## 1) Data Exploration

The excel file was read. The quantitative response variable ('Achieve') and the predictor variables ('Family', 'Peer' and 'School') were being explored. In particular, the below correlation table and heatmap were generated by executing the python scripts.
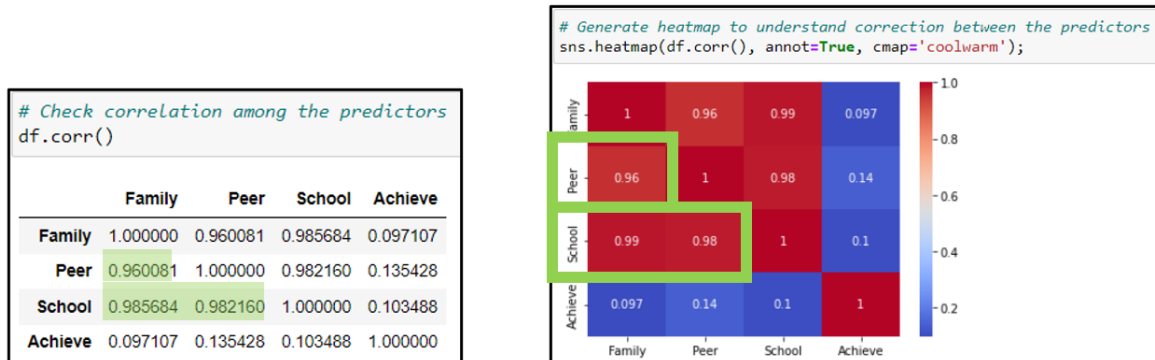


*Diagram 1: Correlation Table and Heatmap generated on the variables to be studied*

** All python scrips used are presented on top of the results

It was observed from both the correlation table (values shaded in green) and heatmap (values boxed in green) that the 'Family', 'Peer' and 'School' predictor variables are strongly correlated.

## 2) Fitting of 3 models

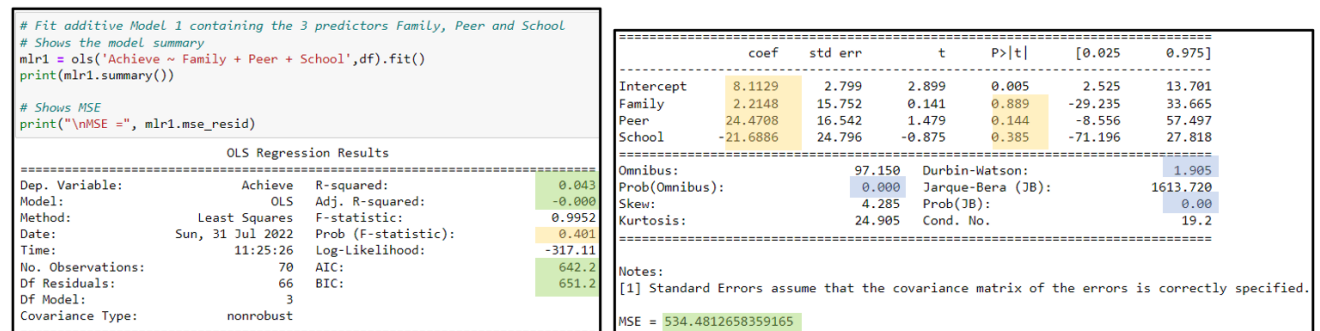## (i) Model 1: Additive MLR model

The below additive MLR model was fitted.



*Diagram 2: Summary and MSE details for MLR Model 1 fitted*

** Values used for goodness of fit and accuracy are shaded in green, values on test statistics and coefficient are shaded in yellow, and values from tests conducted which used in model assumptions are shaded in blue. These shaded values were summarised in Table 1 comparison of models, as well as referenced in this document.

