed a handful of potential outliers.

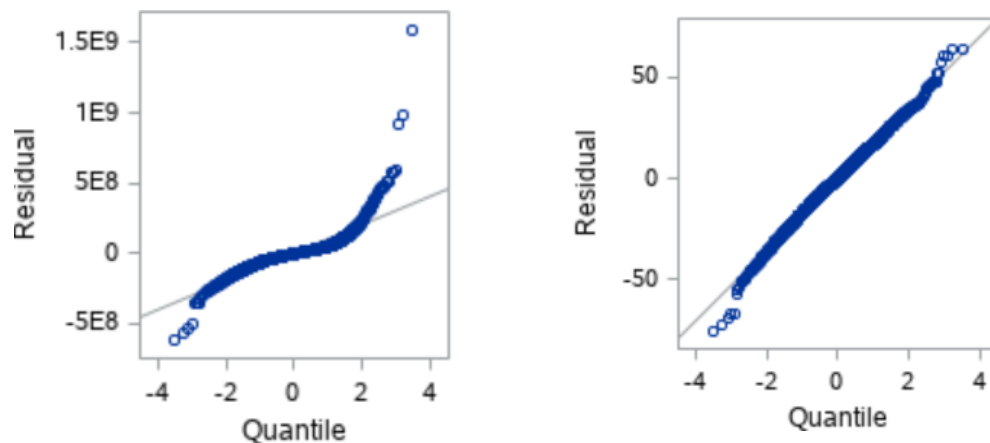


*Figure 1: Scatterplots of Revenue vs. Quantitative Variables.*

Despite the challenges presented with the original data set, we decided to fit a model using OLS regression. This is due to the fact that the quantitative variables do more or less show a linear relationship with revenue and the issues presented by the outlier points and distribution of the data may be fixed through statistical techniques.

Before fitting the model, we needed to “clean-up” the data. Some values had not transferred over correctly, which was stated in the original data’s description, resulting in values equaling zero. This was nonsensical in some cases such as budget equaling zero. These values were removed before fitting the model. Additionally, we split our data into a “train” group and a “test” group so that we could test how well a reduced model would perform on unseen data when compared to the full model.

The initial model did not fit model assumptions, as it was non-normally distributed and had nonconstant variance. To fix this, we examined transformations for revenue in addition to some of the explanatory variables. We ended up performing a fourth root transformation on revenue, a cubed root transformation on budget, a cubed root transformation on popularity and a fourth root transformation on vote count. Other transformations were also attempted but they did not significantly improve the model enough to warrant their inclusion. Shown in *Figure 2* is the QQ plot before and after transformation. Before transformation, the model did not fit the line well and was nonnormally distributed. After transformation, the model more closely fits the line with though the points deviate more towards the ends.



*Figure 2: QQ Plots of Initial Model (left) and Transformed Model (right).*

After these transformations were completed, our model fit assumptions of normality shown above and constant variance. The BF test had a p-value of  $2.8502E-75$  before transformations and the BF test had a p-value of  $0.099744$  after transformation. Included below in *Figure 3* are the residual plots before and after transformations.

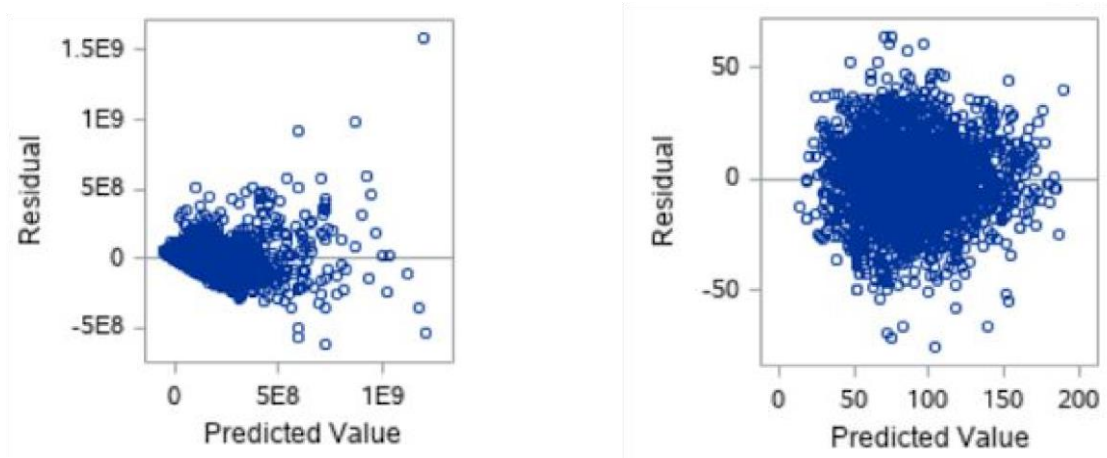


Figure 3: Residual Plots of Initial Model (left) and Transformed Model (right).

