

Defining Problem Statement

- Jamboree is an education provider that helps students crack entrance exams such as GMAT, GRE or SAT
- Jamboree has launched a feature where students/learners can come to their website and check their probability of getting into the IIVV league college. This feature estimates the chances of graduate admission from an Indian perspective.
- Help Jamboree understanding what factors are important in graduate admissions and how these factors are interrelated among themselves. It will also help predict one's chances of admission given the rest of the variables.

Dataset Definition:

- GRE Scores (out of 340)
- TOEFL Scores (out of 120)
- University Rating (out of 5)
- Statement of Purpose and Letter of Recommendation Strength (out of 5)
- Undergraduate GPA (out of 10)
- Research Experience (either 0 or 1)
- Chance of Admit (ranging from 0 to 1) -> Predictor variable

Importing Libraries

```
In [208]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import matplotlib.colors as mcolors
import itertools
import warnings
warnings.filterwarnings('ignore')

In [209]: data = pd.read_csv('Jamboree_Admission.csv')
data

Out[209]:
```

	Serial No.	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
0	1	337	118	4	4.5	4.5	9.65	1	0.92
1	2	324	107	4	4.0	4.5	8.87	1	0.76
2	3	316	104	4	3.0	3.5	8.00	1	0.72
3	4	322	110	3	3.5	2.5	8.67	1	0.80
4	5	314	103	2	2.0	2.0	8.21	0	0.65
...
495	496	332	108	5	4.5	4.0	9.02	1	0.87
496	497	337	117	5	5.0	5.0	9.87	1	0.96
497	498	330	120	5	4.5	5.0	9.56	1	0.93
498	499	312	103	4	4.0	5.0	8.43	0	0.73
499	500	327	113	4	4.5	4.5	9.04	0	0.84

500 rows x 9 columns

```
In [210]: #dropping columns useless for predictions/insights
data.drop(['Serial No.'], axis = 1, inplace = True)

# Some columns names have extra spaces
data.columns = ["GRE", "TOEFL", "University Rating", "SOP", "LOR", "CGPA", "Research", "Admit Chance"]
data.head()
```

```
Out[210]:
```

	GRE	TOEFL	University Rating	SOP	LOR	CGPA	Research	Admit Chance
0	337	118	4	4.5	4.5	9.65	1	0.92
1	324	107	4	4.0	4.5	8.87	1	0.76
2	316	104	4	3.0	3.5	8.00	1	0.72
3	322	110	3	3.5	2.5	8.67	1	0.80
4	314	103	2	2.0	2.0	8.21	0	0.65

```
In [211]: data.shape
(500, 8)
```

```
In [212]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 500 entries, 0 to 499
Data columns (total 8 columns):
#   Column              Non-Null Count  Dtype
---  ---
0   GRE                  500 non-null    int64
1   TOEFL                500 non-null    int64
2   University Rating    500 non-null    int64
3   SOP                  500 non-null    float64
4   LOR                  500 non-null    float64
5   CGPA                 500 non-null    float64
6   Research              500 non-null    int64
7   Admit Chance         500 non-null    float64
dtypes: float64(4), int64(4)
memory usage: 31.4 KB
```

Since we want to implement a linear regression model on our data, which would also require checking of correlation and regression factors, I will not convert features like University rating, SOP and LOR to categorical, even though they can be represented as 0

```
In [213]: data.describe()
```

```
Out[213]:
```

	GRE	TOEFL	University Rating	SOP	LOR	CGPA	Research	Admit Chance
count	500.000000	500.000000	500.000000	500.000000	500.000000	500.000000	500.000000	500.000000
mean	316.472000	107.192000	3.114000	3.374000	3.484000	8.576440	0.560000	0.72174
std	11.295148	6.081868	1.143512	0.991004	0.92545	0.604813	0.496884	0.14114
min	290.000000	92.000000	1.000000	1.000000	1.000000	6.870000	0.000000	0.34000
25%	308.000000	103.000000	2.000000	2.500000	3.000000	8.127500	0.000000	0.63000
50%	317.000000	107.000000	3.000000	3.500000	3.500000	8.560000	1.000000	0.72000
75%	325.000000	112.000000	4.000000	4.000000	4.000000	9.040000	1.000000	0.82000
max	340.000000	120.000000	5.000000	5.000000	5.000000	9.920000	1.000000	0.97000

```
In [214]: for column in ["University Rating", "SOP", "LOR", "Research"]:
```

```
print(f"column : \nUnique Values: {data[column].unique()}, \nUnique Value Counts: {data[column].nunique()}, \nprint(f"-----XX-----")
```