The MLR Model1 equation can be expressed as:

$$\widehat{\text{Achieve}} = 8.1129 + 2.2148 * \text{Family} + 24.4708 * \text{Peer} - 21.6886 * \text{School}$$

After fitted the Model 1, it's residuals were studied to verify if the assumptions needed for the linear regression model were met. It was observed that the following 2 assumptions are NOT complied with:

- Residuals should follow a normal distribution
- There should be little or no multicollinearity among the predictors

## (ii) Model 2: Investigation of log-transformation for response variable 'Achieve'

On assumption where the residuals of model 1 not following a normal distribution, the following verification were carried out:
- Plotted the histogram on residuals, where it should show a bell-shape distribution. The generated distribution is skewed to the right as shown in Diagram 3
- Plotted the normal probability plot (or QQ plot), where it's points should fall roughly on a straight line. The points in the generated plot does not fall on the straight line, as shown in Diagram 3
- Conducted Omnibus and Jarque-Bera normality tests, indicating that sufficient evidence at 5% of significance to conclude that data are not normally distributed.

```
# Generate Histogram of the residuals to check normality visually
sns.distplot(mlr1.resid);
plt.title('Model 1 - Histogram of Residuals')
plt.ylabel('Density')
plt.xlabel('Residuals');
```

```
# Generate qqplot on Residuals
sm.qqplot(mlr1.resid, line='s');
```
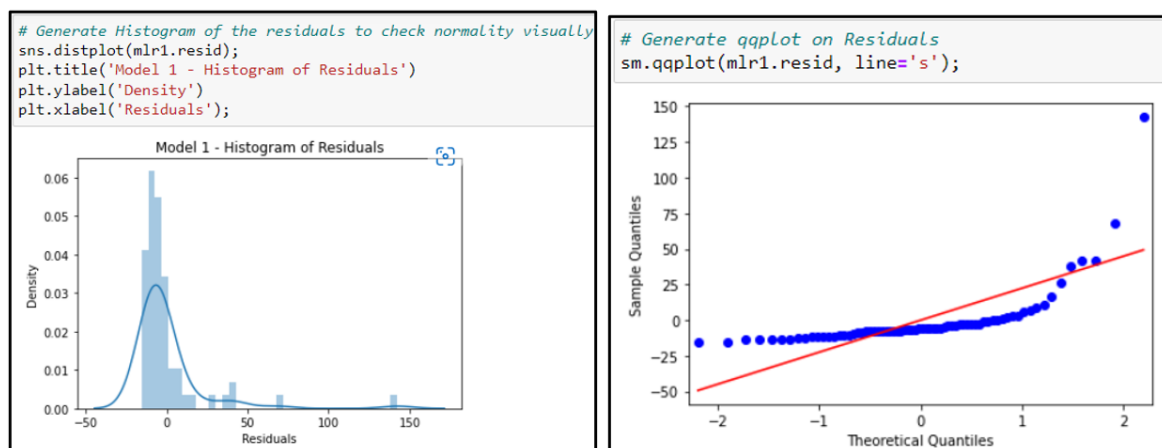


*Diagram 3: Histogram on residuals showing right skewed distribution (left), and QQ plot showing data points not falling on the red line (right), for model 1*

H0: Residuals are normally distributed
H1: Residuals are NOT normally distributed

Level of significance: α = 0.05

Omnibus test's p-value (refer to Diagram 2, Model 1 summary) = 0.000

**Decision:** p-value (which is 0.000) is < α (which is 0.05). Thus, H0 is rejected.

**Conclusion**: There is sufficient evidence at 5% of significance to conclude that data are not normally distributed.

H0: Residuals are normally distributed
H1: Residuals are NOT normally distributed

Level of significance: α = 0.05

Jarque-Bera test's p-value (refer to Diagram 2, Model 1 summary) = 0.00

**Decision:** p-value (which is 0.00) is < α (which is 0.05). Thus, H0 is rejected.

**Conclusion**: There is sufficient evidence at 5% of significance to conclude that data are not normally distributed.

*Diagram 4: Omnibus Normality Test (left) and Jarque-Bera Normality tests (right) for model 1*

To improve the normality distribution, usually the response variable is to be transformed.

