

# Grade Determination due to Socio-Demographic information of Students

The data was obtained in a survey of students math courses in secondary school. It contains a lot of interesting social, gender and study information about students. It is pulled from Kaggle data set source -

<https://www.kaggle.com/datasets/uciml/student-alcohol-consumption>  
(<https://www.kaggle.com/datasets/uciml/student-alcohol-consumption>)

The goal of this analysis and model building analysis is to which socio-demographic properties of students influence the grades of the students and if we can do something about it to influence the grades of the students.

We will start with the data loading and start exploring the dataset and observe what is the type of the data and which model may suit the analysis and explore its performance along with the multiple models and compare its performance along with the naive model or mean models.

If the model works well and is able to identify more accurate results, it will be great to use it and make impact to the life of the students.

```
In [1]: 1 # Import the libraries needed for the project
2 from pathlib import Path
3 import pandas as pd
4 import numpy as np
5 import seaborn as sns
6 import dmba
7 import matplotlib.pyplot as plt
8 %matplotlib inline
9 # Library to show all the columns in jupyter notebook
10 from IPython.display import display
11
12 import ipywidgets as widgets
13 from sklearn import preprocessing
14 from sklearn.linear_model import LinearRegression
15 from sklearn.tree import DecisionTreeClassifier, DecisionTreeRegressor
16 from sklearn.preprocessing import MinMaxScaler
17 from sklearn.model_selection import train_test_split, GridSearchCV
18 from dmba import regressionSummary, exhaustive_search, classificationSummary
19 from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
20 from dmba import backward_elimination, forward_selection, stepwise_selection
21 from dmba import adjusted_r2_score, AIC_score, BIC_score
22 from sklearn.neighbors import NearestNeighbors, KNeighborsClassifier
23 from sklearn.metrics import accuracy_score
24 from sklearn.linear_model import LogisticRegression, LogisticRegressionCV
25 from dmba import classificationSummary, gainsChart, liftChart
26
27 from sklearn.metrics import pairwise
28 from scipy.cluster.hierarchy import dendrogram, linkage, fcluster
29 from sklearn.cluster import KMeans
```

no display found. Using non-interactive Agg backend

## Load the data set to the dataframe and explore the dataset

```
In [2]: 1 # Set the variable to display the columns while running the code for
2 temp_data_load = pd.read_csv("student-por.csv")
3 pd.set_option('display.max_columns', None)
4 pd.set_option('display.max_rows', None)
5 temp_data_load.head()
```

Out [2]:

	school	sex	age	address	famsize	Pstatus	Medu	Fedu	Mjob	Fjob	reason	guardian
0	GP	F	18	U	GT3	A	4	4	at_home	teacher	course	mother
1	GP	F	17	U	GT3	T	1	1	at_home	other	course	father
2	GP	F	15	U	LE3	T	1	1	at_home	other	other	mother
3	GP	F	15	U	GT3	T	4	2	health	services	home	mother
4	GP	F	16	U	GT3	T	3	3	other	other	home	father

## Description of Variable and Data Understanding

The dataset has the below attributes and the information about those attributes is mentioned below.

Our area of interest is to understand the final grades i.e. G3 for the students and predict with more accuracy and see if could influence the grades of students who are lagging behind.

school – student's school (binary: 'GP' – Gabriel Pereira or 'MS' – Mousinho da Silveira)  
 sex – student's sex (binary: 'F' – female or 'M' – male)  
 age – student's age (numeric: from 15 to 22)  
 address – student's home address type (binary: 'U' – urban or 'R' – rural)  
 famsize – family size (binary: 'LE3' – less or equal to 3 or 'GT3' – greater than 3)  
 Pstatus – parent's cohabitation status (binary: 'T' – living together or 'A' – apart)  
 Medu – mother's education (numeric: 0 – none, 1 – primary education (4th grade), 2 – 5th to 9th grade, 3 – secondary education or 4 – higher education)  
 Fedu – father's education (numeric: 0 – none, 1 – primary education (4th grade), 2 – 5th to 9th grade, 3 – secondary education or 4 – higher education)  
 Mjob – mother's job (nominal: 'teacher', 'health' care related, civil 'services' (e.g. administrative or police), 'at\_home' or 'other')

## Exploratory Data Analysis and Data Cleaning

### Transaction insights and understanding data ranges

The uniqueness data analysis on the data shows that most of the attributes are categorical or ordinal with absence and grades being only continuous variables.

Also G1, G2 G3 are three grades with G3 being the final grades.