University Rating :
Unique Values: [4 3 2 5 1],
Unique Value Counts: 5,
% of total: 3 32.4 2 25.2 5 17.8 1 2.0
State: University Rating, dtype: float64
-----XX-----
SOP
Unique Values: [4.5 4. 3. 3.5 2. 5. 1.5 1. 2.5],
Unique Value Counts: 9,
% of total: 4.0 18.8 3.5 12.6 5.0 10.0 2.5 10.0 1.0 2.0
State: SOP, dtype: float64
-----XX-----
LOR
Unique Values: [4.5 3.5 2.5 3. 4. 1.5 2. 5. 1. 1.],
Unique Value Counts: 9,
% of total: 4.0 18.8 3.5 17.2 5.0 12.6 2.5 10.0 1.0 2.0
State: LOR, dtype: float64
-----XX-----
Research :
Unique Values: [1 0],
Unique Value Counts: 2,
% of total: 1 51.0 0 49.0
Name: Research, dtype: float64
-----XX-----

Null Values and Duplicates

```
In [215]: data.isna().sum()
```

```
GRE      0
TOEFL    0
University Rating  0
SOP       0
LOR       0
CGPA      0
Research  0
Admit Chance  0
dtype: int64
```

```
In [216]: data.duplicated().sum()
```

```
GRE TOEFL University Rating SOP LOR CGPA Research Admit Chance
0
```

Outlier detection

```
In [217]: def printoutlier(col):
q1 = np.quantile(col, 0.25)
q2 = np.quantile(col, 0.50)
q3 = np.quantile(col, 0.75)
min_outlier = q1 - 1.5*(IQR)
max_outlier = q3 + 1.5*(IQR)
return col[(col >= max_outlier) | (col <= min_outlier)].count()

print("Outlier Counts")
print("-----XX-----")
for col in data.columns:
sns.boxplot(x=printoutlier(data[col]))
print(col, "n", count_outliers)
```

Outlier Counts
-----XX-----
GRE : 0
count: 0
University Rating : 0
SOP : 0
count: 12
LOR : 0
CGPA : 0
Research : 0
Admit Chance : 2

```
In [218]: sns.countplot(x = data["LOR"])
```

```
Out[218]:
```

- Outliers should be removed if they are truly anomalous values. The outliers in this case seem to represent genuine extreme observations in the data. Therefore, since the outliers are valid and representative of the underlying phenomenon being studied, it is inappropriate to remove them.
- We must also consider the fact that linear regression assumes that residuals are normally distributed and homoskedasticity. Outliers can violate these assumptions. Hence if they impact the model too much, I will remove them at a later stage.

EDA - Univariate and Bi-variate

```
In [219]: continuous_cols = ["GRE", "TOEFL", "CGPA", "Admit Chance"]
discrete_cols = data.columns[~data.columns.isin(continuous_cols)]
```

```
In [220]: num_features = len(data.columns)
fig, axs = plt.subplots(nrows = num_features, ncols = 2, figsize=(6, 3 * num_features))
colors = sns.color_palette("Set1", num_features)
mpltlib_named_colors = list(mcolors.CSS4_COLORS.keys())
```

```
for i, (feature, color) in enumerate(discrete_cols, colors):
row = i
sns.boxplot(data=data, y = feature, ax = axs[row][0], color = color).set(ylabel = "")
sns.histplot(data=data, x = feature, ax = axs[row][1], color = color, edgecolor = "black").set(xlabel = "")
axs[row][0].set_title(f'{feature} Univariate analysis', x = 1.5, y = 1.1, fontsize = 12)
```

```
new_colors = itertools.islice(colors, 1, None)
for (feature, color) in zip(discrete_cols, new_colors):
row = i // 2
sns.boxplot(data=data, y = feature, ax = axs[row][0], color = color).set(ylabel = "")
sns.histplot(data=data, x = feature, ax = axs[row][1], color = color, edgecolor = "black").set(xlabel = "")
axs[row][0].set_title(f'{feature} Univariate analysis', x = 1.5, y = 1.1, fontsize = 12)
```

```
i = i + 1
plt.tight_layout()
plt.subplots_adjust(right = 2)
plt.show()
```