Our model now fits the assumptions required for OLS regression, and therefore we can investigate variable selection to reduce the size of our model. It is vital to note that our model does have some issues with multicollinearity due to the variance inflation factors being greater than zero. This is not surprising due to the nature of some of our variables being closely related. Although this may not be ideal, multicollinearity does not cause problems with predictive ability of our model - it only inflates the beta coefficients making them less interpretable. Due to the full model presenting possibly problematic multicollinearity, when choosing a reduced model, we were conscious to reduce the amount of multicollinearity in the quantitative variables. The variance inflation factors for each variable are shown in Table 2 below.

Variable	Variance Inflation Factor
Budget <sup>1/3</sup>	2.15233
Popularity <sup>1/3</sup>	6.29899
Runtime	1.55886
Vote_Average	2.00431
Vote_Count <sup>1/4</sup>	7.24806
Genre Action	24.04704
Genre Adventure	13.67903
Genre Animation	5.57303
Genre Comedy	25.15483
Genre Crime	7.39798
Genre Documentary	2.00785
Genre Drama	27.85040
Genre Family	2.67122
Genre Fantasy	5.51754
Genre Foreign	1.06352

Genre History	1.88086
Genre Horror	9.46038
Genre Music	1.99833
Genre Mystery	2.17818
Genre Romance	4.34440
Genre Science Fiction	4.26034
Genre Thriller	6.27115
Genre War	1.82280
Genre Western	
Original Language English	1.14340
Original Language Not English	
Not US Production	1.12531
US Production	

*Table 2: Variance Inflation Factors of Full Model Variables.*

### **Reduced Model:**

To create a reduced model, we used stepwise selection and backwards elimination techniques. We looked at different inclusion levels when stepwise was 0.05 and 0.01. Backwards elimination had exclusion values of 0.05 and 0.01. Stepwise selection and backwards elimination on both 0.01 and 0.05 levels found that budget, runtime, vote count and non US production all were significant so these variables were added to our reduced model. Additionally, backwards elimination and stepwise selection found that different genres were significant. Due to the difference in which genres were significant, and the fact that more than one genre was significant, it was decided to include all levels of genre into our reduced model.

The reduced model was fit using the variables found in *Table 3* below.

Variable	Description	Parameter Estimate	Variance Inflation
	Intercept	-9.88834	0
X <sub>1</sub>	Budget <sup>1/3</sup>	0.11918	1.6674
X <sub>2</sub>	Runtime	0.07743	1.34280
X <sub>3</sub>	Vote_Count <sup>1/4</sup>	10.99158	1.45894
X <sub>4</sub>	Genre Action	-0.78443	23.83631
X <sub>5</sub>	Genre Adventure	4.32999	13.62297
X <sub>6</sub>	Genre Animation	10.58384	5.57046
X <sub>7</sub>	Genre Comedy	4.80043	24.96929
X <sub>8</sub>	Genre Crime	-5.67921	7.37118
X <sub>9</sub>	Genre Documentary	8.09680	2.00059
X <sub>10</sub>	Genre Drama	-0.86569	27.76485

X <sub>11</sub>	Genre Family	10.47295	2.66984
X <sub>12</sub>	Genre Fantasy	2.57804	5.49062
X <sub>13</sub>	Genre Foreign	-2.95688	1.06090
X <sub>14</sub>	Genre History	5.26040	1.87773
X <sub>15</sub>	Genre Horror	5.79456	9.30486
X <sub>16</sub>	Genre Music	7.89712	1.99360
X <sub>17</sub>	Genre Mystery	-6.28039	2.17105
X <sub>18</sub>	Genre Romance	1.96881	4.32930
X <sub>19</sub>	Genre Science Fiction	-3.55806	4.23662
X <sub>20</sub>	Genre Thriller	-5.89554	6.21499
X <sub>21</sub>	Genre War	-8.00001	1.81897
X <sub>22</sub>	Genre Western		
X <sub>23</sub>	US_Production = 0	-3.92985	1.04176

Table 3: Variable Description, Parameter Estimate, and VIF for Reduced Model

After choosing which variables to include in our reduced model, we then examined our model to see if it fit model assumptions for OLS regression. *Figure 4* shows the QQ plot of the reduced model. We see that the QQ plot closely follows the line with a few deviations along the ends. This is similar to the QQ plot found for the full model and shows a normal distribution. Additionally, the correlation test of normality found that the model has a value of 0.99837 showing normality assumptions are met.

In addition to meeting assumptions regarding normality, the reduced model also fits assumptions of constant variance. The residual vs. predicted value plot is also shown in *Figure 4*. This plot shows that the residuals do have constant variance although some points fall outside of the large cloud of data. This is expected due to the large dataset. The Brown-Forsythe test confirms that our reduced model does have constant variance with a p-value of 0.072. This p-value is close to the border of being significant, but this expected as large datasets tend to be on the edge of significance due to a few observations affecting the residual plot.