It gets tricky to identify all the grades and hence we will just focus on the final grades.

```
In [3]: 1 temp_data_load.apply(lambda x: x.unique())
```

```
Out[3]: school                                [GP, MS]
sex                                              [F, M]
age                                [18, 17, 15, 16, 19, 22, 20, 21]
address                                [U, R]
famsize                                [GT3, LE3]
Pstatus                                [A, T]
Medu                                [4, 1, 3, 2, 0]
Fedu                                [4, 1, 2, 3, 0]
Mjob    [at_home, health, other, services, teacher]
Fjob    [teacher, other, services, health, at_home]
reason                                [course, other, home, reputation]
guardian                [mother, father, other]
traveltime                [2, 1, 3, 4]
studytime                [2, 3, 1, 4]
failures                [0, 3, 1, 2]
schoolsup                [yes, no]
famsup                [no, yes]
paid                [no, yes]
activities                [no, yes]
nursery                [yes, no]
higher                [yes, no]
internet                [no, yes]
romantic                [no, yes]
famrel                [4, 5, 3, 1, 2]
freetime                [3, 2, 4, 1, 5]
goout                [4, 3, 2, 1, 5]
Dalc                [1, 2, 5, 3, 4]
Walc                [1, 3, 2, 4, 5]
health                [3, 5, 1, 2, 4]
absences    [4, 2, 6, 0, 10, 8, 16, 14, 1, 12, 24, 22, 32,...]
G1    [0, 9, 12, 14, 11, 13, 10, 15, 17, 8, 16, 18, ...]
G2    [11, 13, 14, 12, 16, 17, 8, 10, 15, 9, 7, 6, 1...]
G3    [11, 12, 14, 13, 17, 15, 7, 10, 16, 9, 8, 18, ...]
dtype: object
```

## Datatype Disctionary for data Loaded

Define a dictionary to mannually define the datatypes so we could later change it and make it more consistent to process

```

In [4]: #1 Define column data types
column_dtype_dict = {'school': 'category', 'sex': 'category', 'age': 'category', 'famsize': 'category', 'Pstatus': 'category', 'Medu': 'category', 'Fedu': 'category', 'Mjob': 'category', 'Fjob': 'category', 'reason': 'category', 'guardian': 'category', 'traveltime': 'category', 'studytime': 'category', 'failures': 'category', 'schoolsup': 'category', 'famsup': 'category', 'paid': 'category', 'activities': 'category', 'nursery': 'category', 'higher': 'category', 'internet': 'category', 'romantic': 'category', 'famrel': 'category', 'freetime': 'category', 'goout': 'category', 'Dalc': 'category', 'Walc': 'category', 'health': 'category', 'absences': 'int64', 'G1': 'int64', 'G2': 'int64', 'G3': 'int64'}
temp_data_load = temp_data_load.astype(column_dtype_dict)
temp_data_load.dtypes

```

```

Out[4]: school      object
sex      object
age      int64
address   object
famsize   object
Pstatus   object
Medu      int64
Fedu      int64
Mjob      object
Fjob      object
reason     object
guardian   object
traveltime int64
studytime  int64
failures   int64
schoolsup  object
famsup     object
paid       object
activities object
nursery    object
higher     object
internet   object
romantic   object
famrel     int64
freetime   int64
goout      int64
Dalc       int64
Walc       int64
health     int64
absences   int64
G1         int64
G2         int64
G3         int64
dtype: object

```

## Summarized facts of the data loaded like Mean, unique, counts, Std and percentiles

```
In [5]: 1 temp_data_load.describe(include='all').transpose()
```

Out [5]:

	count	unique	top	freq	mean	std	min	25%	50%	75%	max
<b>school</b>	649	2	GP	423	NaN	NaN	NaN	NaN	NaN	NaN	NaN
<b>sex</b>	649	2	F	383	NaN	NaN	NaN	NaN	NaN	NaN	NaN
<b>age</b>	649.0	NaN	NaN	NaN	16.744222	1.218138	15.0	16.0	17.0	18.0	22.0
<b>address</b>	649	2	U	452	NaN	NaN	NaN	NaN	NaN	NaN	NaN
<b>famsize</b>	649	2	GT3	457	NaN	NaN	NaN	NaN	NaN	NaN	NaN
<b>Pstatus</b>	649	2	T	569	NaN	NaN	NaN	NaN	NaN	NaN	NaN
<b>Medu</b>	649.0	NaN	NaN	NaN	2.514638	1.134552	0.0	2.0	2.0	4.0	4.0
<b>Fedu</b>	649.0	NaN	NaN	NaN	2.306626	1.099931	0.0	1.0	2.0	3.0	4.0
<b>Mjob</b>	649	5	other	258	NaN	NaN	NaN	NaN	NaN	NaN	NaN
<b>Fjob</b>	649	5	other	367	NaN	NaN	NaN	NaN	NaN	NaN	NaN
<b>reason</b>	649	4	course	285	NaN	NaN	NaN	NaN	NaN	NaN	NaN
<b>guardian</b>	649	3	mother	455	NaN	NaN	NaN	NaN	NaN	NaN	NaN
<b>traveltime</b>	649.0	NaN	NaN	NaN	1.568567	0.74866	1.0	1.0	1.0	2.0	4.0
<b>studytime</b>	649.0	NaN	NaN	NaN	1.930663	0.82951	1.0	1.0	2.0	2.0	4.0
<b>failures</b>	649.0	NaN	NaN	NaN	0.22188	0.593235	0.0	0.0	0.0	0.0	3.0
<b>schoolsup</b>	649	2	no	581	NaN	NaN	NaN	NaN	NaN	NaN	NaN
<b>famsup</b>	649	2	yes	398	NaN	NaN	NaN	NaN	NaN	NaN	NaN
<b>paid</b>	649	2	no	610	NaN	NaN	NaN	NaN	NaN	NaN	NaN
<b>activities</b>	649	2	no	334	NaN	NaN	NaN	NaN	NaN	NaN	NaN
<b>nursery</b>	649	2	yes	521	NaN	NaN	NaN	NaN	NaN	NaN	NaN
<b>higher</b>	649	2	yes	580	NaN	NaN	NaN	NaN	NaN	NaN	NaN
<b>internet</b>	649	2	yes	498	NaN	NaN	NaN	NaN	NaN	NaN	NaN
<b>romantic</b>	649	2	no	410	NaN	NaN	NaN	NaN	NaN	NaN	NaN
<b>famrel</b>	649.0	NaN	NaN	NaN	3.930663	0.955717	1.0	4.0	4.0	5.0	5.0
<b>freetime</b>	649.0	NaN	NaN	NaN	3.180277	1.051093	1.0	3.0	3.0	4.0	5.0
<b>goout</b>	649.0	NaN	NaN	NaN	3.1849	1.175766	1.0	2.0	3.0	4.0	5.0
<b>Dalc</b>	649.0	NaN	NaN	NaN	1.502311	0.924834	1.0	1.0	1.0	2.0	5.0
<b>Walc</b>	649.0	NaN	NaN	NaN	2.280431	1.28438	1.0	1.0	2.0	3.0	5.0
<b>health</b>	649.0	NaN	NaN	NaN	3.53621	1.446259	1.0	2.0	4.0	5.0	5.0
<b>absences</b>	649.0	NaN	NaN	NaN	3.659476	4.640759	0.0	0.0	2.0	6.0	32.0
<b>G1</b>	649.0	NaN	NaN	NaN	11.399076	2.745265	0.0	10.0	11.0	13.0	19.0
<b>G2</b>	649.0	NaN	NaN	NaN	11.570108	2.913639	0.0	10.0	11.0	13.0	19.0
<b>G3</b>	649.0	NaN	NaN	NaN	11.906009	3.230656	0.0	10.0	12.0	14.0	19.0