GRE Univariate analysis

TOEFL Univariate analysis

CGPA Univariate analysis

Admit Chance Univariate analysis

University Rating Univariate analysis

SOP Univariate analysis

LOR Univariate analysis

Research Univariate analysis

Univariate analysis of Admit chance w.r.t discrete variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables

Univariate analysis of Admit chance w.r.t continuous variables


```
explained_variance: 0.8222
mean_squared_log_error: 0.0014
r2: 0.8222
adjusted_r2: 0.819
MAE: 0.0426
MSE: 0.0036
RMSE: 0.06
```

Out[230]:

Co-efficients	
const	0.719850
GRE	0.023553
TOEFL	0.018789
University Rating	0.012394
SOP	-0.002560
LOR	0.016078
CGPA	0.067937
Research	0.010428

L1 Regularization (Lasso Regression) - Feature Selection

In [231]:

```
from sklearn.linear_model import Lasso

model = Lasso(alpha = 0.05)
model.fit(xTrain_Scaled, yTrain)
y_preds = model.predict(xTrain_Scaled)
regression_results(yTrain, y_preds, xTrain_Scaled.shape[1])
l1_coefs = pd.DataFrame(data = np.insert(model.coef_, 0, model.intercept_),
                        index = ["Constant"] + xTrain_Scaled_df.columns.tolist(), #alternatively I can do this too
                        columns = ["Co-efficients"])

l1_coefs

explained_variance: 0.6648
mean_squared_log_error: 0.0025
r2: 0.6648
adjusted_r2: 0.6588
MAE: 0.0645
MSE: 0.0068
RMSE: 0.0824
```

Out[231]:

Co-efficients	
Constant	0.719850
GRE	0.011599
TOEFL	0.001595
University Rating	0.000000
SOP	0.000000
LOR	0.000000
CGPA	0.064377
Research	0.000000

L2 Regularization (Ridge Regression)

In [232]:

```
from sklearn.linear_model import Ridge

model = Ridge(alpha = 1)
model.fit(xTrain_Scaled, yTrain)
y_preds = model.predict(xTrain_Scaled)
regression_results(yTrain, y_preds, xTrain_Scaled.shape[1])
pd.DataFrame(data = np.insert(model.coef_, 0, model.intercept_),
             index = ["Constant"] + xTrain_Scaled_df.columns.tolist(),
             columns = ["Co-efficients"])

explained_variance: 0.8222
mean_squared_log_error: 0.0014
r2: 0.8222
adjusted_r2: 0.819
MAE: 0.0426
MSE: 0.0036
RMSE: 0.06
```

Out[232]:

Co-efficients	
Constant	0.719850
GRE	0.023720
TOEFL	0.018922
University Rating	0.012404
SOP	-0.002322
LOR	0.016121
CGPA	0.067303
Research	0.010448

Testing Assumptions of Linear Regression

Sr.No	Assumption	Test
1	Predictors (x) are independent (no multicollinearity) and observed with negligible error	Multicollinearity check by VIF score (variables are dropped one-by-one till none has VIF>5)
2	Linearity of variables - There is a linear relationship between the predictors (x) and the outcome (y)	Residual vs Actual Plot should not have a pattern
3	The mean of residuals is nearly zero	Calculate Mean of Residuals
4	Residual Errors have constant variance - No Heteroscedasticity	Residual vs Actual Plot should not have a pattern
5	Combining pt4 and pt5 -> Residuals must be normally distributed (bell-curve)	Almost bell-shaped curve in residuals distribution, points in QQ plot are almost all on the line

Notes:

- In the following article: [V.I.F statsmodel add_constant](#), statsmodel implementation of VIF expects the presence of a constant in the matrix of explanatory variables. Therefore, we need to use add_constant from statsmodels to add the required constant to the dataframe before passing its values to the function.
- Features having high VIF (>5) means that they can be **predicted by other independent variables in the dataset**.
- Removing features after VIF check will **never lead to an increase in R2 Score**. This is because we are decreasing model complexity when we remove features and, the model's learning will always **remain same or less than what was before feature removal**.
- The square root of the variance inflation factor indicates how much larger the standard error increases compared to if that variable had 0 correlation to other predictor variables in the model.
- Example:** If the variance inflation factor of a predictor variable were 5.27 ($\sqrt{5.27} = 2.3$), this means that the standard error for the coefficient of that predictor variable is 2.3 times larger than if that predictor variable had 0 correlation with the other predictor variables.
- There are some **hypothesis tests for checking Heteroskedasticity** - [Breusch-Pagan](#) and [White tests](#), [Full implementation](#)