The following python script to compute the correlation of response variable vs different transformed values is executed. The results showed that the log (lnx) transformation generates the closest score to 1.

```
# Compute correlations of Y vs. different transformed X'
from scipy.stats import pearsonr
transformations = {'x':lambda x: x, 'lnx':lambda x: np.log(x), '1/x':lambda x: 1/x,
                   'sqrt(x)':lambda x: np.sqrt(x), 'exp(-x/1000)':lambda x: np.exp(-x/1000)}
for tf in transformations:
    r_coef, r_pval = pearsonr(transformations[tf](df.Achieve), df.Family+df.Peer+df.School)
    print('Transformation: %s   Pearsons r: %0.3f   r_p_value: %0.3f' %(tf, r_coef, r_pval))

Transformation: x    Pearsons r: 0.112    r_p_value: 0.357
Transformation: lnx   Pearsons r: 0.429    r_p_value: 0.000
Transformation: 1/x   Pearsons r: -0.233   r_p_value: 0.052
Transformation: sqrt(x)   Pearsons r: 0.234   r_p_value: 0.051
Transformation: exp(-x/1000)   Pearsons r: -0.117   r_p_value: 0.335
```

*Diagram 5: Compute different Pearson r scores for correlation between the response and the transformed variables.*

Thus, this log transformation on the response variable was chosen and model 2 was fitted.



*Diagram 6: Summary and MSE details for MLR Model 2 using log-transformation of response variable is fitted*

The MLR Model 2 equation can be expressed as:

$$log(\widehat{Achieve}) = -0.07 + 1.1013 * Family + 2.3221 * Peer - 2.2810 * School$$

Next, the histogram and QQ plot were generated for Model 2 to verify resolution of the residual normality. In the below diagram, both plots demonstrated normal distribution after the log linear transformation.

Based on the model summary, both p-values for Omnibus and Jarque-Bera normality tests (=0.756 and 0.749 respectively) are > α (= 0.05). Thus, there is insufficient evidence at 5% of significance to conclude that the data are not normally distributed.
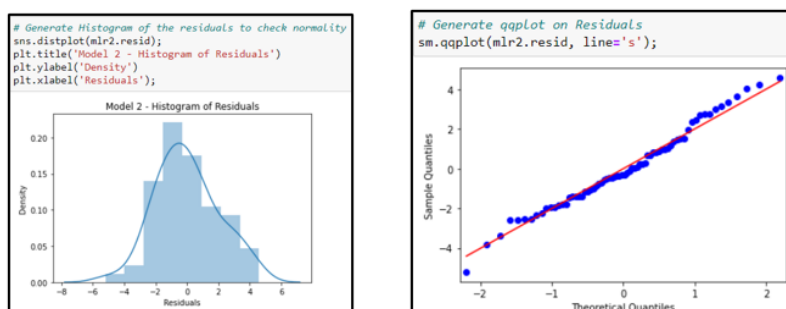


*Diagram 7: Histogram on residuals showing normal bell shape (left), and QQ plot showing data points falling on the red line (right), for model 2*

## (iii) Model 3: Investigate multicollinearity and perform Principal Component Regression

The multicollinearity assumption for linear regression model is also NOT complied with.

The followings indicate the presence of data-based multicollinearity:
- High (strong) correlations between the predictors 'Family', 'Peer' and 'School', as indicated in the correlation matrix and heatmap in the Diagram 1.
- Predictors 'Family', 'Peer' and 'School' have high Variance Inflation Factor (VIF) of 37, 30 and 83 respectively, which are much higher than 5 (=highly correlated).

```
# Generate VIF for each predictor and term
print ("VIF for Model 1 :")
for i in range(3):
    predictor = mlr1.model.exog_names[i+1]
    VIF = vif(mlr1.model.exog, i+1)
    print ("  VIF for predictor",predictor,'is',VIF)

VIF for Model 1 :
    VIF for predictor Family is 37.580639932204285
    VIF for predictor Peer is 30.21165657710285
    VIF for predictor School is 83.15543705009657
```

*Diagram 8: VIF for Model 1*

To resolve the multicollinearity problem, Principal Component Analysis (PCA) can be used to find the "best" linear combinations for the 3 highly corrected predictors, followed by adapting the desired PC (replacing the highly correlated predictors) in the model.