## Limiting the Scope of Analysis for Dependent Variable and Defining the dependent Variable

As we discussed above we are going to limit our analysis to G3 only we can drop rest of the grade columns and focus on the final grades

```
In [6]: 1 temp_data_load = temp_data_load.drop(['G1', 'G2'], axis =1)
        2 temp_data_load.head()
```

Out [6]:

	school	sex	age	address	famsize	Pstatus	Medu	Fedu	Mjob	Fjob	reason	guardian
0	GP	F	18	U	GT3	A	4	4	at_home	teacher	course	mother
1	GP	F	17	U	GT3	T	1	1	at_home	other	course	father
2	GP	F	15	U	LE3	T	1	1	at_home	other	other	mother
3	GP	F	15	U	GT3	T	4	2	health	services	home	mother
4	GP	F	16	U	GT3	T	3	3	other	other	home	father

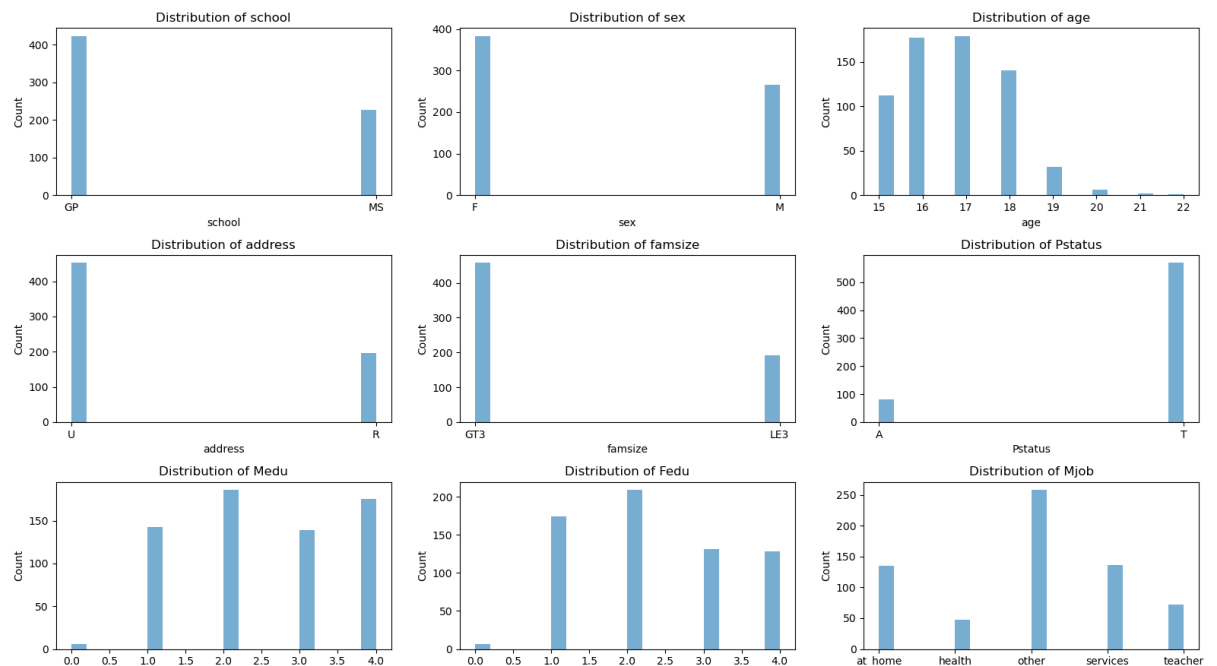
## Distribution Analysis of the data

Here we will observe the trend in data fields and see if we have any insights which are not obvious and give us any information that could be useful in identifying the dataset.

```

In [7]: 1 # Calculate the number of rows needed to display the subplots
2 num_cols = len(temp_data_load.columns)
3 num_rows = (num_cols - 1) // 3 + 1
4
5 # Create a figure and subplots for each column
6 fig, axes = plt.subplots(nrows=num_rows, ncols=3, figsize=(16, 3*num
7
8 # Plot the distribution graphs for each column
9 for i, column in enumerate(temp_data_load.columns):
10     row_idx = i // 3
11     col_idx = i % 3
12     axes[row_idx, col_idx].hist(temp_data_load[column], bins=20, den
13     axes[row_idx, col_idx].set_title(f'Distribution of {column}')
14     axes[row_idx, col_idx].set_xlabel(column)
15     axes[row_idx, col_idx].set_ylabel('Count')
16
17 # Adjust the spacing between subplots for better presentation
18 plt.tight_layout()
19
20 # Show the plots
21 plt.show()
22

```



## One point brief summary of observations

1. School – GP has almost twice the students against the MS