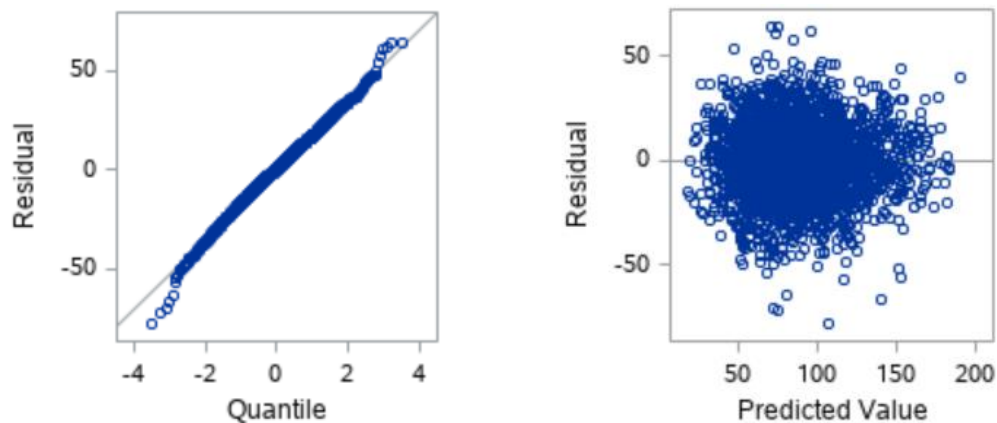


Figure 4: QQ plot for reduced model (left) and Residual plot for reduced model (right).

The final reduced model has the following equation:

$$\text{Revenue}^{1/4} = -9.888 + 0.119X_1 + 0.077X_2 + 10.992X_3 - 0.784X_4 + 4.330X_5 + 10.584X_6 + 4.800X_7 - 5.679X_8 + 8.097X_9 - 0.866X_{10} + 10.473X_{11} + 2.578X_{12} - 2.957X_{13} + 5.260X_{14} + 5.795X_{15} + 7.897X_{16} - 6.280X_{17} + 1.969X_{18} - 3.558X_{19} - 5.896X_{20} - 8.000X_{21} + 0X_{22} - 3.930X_{23}$$

In the model we chose to keep every genre in an effort to distinguish which kind of movie is most profitable. We interpret our equation so that any value shown in the reduced linear model with a positive multiplier suggests an increase of the fourth root of revenue, while a negative suggests a decrease. For example, for every unit increase in the runtime, we expect the fourth root of revenue of that movie to increase by 0.077 on average, holding all other variables constant. Also, when a movie is not a US production, we expect that it will decrease revenue by 3.930 on average, holding all other variables constant. For the genre, we consider that a movie can have only one main genre. Each score is based on Western being the dummy variable of the genre, which does not influence the revenue of a movie in our model. Due to the significant amount of multicollinearity associated with the genre variable, we cannot accurately state the affect that a specific genre will have on the fourth root of revenue.

Multicollinearity is present in our final model. The presence of multicollinearity in our model is mostly due to the genre variables having high variance inflation factors. The other quantitative variables – budget, runtime, and vote\_count had relatively low variance inflation values which had an average close to 1 as seen in *Table 3*. The variable describing that a movie was not a US production also had a low variance inflation factor close to one. While some of the genre values did have high variance inflation factors all genres were included in our model to make it more interpretable. If we only included the genres that did not have high variance inflation factors, our model would not be easy to interpret and due to the nature of the genre variable, multicollinearity was unavoidable. Even though multicollinearity is present, it does not affect the predictive power of our model. It only makes the coefficients related to genre difficult to interpret with their effect on revenue<sup>1/4</sup>.

When creating this reduced model, we examined possible interaction terms. We examined interaction terms with all the quantitative variables including the following:  $X_1 * X_2$ ,  $X_1 * X_3$ ,  $X_2 * X_3$ . From this analysis, we found that the interaction term of  $X_1 * X_3$  was significant and all other interaction terms were not significant. Due to this values significance, we examined the effect it had on our reduced model. We found that when the interaction term was included in our model, it resulted in a higher MSPE then when our model did not have the interaction term. For this reason, we chose not to include any of the interaction terms in our final model.

Influential points were examined in the reduced model by examining the Cook's D plot and DIFFITS shown in *Figure 5* and *Figure 6*. From these figures we can see that there are a good amount of observations that could be considered influential points. However, due to no

one point having a much larger Cook's D or DIFFITS than the others, and the fact that our data set is so large, there is not enough evidence to justify the removal of these influential points. Removing valid values due to their high influence is seen as a last-ditch effort and because our values are not too influential, this is not necessary.

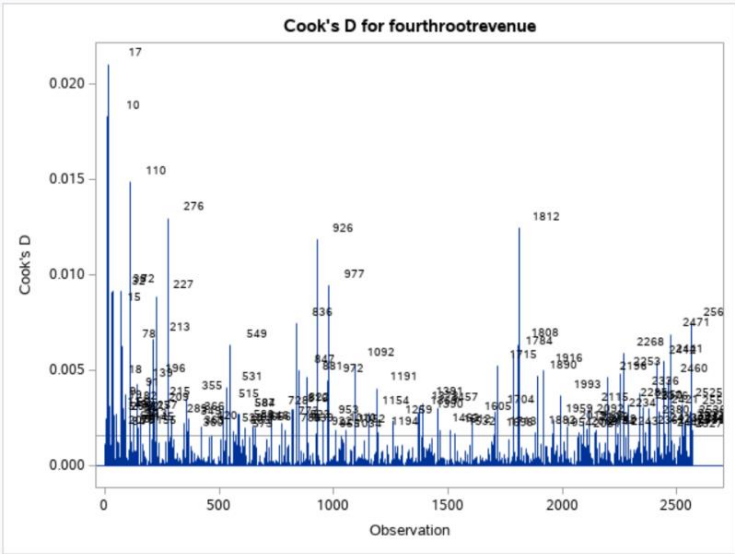


Figure 5: Cook's D plot for reduced model.