The following python script is executed to find the PCA, compute the PC values and then added the PC values to the data frame.

```
# Create the PCA instance with 3 components, as p=3
pca3 = PCA(n_components=3)

# Then fit highly corrected predictor Family, Peer, School data to PCA,
# i.e., to determine linear coefficients.
# Variables will be automatically centred by default.
pca3.fit(df.iloc[:,0:3])

print('Proportion of total variance:', pca3.explained_variance_ratio_)

print('\nPrincipal components are:\n', pca3.components_)

Proportion of total variance: [0.98445476 0.01286135 0.00268389]

Principal components are:
[[ 0.61732371  0.5241853   0.5866355 ]
 [-0.67025631  0.74086406  0.04332342]
 [ 0.41190766  0.41994072 -0.80869145]]
```

```
# Compute and transform the original predictors to PC1,PC2 & PC3 scores.
newX = pca3.transform(df.iloc[:,0:3])

# Append PC1, PC2 & PC3 scores to (original) data frame
df['pc1'] = newX[:,0]
df['pc2'] = newX[:,1]
df['pc3'] = newX[:,2]
df.tail()
```

| | Family | Peer | School | Achieve | pc1 | pc2 | pc3 |
|---|---|---|---|---|---|---|---|
| 65 | 0.53940 | 0.16182 | 0.33477 | 2.638842 | 0.540716 | -0.229742 | 0.005424 |
| 66 | 0.22491 | 0.74800 | 0.66182 | 23.592527 | 0.845700 | 0.429496 | -0.142438 |
| 67 | 1.48244 | 1.47079 | 1.54283 | 0.148375 | 2.517711 | 0.160286 | -0.033389 |
| 68 | 2.05425 | 1.80369 | 1.90066 | 1.907856 | 3.255120 | 0.039163 | 0.052569 |
| 69 | 1.24058 | 0.64484 | 0.87372 | 0.172191 | 1.542930 | -0.318511 | 0.061241 |

*Diagram 9: Performed PCA, transformed the linear combinations into PC and added them to the data frame.*

PC1 accounts for 98.4% of the total variability, and thus Model 3 was fitted using just the PC1 as follows:

```
# Fit SLR model using PC1 score as the predictor.
slr3 = ols('Achieve ~ pc1', df).fit()
print(slr3.summary())
print("\nMSE =", slr3.mse_resid)
```

```
                    OLS Regression Results
==============================================================================
Dep. Variable:          Achieve    R-squared:                       0.012
Model:                      OLS    Adj. R-squared:                 -0.002
Method:           Least Squares    F-statistic:                    0.8444
Date:          Sun, 31 Jul 2022    Prob (F-statistic):              0.361
Time:                  14:52:31    Log-Likelihood:                -318.23
No. Observations:            70    AIC:                             640.5
Df Residuals:                68    BIC:                             645.0
Df Model:                     1
Covariance Type:      nonrobust
==============================================================================
```

```
==============================================================================
                 coef    std err          t      P>|t|      [0.025      0.975]
------------------------------------------------------------------------------
Intercept      8.6636      2.766      3.132      0.003       3.144      14.183
pc1            1.4711      1.601      0.919      0.361      -1.724       4.666
==============================================================================
Omnibus:                    104.048    Durbin-Watson:                   1.904
Prob(Omnibus):                0.000    Jarque-Bera (JB):             2126.795
Skew:                         4.670    Prob(JB):                         0.00
Kurtosis:                    28.337    Cond. No.                         1.73
==============================================================================

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

MSE = 535.5771011689628
```