V.I.F Check

In [233]:

```
from statsmodels.stats.outliers_influence import variance_inflation_factor
from statsmodels.tools.tools import add_constant
```

In [240]:

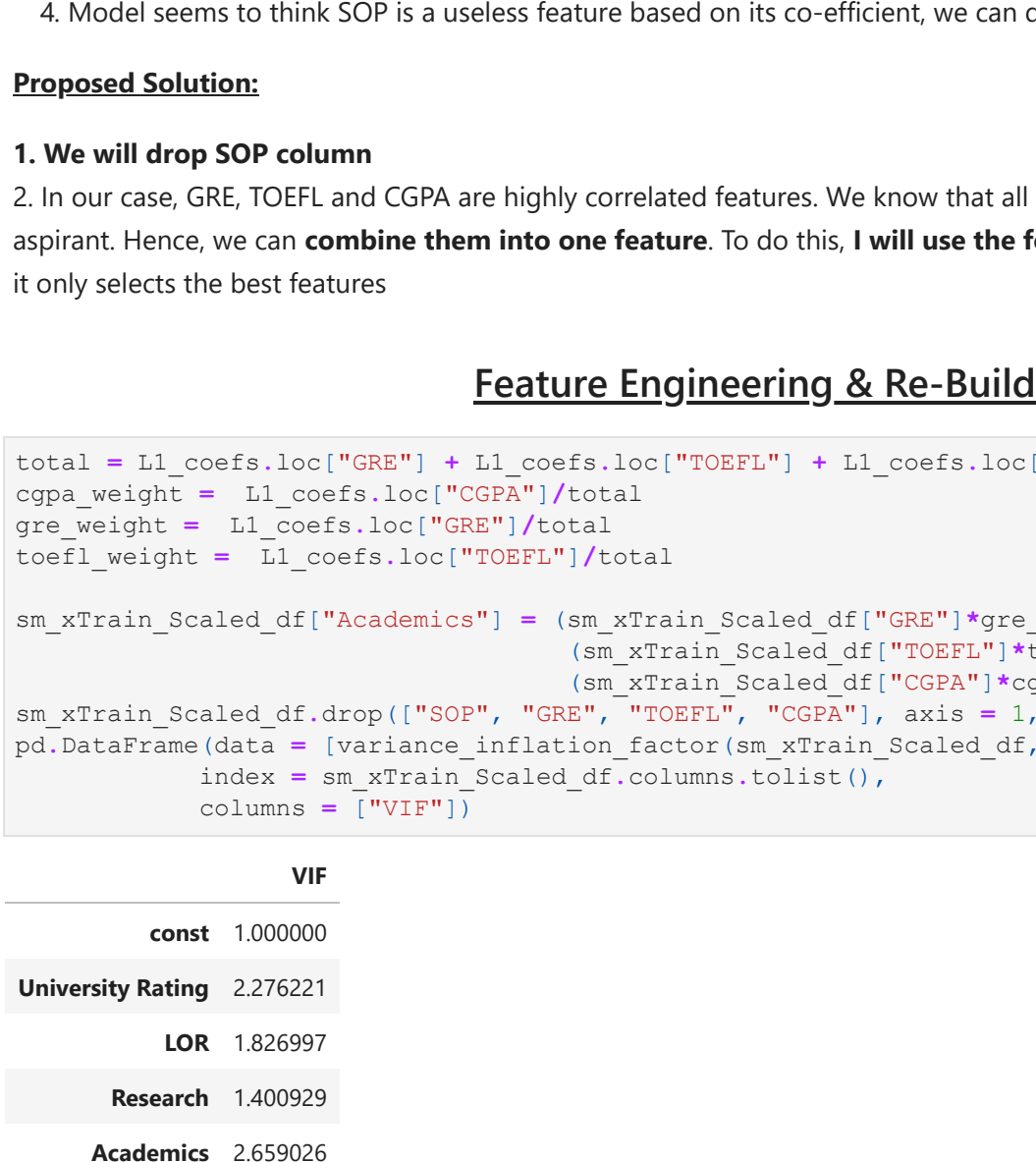
```
# calculating VIF for each feature
# variance_inflation_factor takes input as [all values in df, i'th column for which we are calc VIF]
pd.DataFrame(data = [variance_inflation_factor(sm_xTrain_Scaled_df, i) for i in range(0, sm_xTrain_Scaled_df.shape[1])],
             index = sm_xTrain_Scaled_df.columns.tolist(),
             columns = ["VIF"])
```