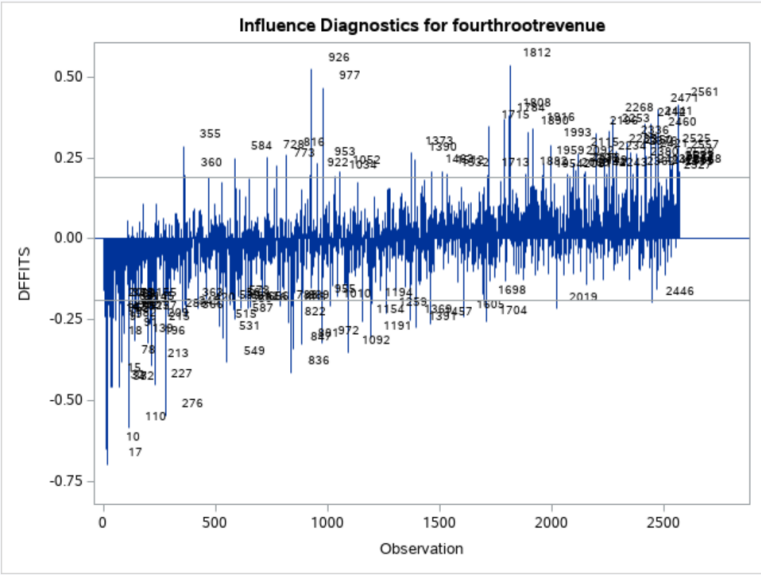


Figure 6: DIFFITS plot for reduced model.

Outliers in the reduced model were examined using the studentized residual plot shown in Figure 7. From this plot we can see that there is a good amount of observations that are



considered outliers and some that do show to have leverage on our model. Although, these points do exist, they are closely clustered so removal of one would not make sense unless we removed an entire cluster. Also, although these points exist, our model fits the requirements of OLS so we do not have a strong enough reason to remove them.

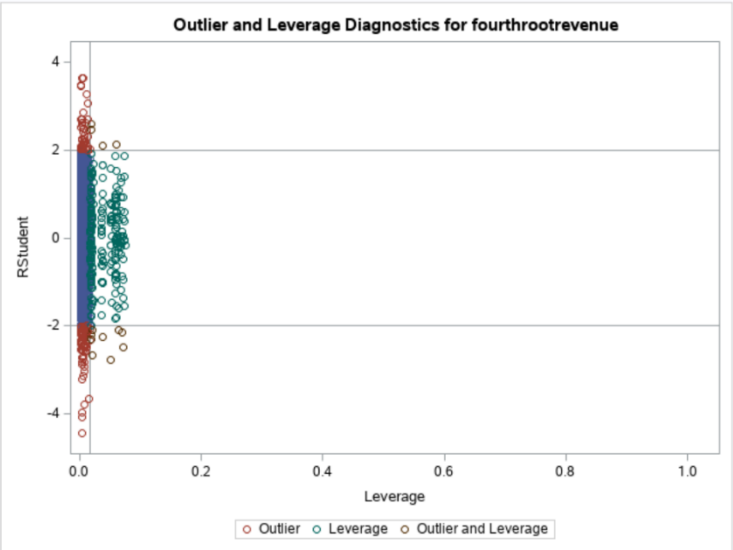


Figure 7: Studentized residual plot for reduced model.

The reduced model was compared to the full model to test how well it would predict new data. This was done by calculating MSPR values using the test set. We found the following results in Table 4 for the reduced model, full model, and null model (model that only included the intercept). We see that the reduced model has an MSPR value similar to the full model meaning that the predictive power is similar between the two models. Additionally, both the full model and the reduced model outperformed the null model when predicting on new data.

Model	Full Model	Null Model	Reduced Model
MSPR	345.53	1263.46	344.82

Table 4: MSPR for full, null, and reduced models.

Overall, the use of OLS on our dataset was able to look at our model in a structured setting that could easily be interpreted. However, it presented challenges when working with multicollinearity and we are only able to look at two-way interaction terms. To overcome these limitations, we decided to explore the use of a regression tree to model our dataset.

**Alternative Approach:**

As an alternative approach to OLS regression we looked at a regression tree of our data. Regression trees have the benefit of looking at high power interactions that are difficult to model

in OLS regression. We are interested in looking at how these high-power interactions can affect the predicted revenue of a movie, so we chose to explore this approach.

In our regression tree we decided to use the fourth root of revenue so that we could more easily compare the regression tree to our linear model, and we could work with smaller values of revenue. Our regression tree suggested 40 nodes; the subset regression tree is shown below in Figure 8.