*Diagram 10: Summary and MSE details for SLR Model 3 using PC1*

The MLR Model 3 equation can be expressed as:

$$\widehat{\text{Achieve}} = 8.6636 + 1.4711 * PC1, \text{ where}$$

PC1 = 0.61732371*($X_{\text{Family}}$  $X_{\text{Family Center}}$) + 0.5241853*($X_{\text{Peer}}$ - $X_{\text{Peercenter}}$) + 0.5866355*($X_{\text{School}}$ - $X_{\text{School Center}}$)

PC1 = 0.61732371*($X_{\text{Family}}$ - 0.049383) + 0.5241853*($X_{\text{Peer}}$ - 0.046314 ) + 0.5866355*($X_{\text{School}}$ − 0.031906)

The VIF for Model 3 is now 1, verifying that that the multicollinearity is resolved.

```
# Generate VIF for PC1
print ("VIF for Model 3 :")
print(slr3.model.exog_names[1], vif(slr3.model.exog,1))

VIF for Model 3 :
pc1 1.0
```

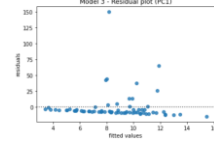*Diagram 11: VIF for Model 3*

## 3) Comparison of the 3 Models

Refer to below table 1 for details of comparison.

In summary:
- All the 3 models do not provide good model fit as the coefficient of determination is low, and the test statistics indicate that the predictors are statistically insignificant

- Of the 3 models, Model 3 might be the better choice due to its more favourable AIC and BIC values, and with absence of multicollinearity

| | Metrics | Model 1 | Model 2 | Model 3 | Explanation |
|---|---|---|---|---|---|
| (a) Goodness of fit | (i) Coefficient of determination, $R^2$ or adjusted $R^2$ | $R^2 = 0.043$<br><br>Adjusted $R^2 = -0.000$ | $R^2 = 0.206$<br><br>Adjusted $R^2 = 0.170$ | $R^2 = 0.012$ | • For SLR, $R^2$ is more useful. For MLR, adjusted $R^2$ is more useful, as a penalty factor is applied to insignificant predictors.<br>• All the 3 models have low values, indicating that the predictors do not account much proportion of the total variation in the response variable. |
| | (ii) AIC, BIC | • AIC = 642.2<br>• BIC = 651.2 | • AIC = 304.4<br>• BIC = 313.4 | • AIC = 640.5<br>• BIC = 645.0 | • The lower AIC and BIC, the better is the model fit.<br>• Model 3 is a better model fit than Model 1.<br>• It is not meaningful to compare Model 2 against Model 1 and 3, as it is of different scaling. |
| (b) Accuracy for prediction | (i) MSE | 534.481 | 4.286 | 535.577 | • Mean Square Error (MSE) measures the prediction accuracy of model. The lower the MSE, the higher is the prediction accuracy.<br>• Thus, model 1 is slightly more accurate for prediction than model 3 |
| (c) Statistical Tests for prediction | (i) F-test on overall model | 0.401 | 0.002 | 0.361 | As the p-value for Model 1 and 3 (=0.401 and 0.361 respectively) is > α (=0.05), do not reject H0.<br>There is insufficient evidence at 5% of significance to conclude that model contains at least one significant predictor. |
| | (ii) t-tests on coefficients | p-values are 0.889 (Family), 0.144 (Peer), 0.385 (School) | p-values are 0.438 (Family), 0.122 (Peer), 0.308 (School) | p-value for PC1 is 0.361 | In all 3 models, each predictors' p-values are > α (=0.05), thereby H0 is not rejected.<br>Thus, for all the predictors in the 3 models, there is insufficient evidence at 5% of significance to conclude that each predictor is statistically significant. |