Out[240]:

VIF	
const	1.000000
GRE	4.662889
TOEFL	3.920767
University Rating	2.803090
SOP	3.082703
LOR	2.006346
CGPA	5.017933
Research	1.485910

In [251]:

```
sns.pairplot(data = sm_xTrain_Scaled_df, kind = "reg", vars = ["GRE", "TOEFL", "CGPA"], corner=True,
              diag_kind="hist", diag_kind="color": "forestgreen", plot_kws={'line_kws':{'color':'red'}, 'color':'red'}, 'color':
              plt.show()
```



Residual check

In [236]:

```
residuals = yTrain - y_preds #assumption check is always done on training data
```

In [237]:

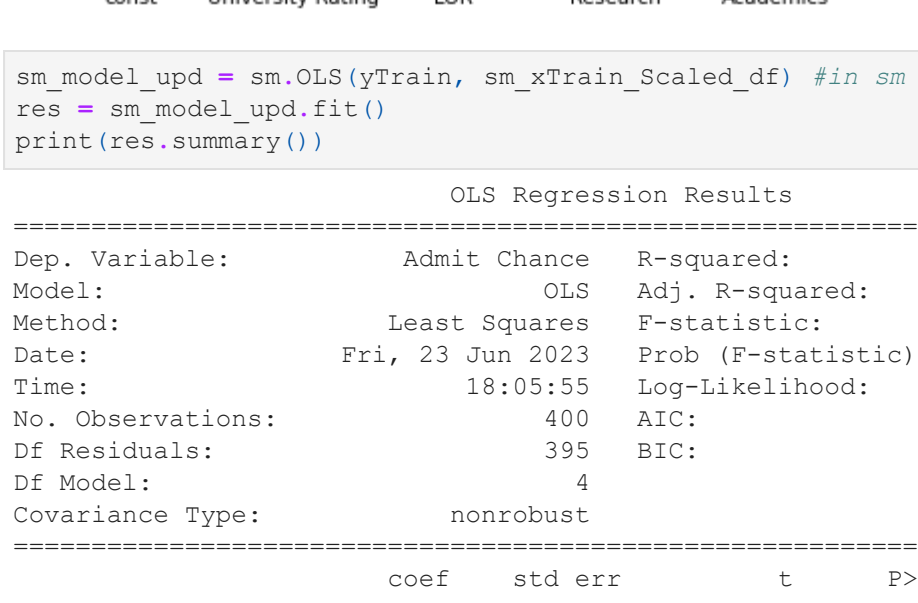
```
print("Mean of Residuals:", np.round(residuals.mean(), 3))
print("Std.dev Residuals:", np.round(residuals.std(), 1))
```

Mean of Residuals: -0.0

Std.dev Residuals: 0.1

In [238]:

```
plt.figure(figsize = (7, 5))
sns.scatterplot(x = residuals, y = residuals, scatter_kws = {"color": "black", "alpha": 0.5},
               line_kws = {"color": "red", "ci = 95}).set(xlabel = "Predicted Values", ylabel = "Residuals")
plt.show()
```



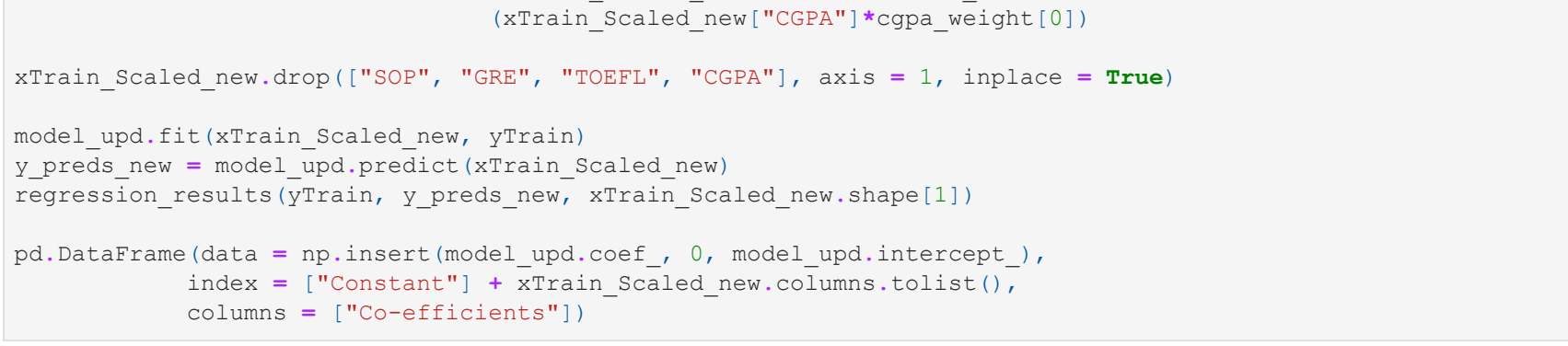
In [239]:

```
fig, axs = plt.subplots(nrows = 1, ncols = 2, figsize = (7,5))
sns.residplot(x = residuals, y = residuals, line_kws = {"color": "red", "ci = 95"}, ax = axs[0])
# s = standardized line, the expected order statistics are scaled by the standard deviation of the
# given sample and have the mean added to them
fig.supylabel("Normality of error terms/residuals", x = 1, y = 1, fontsize = 16)
```