Our regression tree allows us to look at high power interaction terms. We can see that when a movie has a vote count greater than 413, the budget becomes important in determining the revenue. Also, once we look at budget on the second node, we see that a movie with a budget less than \$72,200,000 depends on the budget again for determining revenue while if the budget is above 72,200,000 a movie depends on its vote count.

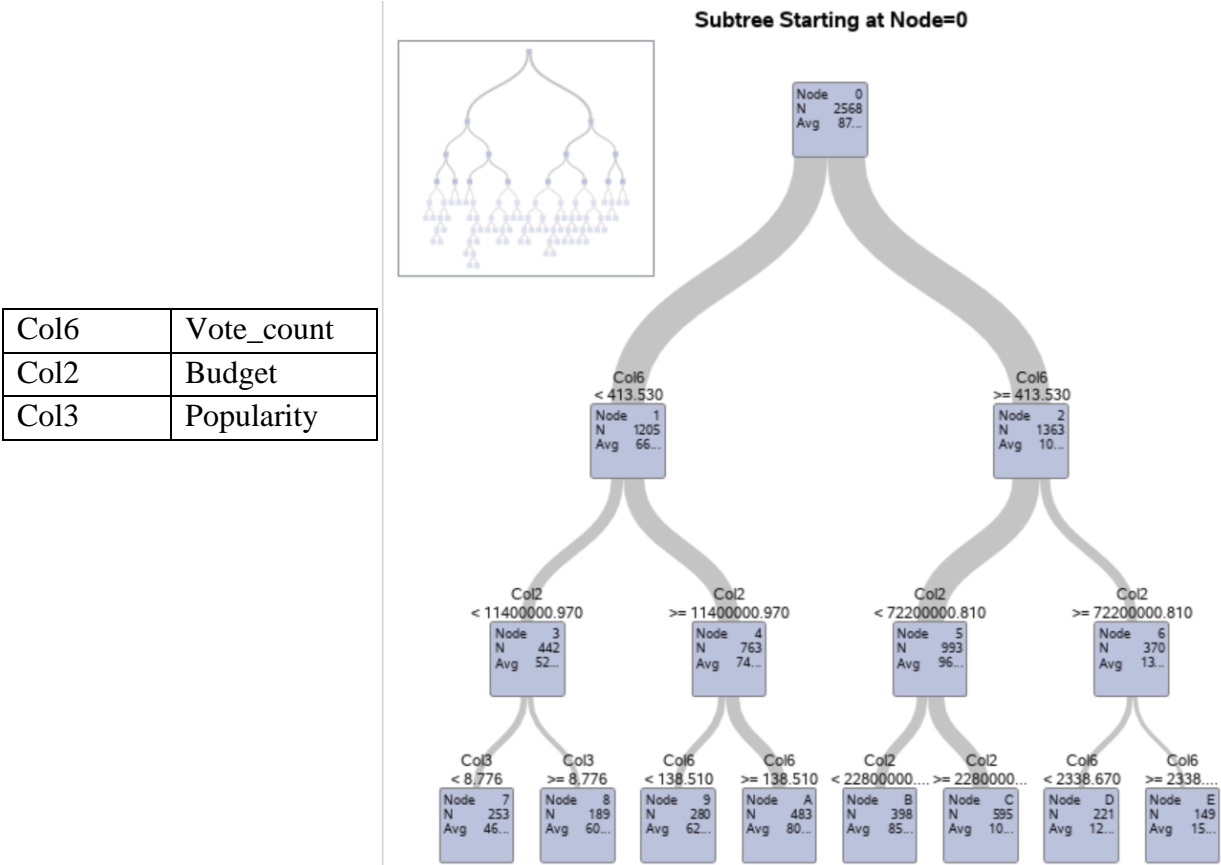


Figure 8: Descriptor of Variables (left) and Regression Subtree Starting at Node = 0 (right).

The relative importance of each rule is shown in the following Figure 9. We can interpret this to mean that vote count and budget are the most important variables in our regression tree. This means that as you decide to make a movie, you would select each branch to determine which kind of movie is most profitable, with revenue increasing as you proceed further down the tree and to the right.

Variable Importance				
Variable	Variable Label	Training		Count
		Relative	Importance	
Col6	vote_count	1.0000	1160.8	12
Col2	budget	0.7215	837.5	13
Col5	vote_average	0.1353	157.0	8
Col3	popularity	0.1281	148.8	1
Col4	runtime	0.0918	106.5	4
Col8	genre Adventure	0.0715	82.9700	1
Col27	us_production 0	0.0713	82.7363	1
Col10	genre Comedy	0.0673	78.1470	1
Col7	genre Action	0.0450	52.1827	1
Col21	genre Romance	0.0446	51.7872	1
Col11	genre Crime	0.0431	50.0126	1
Col24	genre War	0.0342	39.6791	1

Figure 9: Variable importance from Regression Tree

Regression trees are useful to look at high power interactions while linear regression models can easily be interpreted. Both these models have advantages and disadvantages. In terms of our data, we examined how the regression tree performed in comparison to our OLS model. The summary of this is found in *Table 5* below. In this table we see that our reduced model had an MSPR of 344.82 while our regression tree had a value of 422.07. This shows that our regression tree was not able to better predict on new data as our reduced model was so it is not preferable to linear regression in the case of our dataset.