2. Number of female students are slightly higher than Male students
3. There are only few students are older than 19–22 and most of the students fall in the range of 15–18
4. Most of the students are Urban
5. Many of the students have larger families in comparison to smaller families
6. Most of the students parents are staying together
7. With the absence counts in low ranges it means that most of the students prefer to come to college or take fewer leaves.
8. Weekend Alcohol and Daily alcohol consumption is always true and in certain cases the students drink 5 times weekend as well.
9. Students often go out mean mean around 3 times a week.

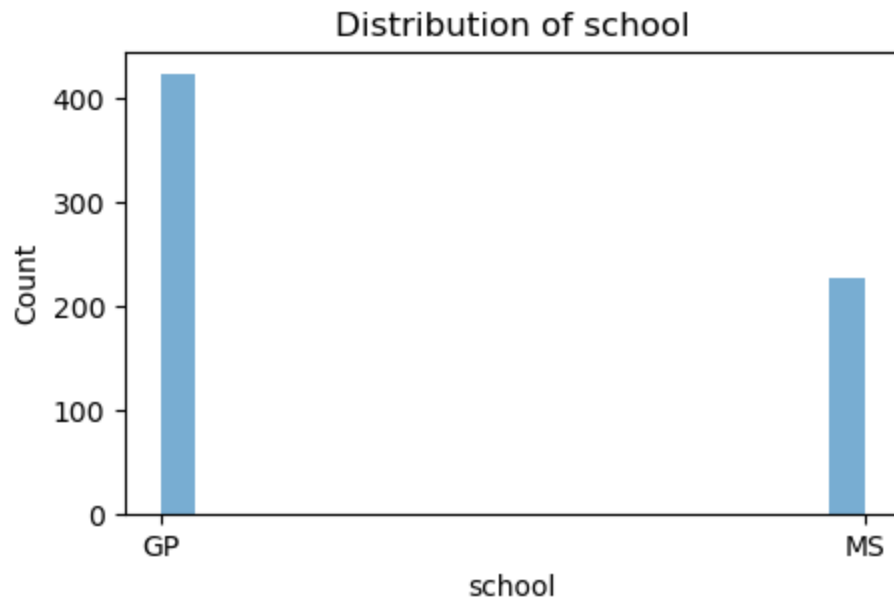
**Below is the widget model to show us the above graphs in more clarity and you can choose the attribute.**

**Use the Widget to Select a column for its Graph for better view and option to Select**

```
In [8]: 1 def plot_distribution(column):
2         plt.figure(figsize=(5, 3))
3         plt.hist(temp_data_load[column], bins=20, density=False, alpha=0.5)
4         plt.title(f'Distribution of {column}')
5         plt.xlabel(column)
6         plt.ylabel('Count')
7         plt.show()
8
9         # Create a selection widget to choose the column
10        column_selector = widgets.Dropdown(
11            options=temp_data_load.columns,
12            description='Choose a column:',
13            disabled=False,
14        )
15
16        # Apply the plotting function when the widget value changes
17        widgets.interactive(plot_distribution, column=column_selector)
```

Out [8]:

Choose a c...

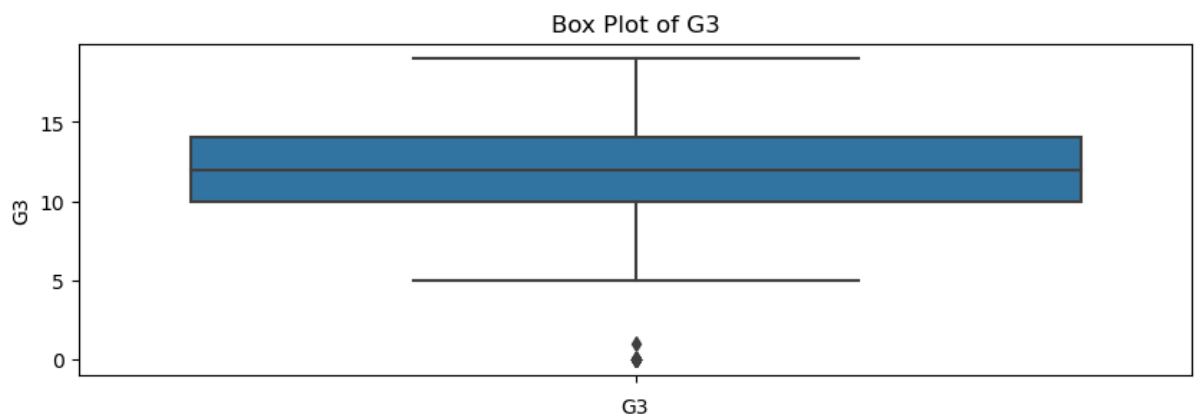


## G3 Grade distribution using Boxplot

It seems we have some outliers in the data. It means only few students have very low grades. This will make the matrix have bias towards those grades and our will have shift in the trend due to these outliers.