plt.subplots_adjust(right = 2)

plt.show()

Normality of error terms/residuals



In [240]:

```
from statsmodels.stats.diagnostic import het_breuschpagan
from statsmodels.stats.diagnostic import het_white

white_test = het_white(residuals, sm_xTrain_Scaled_df)
bp_test = het_breuschpagan(residuals, sm_xTrain_Scaled_df)

labels = ['LM Statistic', 'LM-Test p-value', 'F-Statistic', 'F-Test p-value']
white_res = dict(zip(labels, bp_test))
bp_res = dict(zip(labels, white_test))

if (white_res['F-Test p-value'] < 0.05 or bp_res['F-Test p-value'] < 0.05):
    print("Heteroskedastic")
    print("White_res: F-Test p-value = {}".format(white_res['F-Test p-value']))
    print("bp_res: F-Test p-value = {}".format(bp_res['F-Test p-value']))
else:
    print("Not Heteroskedastic")
```

Heteroskedastic

white_res: F-Test p-value = 0.0022201238647138926

bp_res: F-Test p-value = 0.0017854422305970082

Insights:

- As we saw in the heatmap, CGPA and GRE are highly correlated with each other. This is one of the major reasons they have a high VIF.
- To deal with high VIF values we can use feature engineering to combine correlated variables into a single correlated variable or drop the highest VIF valued features one-by-one.
- Performing L1 Regularization, we are seeing that GRE, TOEFL, and CGPA are the only features selected. This means we should not drop them.
- Model seems to think SOP is a useless feature based on its co-efficient, we can drop it and check results

Proposed Solution:

- We will drop SOP column**
- In our case, GRE, TOEFL and CGPA are highly correlated features. We know that all these features test the academic proficiency of an aspirant. Hence, we can **combine them into one feature**. To do this, **I will use the feature importances given by L1 regularization** as it only selects the best features

Feature Engineering & Re-Building Model

In [241]:

```
total = l1_coefs.loc["GRE"] + l1_coefs.loc["TOEFL"] + l1_coefs.loc["CGPA"]
gre_weight = l1_coefs.loc["CGPA"]/total
toefl_weight = l1_coefs.loc["GRE"]/(total)
toefl_weight = l1_coefs.loc["TOEFL"]/(total)

sm_xTrain_Scaled_df["Academics"] = (sm_xTrain_Scaled_df["GRE"]*gre_weight[0]) + \
                                   (sm_xTrain_Scaled_df["TOEFL"]*toefl_weight[0]) + \
                                   (sm_xTrain_Scaled_df["CGPA"]*cpga_weight[0])
sm_xTrain_Scaled_df.drop(["SOP", "GRE", "TOEFL", "CGPA"], axis = 1, inplace = True)
pd.DataFrame(data = [variance_inflation_factor(sm_xTrain_Scaled_df, i) for i in range(0, sm_xTrain_Scaled_df.shape[1])],
             index = sm_xTrain_Scaled_df.columns.tolist(),
             columns = ["VIF"])
```

Out[241]:

VIF	
const	1.000000
University Rating	2.726221
LOR	1.826997
Academics	1.400929
Research	2.659026

In [242]:

```
plt.figure(figsize = (9,6))
sns.heatmap(data = sm_xTrain_Scaled_df.corr(), cmap="Blues", annot = True)
plt.show()
```