Model	Reduced Model	Regression Tree
MSPR	344.82	422.07

Table 5: Comparison of Reduced Model and Regression Tree.

## Conclusion:

Predicting a movie's revenue is of important interest to producers and directors alike. In this work we examined how revenue for a movie could be predicted using a few explanatory variables with linear regression and an alternative approach of a regression tree.

Using these techniques, we found that longer movies with a higher vote count and a higher budget provide the largest revenue. Our tree model shows that vote count and budget are the most important variables, so this further confirms our OLS model.

In the future, this data could be used to find cultural trends in movie success over time, if the data was measured season by season. Also, this data could model expected revenues of new movies to help determine the best budget to maximize potential profit.

To improve the validity of our model we would want to look a more possible explanatory variables as we were limited to only examining a few. Also, although our linear regression model did outperform the regression tree, examining a more in-depth neural network would be advantageous in the case of this data. A neural network could look at high power interaction

274 terms with many different combinations and could have better predictive power than simple OLS  
275 regression. OLS regression is a powerful tool because it can easily be interpreted but it does  
276 have many limitations. Predicting the success of a movie is a complicated problem due to the  
277 influence of many different factors.

```

278                                     Appendix
279  /* This first line of code will need to be changed */
280  FILENAME REFFILE '/home/u45031672/my_courses/STAT 5100/Final
281  Project/melissa_movies_update_edited.csv';
282  PROC IMPORT DATAFILE=REFFILE replace
283          DBMS=CSV
284          OUT=WORK.melissa_movies_update_edited;
285          GETNAMES=YES;
286  RUN;
287  /*Examine Scatterplots, Boxplots and Histograms for quantitative variables.
288  proc sgplot data=melissa_movies_update_edited;
289          scatter x=budget y=revenue;
290  run;
291  proc univariate data=melissa_movies_update_edited nonprint;
292  histogram budget;
293  run;
294  proc sgplot data=melissa_movies_update_edited;
295          vbox budget;
296  run;
297  proc sgplot data=melissa_movies_update_edited;
298          scatter x=popularity y=revenue;
299  run;
300  proc univariate data=melissa_movies_update_edited nonprint;
301  histogram popularity;
302  run;
303  proc sgplot data=melissa_movies_update_edited;
304          vbox popularity;
305  run;
306  proc sgplot data=melissa_movies_update_edited;
307          scatter x=runtime y=revenue;
308  run;
309  proc univariate data=melissa_movies_update_edited nonprint;
310  histogram runtime;
311  run;
312  proc sgplot data=melissa_movies_update_edited;
313          vbox runtime;
314  run;
315  proc sgplot data=melissa_movies_update_edited;
316          scatter x=vote_average y=revenue;
317  run;

```

```

318 proc sgplot data=melissa_movies_update_edited;
319     vbox vote_average;
320 run;
321 proc univariate data=melissa_movies_update_edited nonprint;
322 histogram vote_average;
323 run;
324 proc sgplot data=melissa_movies_update_edited;
325     scatter x=vote_count y=revenue;
326 run;
327 proc univariate data=melissa_movies_update_edited nonprint;
328 histogram vote_count;
329 run;
330 proc sgplot data=melissa_movies_update_edited;
331     vbox vote_count;
332 run;
333 data melissa_movies_update_edited; set melissa_movies_update_edited;
334 if budget in (0) then delete;
335 /*Remove if vote_average or vote_count equal zero to do variable tranfromation*/
336 if vote_average in (0) then delete;
337 if vote_count in (0) then delete;
338 run;
339 proc glmmod data=melissa_movies_update_edited outdesign=GLMDesign outparm=GLMParm
340 NOPRINT;
341     class release_date genre us_production;
342     model revenue=budget popularity runtime vote_average vote_count genre english
343 us_production;
344 run;
345
346 /* Separate Into Training and Test Sets.
347 Only Fit Models to the Training Set. The variable
348 "Selected" separates training (0) from test (1) */
349 proc surveyselect data=GLMDesign seed=12345 out=movie
350     rate=0.2 outall; /* Withold 20% for validation */
351 run;
352 data train; set movie;
353 if Selected = 0;
354 run;
355 data test; set movie;
356 if Selected = 1;
357 run;