We can remove those records and do the data analysis on the pending records.

```
In [9]: 1 # boxplot the G3 and understand the distribution of the data. We can
2
3 plt.figure(figsize=(10, 3))
4 sns.boxplot(y=temp_data_load['G3'])
5 plt.title('Box Plot of G3')
6 plt.xlabel('G3')
7 plt.show()
```



## Remove the outlier G3 records from the data

Let's drop if there is any outlier in the dependent data to avoid any skewenees in the model.

```
In [10]: 1 # Function to remove outliers based on box plot
2 def remove_outliers(df, column, multiplier=1.5):
3     Q1 = temp_data_load[column].quantile(0.25)
4     Q3 = temp_data_load[column].quantile(0.75)
5     IQR = Q3 - Q1
6
7     lower_bound = Q1 - multiplier * IQR
8     upper_bound = Q3 + multiplier * IQR
9
10    return temp_data_load[(temp_data_load[column] >= lower_bound) &
11                           (temp_data_load[column] <= upper_bound)]
12 # Specify columns for outlier removal
13 columns_to_check = ['G3']
14
15 # Remove outliers for each column
16 for column in columns_to_check:
17     new_temp_data = remove_outliers(temp_data_load, column).reset_index(drop=True)
18
19 print(len(new_temp_data))
20 temp_data_load = new_temp_data.drop('index', axis = 1)
21 print(len(temp_data_load))
22 temp_data_load.head()
```

633

633

Out[10]:

	school	sex	age	address	famsize	Pstatus	Medu	Fedu	Mjob	Fjob	reason	guardian
0	GP	F	18	U	GT3	A	4	4	at_home	teacher	course	mother
1	GP	F	17	U	GT3	T	1	1	at_home	other	course	father
2	GP	F	15	U	LE3	T	1	1	at_home	other	other	mother
3	GP	F	15	U	GT3	T	4	2	health	services	home	mother
4	GP	F	16	U	GT3	T	3	3	other	other	home	father

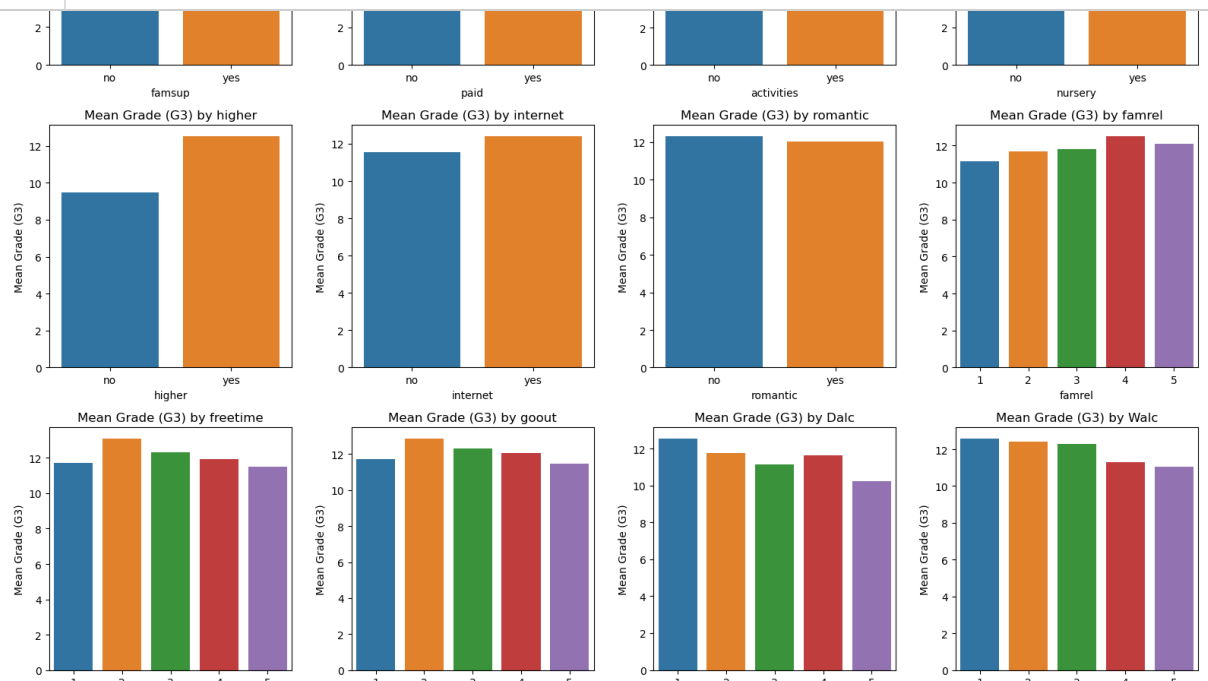
## Mean Grade for Each of the Attributes to understand the trend