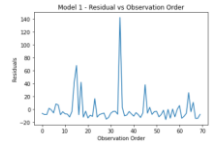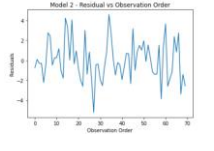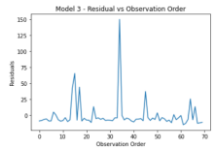| | Metrics | Model 1 | Model 2 | Model 3 | Explanation |
|---|---|---|---|---|---|
| (d) Meeting assumption for linear regression model | (i) Residuals follows a normal distribution<br>- Via plotting histogram on residuals<br>- Via QQ plot<br>- Conduct Omnibus test<br>- Conduct Jarque-Bera (JB) test | Histogram:<br><br><br>QQ Plot:<br><br><br>Prob(Omnibus) = 0.000<br><br>Prob(JB)= 0.00 | Histogram:<br><br><br>QQ Plot:<br><br><br>Prob(Omnibus) = 0.756<br><br>Prob(JB)= 0.749 | Histogram:<br><br><br>QQ Plot:<br><br><br>Prob(Omnibus) = 0.000<br><br>Prob(JB)= 0.00 | For Model 1 and 3:<br>Both the residual histogram and QQplot shows it is skewed to the right.<br><br>The omnibus and Jarque-Bera tests both have p-values < α (=0.05). Thus, H0 is rejected. There is sufficient evidence at 5% of significance to conclude that data are not normally distributed.<br><br>For Model 2:<br>The residuals follow normal distribution. |
| | (ii) Residuals have constant variance (homoscedastic)<br>- Via plotting residuals vs fitted value<br>- Conduct Breusch Pagan test | Plot:<br><br><br>P-value for Breusch Pagan 0.302 | Plot:<br><br><br>P-value for Breusch Pagan 0.944 | Plot:<br><br><br>P-value for Breusch Pagan 0.345 | Model 1, 2 and 3 have Breusch Pagan p-values (=0.302, 0.944 and 0.345 respectively) which are > α (=0.05).<br>Thus, do not reject H0.<br><br>There is insufficient evidence at 5% of significance to conclude that the residuals are heteroscedastic. |

| | Metrics | Model 1 | Model 2 | Model 3 | Explanation |
|---|---|---|---|---|---|
| | (iii) Residuals are independent<br>- Plot residuals vs observed order<br>- Conduct Durbin-Watson Test | Plot:<br><br>`plt.plot(mlr1.resid)`<br>`plt.xlabel('Observation Order')`<br>`plt.ylabel('Residuals')`<br>`plt.title ('Model 1 - Residual vs Observation Order')`<br><br>_Model 1 - Residual vs Observation Order plot_<br><br>Durbin-Watson 1.905 | Plot:<br><br>`plt.plot(slr2.resid)`<br>`plt.xlabel('Observation Order')`<br>`plt.ylabel('Residuals')`<br>`plt.title ('Model 2 - Residual vs Observation Order');`<br><br>_Model 2 - Residual vs Observation Order plot_<br><br>Durbin-Watson 1.791 | Plot:<br><br>`plt.plot(slr3.resid)`<br>`plt.xlabel('Observation Order')`<br>`plt.ylabel('Residuals')`<br>`plt.title ('Model 3 - Residual vs Observation Order')`<br><br>_Model 3 - Residual vs Observation Order plot_<br><br>Durbin-Watson 1.904 | With the plot showing random pattern, and the Durbin-Watson value close to 2, there is no autocorrelation for all 3 models. |
| | (iv) Little or no multicollinearity among the predictors<br>- Observe Correlation matrix<br>- Derive High Variance Inflation Factor (VIF) | VIF for predictors:<br><br>• Family is 37.581<br>• Peer is 30.212<br>• School is 83.155 | VIF for predictors:<br><br>• Family is 37.581<br>• Peer is 30.211<br>• School is 83.155 | VIF for PC1 is 1 | The predictors in Model 1 and Model 2 have high VIF (VIF is high if > 5). Thus, Model 1 and Model 2 have multicollinearity problem.<br><br>There is no multicollinearity problem with Model 3. |

_Table 1: Comparison of the 3 models fitted_
** All figures in above table are rounded to 3 decimal places.

**<< End of Report >>**