```

```

358
359 /*Crude Regression Model*/
360 proc reg data=train
361     plots =(CooksD RStudentByLeverage DFFITS DFBETAS);
362     model revenue = COL1-COL28/vif;
363     output out=out0 r=resid p=pred;
364     store regModel;
365 run;
366 %resid_num_diag(dataset=out0, datavar=resid, label ='Residual',
367 predvar=pred, predlabel = 'Predicted Value Initial Model');
368 run;
369 /*Transformation for each variable*/
370 /*COL2 lambda equals 0.35*/
371 proc transreg data=train;
372     model boxcox(COL2/lambda=-0.6 to 0.6 by 0.05)
373         =identity(revenue);
374     title1 'Box-Cox Transformation';
375 run;
376 /*COL3 lambda equals 0.3*/
377 proc transreg data=train;
378     model boxcox(COL3/lambda=-0.6 to 0.6 by 0.05)
379         =identity(revenue);
380     title1 'Box-Cox Transformation';
381 run;
382 /*COL4 lambda equals -0.65*/
383 proc transreg data=train;
384     model boxcox(COL4/lambda=-2 to 2 by 0.05)
385         =identity(revenue);
386     title1 'Box-Cox Transformation';
387 run;
388 /*COL5 lambda equals 1.85 This tranfomation is not included as it doesn't make sense*/
389 proc transreg data=train;
390     model boxcox(COL5/lambda=-2 to 2 by 0.05)
391         =identity(revenue);
392     title1 'Box-Cox Transformation';
393 run;
394 /*COL6 lambda equals 0.25*/
395 proc transreg data=train;
396     model boxcox(COL6/lambda=-0.6 to 0.6 by 0.05)
397         =identity(revenue);

```

```

398         title1 'Box-Cox Transformation';
399 run;
400
401 /*Fit data using interpretable transformations that have significant effect*/
402 proc transreg data=train;
403     model boxcox(revenue/lambda=-0.2 to 0.4 by 0.05)
404         =identity(COL1 cubedrootCOL2 cubedrootCOL3 COL4 COL5 fourthrootCOL6
405         COL7-COL28);
406     title1 'Box-Cox Transformation';
407 run;
408 data train; set train;
409     fourthrootrevenue = (revenue)**(1/4);
410     cubedrootCOL2 = (COL2)**(1/3);
411     cubedrootCOL3 = (COL3)**(1/3);
412     fourthrootCOL6 = (COL6)**(1/4);
413     cubedrootCOL2_fourthrootCOL6 =cubedrootCOL2*fourthrootCOL6;
414 run;
415 proc reg data=train plots =(CooksD RStudentByLeverage DFFITS DFBETAS);
416     model fourthrootrevenue = cubedrootCOL2 cubedrootCOL3 COL4 COL5
417     fourthrootCOL6 COL7-COL28 /vif;
418     output out=out6 r=resid p=pred;
419     title1 'Simple model for Tranfomed Data';
420 store intialmodel;
421 run;
422 %resid_num_diag(dataset=out6, datavar=resid, label ='Residual',
423 predvar=pred, predlabel = 'Predicted Value Tranformed');
424 run;
425
426 /*Variable selection*/
427 /*Stepwise Selection*/
428 proc reg data=train;
429     model fourthrootrevenue = COL1 cubedrootCOL2 cubedrootCOL3 COL4 COL5
430     fourthrootCOL6 COL7-COL28
431     /selection=stepwise slentry=.05 slstay=.05;
432     title1 'Stepwise Selection';
433 run;
434 proc reg data=train;
435     model fourthrootrevenue = COL1 cubedrootCOL2 cubedrootCOL3 COL4 COL5
436     fourthrootCOL6 COL7-COL28
437     /selection=stepwise slentry=.01 slstay=.01;

```



```

438         title1 'Stepwise Selection';
439 run;
440
441 /*Backwards Elimination*/
442 proc reg data=train;
443     model fourthrootrevenue = COL1 cubedrootCOL2 cubedrootCOL3 COL4 COL5
444         fourthrootCOL6 COL7-COL28
445         /selection=backward slstay=0.05;
446     title1 'Backward Elimination';
447 run;
448 proc glmselect data=train plots=(criterion ase);
449     model fourthrootrevenue = COL1 cubedrootCOL2 cubedrootCOL3 COL4 COL5
450         fourthrootCOL6 COL7-COL28 /
451         selection=backward slstay=0.5;
452     title1 'Backwards Variable Selection';
453 run;
454
455 /*Model with variable selection*/
456 proc reg data=train plots (label) =(CooksD RStudentByLeverage DFFITS DFBETAS);
457     model fourthrootrevenue = cubedrootCOL2 COL4 fourthrootCOL6 COL7-COL24
458         COL27 /vif;
459     output out=out6 r=resid p=pred;
460     title1 'Simple model for Reduced variables'
461 store mymodel6;
462 run;
463 %resid_num_diag(dataset=out6, datavar=resid, label ='Residual',
464 predvar=pred, predlabel = 'Predicted Value Reduced');
465 run;
466 /*****Look at interaction terms*****/
467 /**BUDGET and VOTE COUNT**/
468 /* Define higher-order predictors */
469 data train; set train;
470 cubedrootCOL2_fourthrootCOL6 =cubedrootCOL2*fourthrootCOL6;
471 sv2 = cubedrootCOL2**2;
472 fr2 = fourthrootCOL6**2;
473 run;
474 proc reg data=train;
475     model fourthrootrevenue = cubedrootCOL2 fourthrootCOL6 cubedrootCOL2_fourthrootCOL6
476     /vif;
477     title1 "Interaction model Budget and Vote Count";