1. Not surprisingly mean grade are lower for students who have age beyond 19 years. This may be because of the fact that there might be some past failures and struggling to keep up with the grades. We can explicitly validate these with the data set.
2. Surprisingly the Mean grade of people with very high absences are high which is odd. Further data exploration is needed. For now we are restricting ourselves to not get overwhelmed with the data analysis.
3. Also it looks like the mean grade have negative trend with the alcohol consumption

```

In [11]: 1 # Calculate the number of rows and columns needed for the subplots
2 num_cols = len(temp_data_load.columns)
3 num_rows = (num_cols - 1) // 4 + 1
4
5 # Create a figure and subplots for each column
6 fig, axes = plt.subplots(nrows=num_rows, ncols=4, figsize=(16, 4*num
7
8 # Iterate through each column and create subplots
9 for i, column in enumerate(temp_data_load.columns):
10     row_idx = i // 4
11     col_idx = i % 4
12
13     # Calculate the mean of G3_mat for each category in the current column
14     mean_grades_by_column = temp_data_load.groupby(column)['G3'].mean()
15     #mean_grades_by_column = temp_data_load['G3']
16
17     # Plot the bar plot for G3_mat mean by category in the current column
18     sns.barplot(x=mean_grades_by_column.index, y=mean_grades_by_column)
19     axes[row_idx, col_idx].set_title(f'Mean Grade (G3) by {column}')
20     axes[row_idx, col_idx].set_xlabel(column)
21     axes[row_idx, col_idx].set_ylabel('Mean Grade (G3)')
22
23 # Adjust spacing between subplots for better presentation
24 plt.tight_layout()
25
26 # Show the plots
27 plt.show()

```



```

In [12]: 1 # Analysis of Age =22 records
2 temp_data_load[temp_data_load['age']==22]

```

Out[12]:

	school	sex	age	address	famsize	Pstatus	Medu	Fedu	Mjob	Fjob	reason	guardia
277	GP	M	22	U	GT3	T	3	1	services	services	other	moth

## Analyze the Age 22 record

Checking the record for age = 22 to see why he has lower grade and it seems he has lower grades due to multitude of factors like higher going out count, daily and weekend alcohol consumption, his reason of joining the school, higher absences and higher amount of past failures. All these factors signals not good grades. However actuation and cautation cannot be concluded here so we will leave it here.

In [13]: `1 temp_data_load.head()`

Out[13]:

	school	sex	age	address	famsize	Pstatus	Medu	Fedu	Mjob	Fjob	reason	guardian
0	GP	F	18	U	GT3	A	4	4	at_home	teacher	course	mother
1	GP	F	17	U	GT3	T	1	1	at_home	other	course	father
2	GP	F	15	U	LE3	T	1	1	at_home	other	other	mother
3	GP	F	15	U	GT3	T	4	2	health	services	home	mother
4	GP	F	16	U	GT3	T	3	3	other	other	home	father

## Analyze the Higher Absence Records

Looking at this information nothing could be concluded however only common thing is there access to internet, common goal of higher education and family education support and quality of family relationship being very high in most of the cases.

Yes there are few have difference however they may be able compensate on other front.

This analysis is observational and and its causation could not be established yet. So we will leave it here for now.

In [14]: `1 temp_data_load[temp_data_load['absences']>=20]`

Out[14]:

	school	sex	age	address	famsize	Pstatus	Medu	Fedu	Mjob	Fjob	reason	guardian
150	GP	F	15	U	GT3	A	3	3	services	services	home	rr
155	GP	M	17	U	GT3	T	2	1	other	other	home	rr
195	GP	F	17	U	LE3	T	3	3	other	other	reputation	rr
210	GP	F	17	U	GT3	T	4	4	services	teacher	home	rr
215	GP	F	17	R	GT3	T	2	2	other	other	reputation	rr
254	GP	M	18	U	GT3	T	2	2	other	at_home	course	
323	GP	M	17	U	LE3	A	4	1	services	other	home	rr
411	GP	M	21	R	LE3	T	1	1	at_home	other	course	



## Cleaning the data and defining the datatypes.

In the below step we will define the datatype, convert the values to numerical values to be fed in data models and also create the dummies for categorical columns.

```
In [15]: 1 temp_data_load['school'] = temp_data_load['school'].astype('category')
2 temp_data_load['sex'] = temp_data_load['sex'].astype('category')
3 temp_data_load['address'] = temp_data_load['address'].astype('category')
4 temp_data_load['famsize'] = temp_data_load['famsize'].astype('category')
5 temp_data_load['Pstatus'] = temp_data_load['Pstatus'].astype('category')
6 temp_data_load['Mjob'] = temp_data_load['Mjob'].astype('category')
7 temp_data_load['Fjob'] = temp_data_load['Fjob'].astype('category')
8 temp_data_load['reason'] = temp_data_load['reason'].astype('category')
9 temp_data_load['guardian'] = temp_data_load['guardian'].astype('category')
10 temp_data_load.head()
11
12
```

Out[15]:

	school	sex	age	address	famsize	Pstatus	Medu	Fedu	Mjob	Fjob	reason	guardian
0	GP	F	18	U	GT3	A	4	4	at_home	teacher	course	mother
1	GP	F	17	U	GT3	T	1	1	at_home	other	course	father
2	GP	F	15	U	LE3	T	1	1	at_home	other	other	mother
3	GP	F	15	U	GT3	T	4	2	health	services	home	mother
4	GP	F	16	U	GT3	T	3	3	other	other	home	father

```
In [16]: 1 # Create the dummies for categorical columns
2 df = pd.get_dummies(temp_data_load, prefix_sep='_', drop_first=True)
3 df.head()
4
```

Out[16]:

	age	Medu	Fedu	travelttime	studytime	failures	famrel	freetime	goout	Dalc	Walc	health
0	18	4	4	2	2	0	4	3	4	1	1	3
1	17	1	1	1	2	0	5	3	3	1	1	3
2	15	1	1	1	2	0	4	3	2	2	3	3
3	15	4	2	1	3	0	3	2	2	1	1	5
4	16	3	3	1	2	0	4	3	2	1	2	5

## Normalize the model the attributes to the Standard value for any continuous variables

```
In [17]: Do the scaler fitting of attributes Age, Medu, Fedu, Travel Time, StudyTime,
1
2
3 scaler_fiting_cols_list = ['age', 'Medu', 'Fedu', 'traveltime', 'studytime', 'failures',
4
5                               'goout', 'Dalc', 'Walc', 'health', 'absences']
6
7 cols_list = ['z_' + col for col in scaler_fiting_cols_list]
8
9 scaler = preprocessing.StandardScaler()
10 scaler.fit(df[scaler_fiting_cols_list])
11
12 df_norm = pd.concat([pd.DataFrame(scaler.transform(df[scaler_fiting_cols_list])), df[['goout', 'Dalc', 'Walc', 'health', 'absences']]])
13 df_norm = df_norm.drop(columns = scaler_fiting_cols_list)
14 df_norm.head()
```

Out[17]:

	z_age	z_Medu	z_Fedu	z_traveltime	z_studytime	z_failures	z_famrel	z_freetime	z_goout
0	1.054848	1.306950	1.524682	0.581688	0.070202	-0.357697	0.071834	-0.161563	0.711568
1	0.230504	-1.341820	-1.199092	-0.752401	0.070202	-0.357697	1.129292	-0.161563	-0.054256
2	-1.418185	-1.341820	-1.199092	-0.752401	0.070202	-0.357697	0.071834	-0.161563	-1.009504
3	-1.418185	1.306950	-0.291167	-0.752401	1.271221	-0.357697	-0.985625	-1.117353	-1.009504
4	-0.593841	0.424026	0.616758	-0.752401	0.070202	-0.357697	0.071834	-0.161563	-1.009504