```

```

478 run;
479 proc reg data=train;
480     model fourthrootrevenue = cubedrootCOL2 fourthrootCOL6
481 cubedrootCOL2_fourthrootCOL6 sv2 fr2 /vif;
482     highercheck: test cubedrootCOL2_fourthrootCOL6=sv2=fr2=0;
483     title1 'Check for higher-order predictors Budget and Vote Count';
484 run;
485 proc reg data=train;
486     model fourthrootrevenue = cubedrootCOL2 fourthrootCOL6/vif;
487     title1 'Lower-order model';
488 run;
489 /* Now look at higher-order variables with standardized data */
490 proc stdize data=train out=std_train
491     method=std mult=.0197372692;
492 run; /* Note that mult = 1/sqrt(n-1) */
493 data std_train; set std_train;
494 cubedrootCOL2_fourthrootCOL6 =cubedrootCOL2*fourthrootCOL6;
495 sv2 = cubedrootCOL2**2;
496 fr2 = fourthrootCOL6**2;
497 run;
498 proc reg data=std_train;
499     model fourthrootrevenue = cubedrootCOL2 fourthrootCOL6
500 cubedrootCOL2_fourthrootCOL6 / vif;
501     title1 'Check for interaction (standardized scale) Budget and Vote Count';
502 run;
503 proc reg data=std_train;
504 model fourthrootrevenue = cubedrootCOL2 fourthrootCOL6 cubedrootCOL2_fourthrootCOL6
505 sv2 fr2 /vif;
506 highercheck: test cubedrootCOL2_fourthrootCOL6=sv2=fr2=0;
507 title1 'Check for higher-order predictors (standardized scale) Budget and Vote Count';
508 run;
509
510 /*Model with variable selection AND interaction term*/
511 proc reg data=train plots =(CooksD RStudentByLeverage DFFITS DFBETAS);
512     model fourthrootrevenue = COL1 cubedrootCOL2 COL4 fourthrootCOL6
513 cubedrootCOL2_fourthrootCOL6 COL7-COL25 COL27 /vif;
514     output out=out13 r=resid p=pred;
515     title1 'Simple model for Reduced variables'
516 store mymodel13;
517 run;

```

```

518 %resid_num_diag(dataset=out13, datavar=resid, label='Residual',
519 predvar=pred, predlabel='Predicted Value Reduced');
520 run;
521
522 /*Add in transformations to test data to calculate MSPR*/
523 data test; set test;
524     fourthrootrevenue = (revenue)**(1/4);
525     cubedrootCOL2 = (COL2)**(1/3);
526     cubedrootCOL3 = (COL3)**(1/3);
527     fourthrootCOL6 = (COL6)**(1/4);
528     cubedrootCOL2_fourthrootCOL6 = cubedrootCOL2*fourthrootCOL6;
529 run;
530
531 /*MSPR for full model*/
532 proc plm restore=initialmodel;
533     score data=test out=newTest predicted;
534     run;
535 data newTest; set newTest;
536 MSE = (fourthrootrevenue - Predicted)**2;
537 run;
538 proc means data = newTest;
539 var MSE;
540 run;
541
542 /******MSPR for null model*****/
543 proc reg data=train
544     plots =(Cooksd RStudentByLeverage DFFITS DFBETAS);
545     model fourthrootrevenue = ;
546     output out=out2 r=resid p=pred;
547     store modelintercept;
548 run;
549
550 proc plm restore=modelintercept;
551     score data=test out=newTest10 predicted;
552     run;
553 data newTest10; set newTest10;
554 MSE = (fourthrootrevenue - Predicted)**2;
555 run;
556 proc means data = newTest10;
557 var MSE;

```

```

558 run;
559
560 /*MSPR for reduced model NO Interaction*/
561 proc plm restore=mymodel6;
562   score data=test out=newTest predicted;
563   run;
564   data newTest; set newTest;
565   MSE = (fourthrootrevenue - Predicted)**2;
566   run;
567   proc means data = newTest;
568   var MSE;
569   run;
570
571 /*MSPR for reduced model Yes Interaction*/
572 proc plm restore=mymodel13;
573   score data=test out=newTest predicted;
574   run;
575   data newTest; set newTest;
576   MSE = (fourthrootrevenue - Predicted)**2;
577   run;
578   proc means data = newTest;
579   var MSE;
580   run;
581
582 /*Regression Tree*/
583 proc hpsplit data=train seed=123 maxdepth=10 maxbranch=2;
584   model fourthrootrevenue=COL1-COL28;
585   output out=out20;
586   code file='/home/u45031672/my_courses/STAT 5100/Final Project/tree2.sas';
587   /* This saves the tree to a file (need to change the path) */
588   run;
589   /**Call the test data and include the tree, this will make predictions on the tree */
590   data scored;
591   set test;
592   %include '/home/u45031672/my_courses/STAT 5100/Final Project/tree2.sas';
593   run;
594   /* Now calculate the MSPR as we did in OLS */
595   data testTree;
596   set scored;
597   ASE = (fourthrootrevenue - P_fourthrootrevenue)**2;

```

```
598  run;
599  proc means data = testTree;
600  var ASE;
601  run;
```