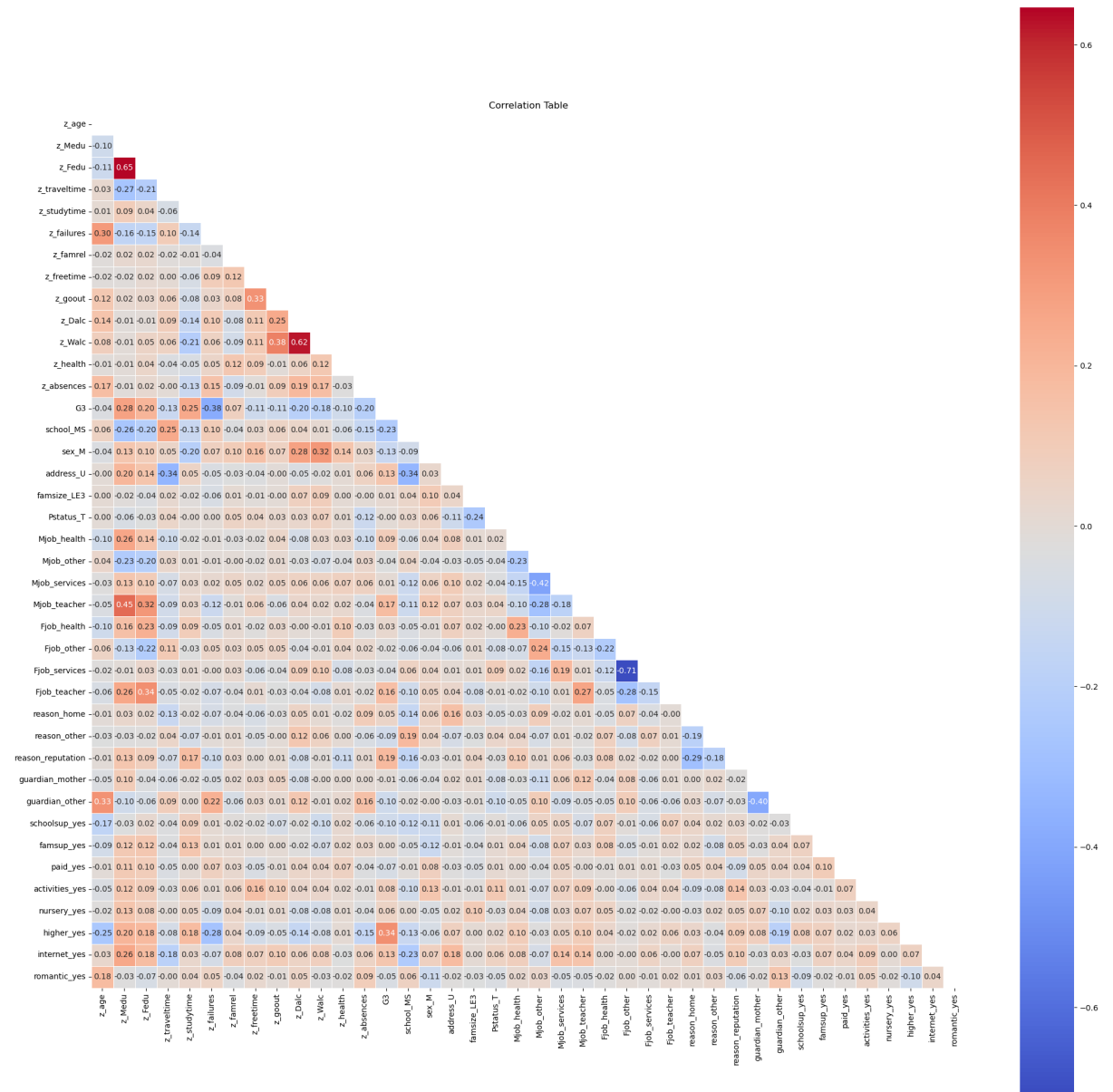
## Check the Corr matrix

Let's see the corealtion among the attributes to identify if we have correlation among the attributes and which attributes are related.

```

In [18]: 1 correlation_matrix = df_norm.corr()
2
3 Create a mask for the upper half of the correlation table
4 mask = np.triu(np.ones_like(correlation_matrix, dtype=bool))
5
6 Plot the correlation table as a heatmap with the upper half masked
7 plt.figure(figsize=(25, 25))
8 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f",
9             plt.title('Correlation Table'))
10 plt.show()

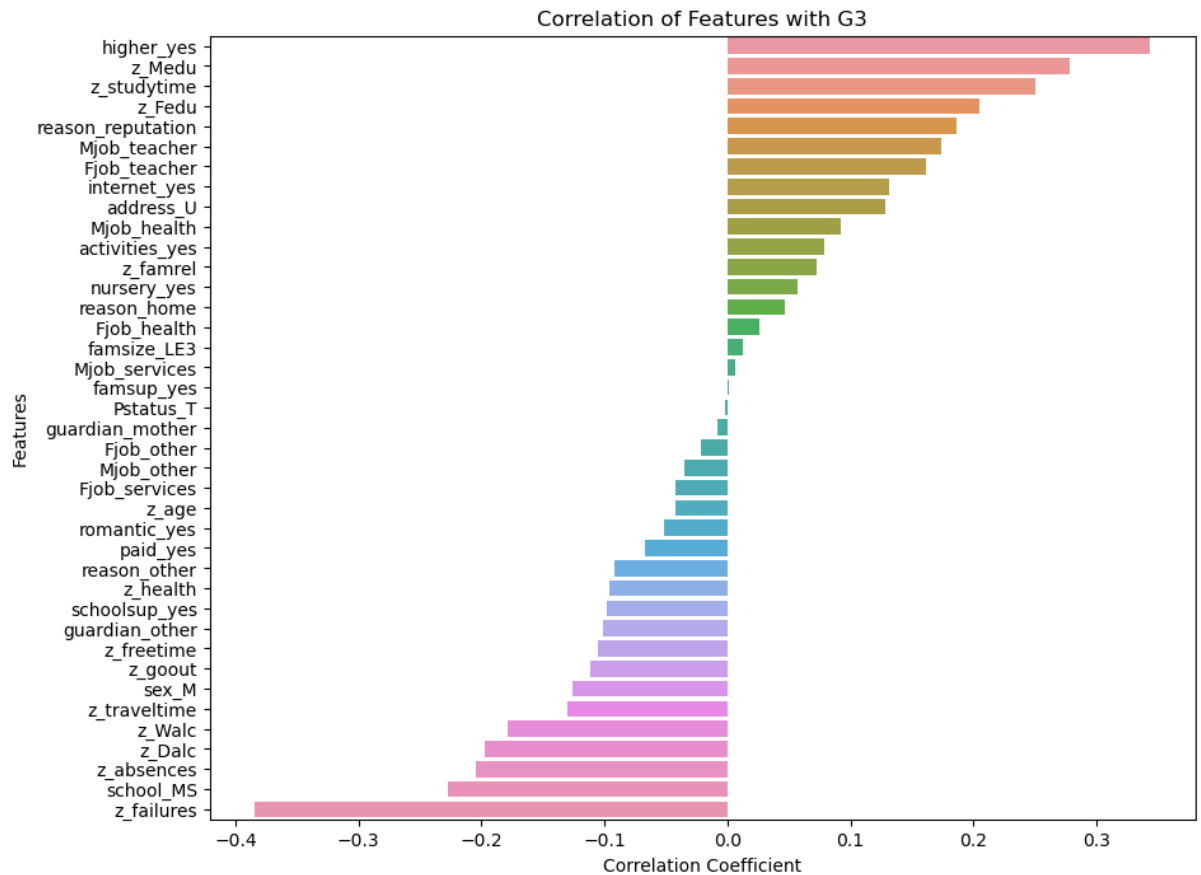
```



## Lets Plot the indexes in sorted order to visualize high impact attributes to G3

From the graph it seems higher education aspiration, Mother's education, Studytime, Father's education, studying for reputation, mother and fathers occupation as teacher, internet access and urban have positive correlation and past failures, school, absences, alcohol consumption,

```
In [19]: 1 correlations = df_norm.corr()['G3'].drop('G3').sort_values(ascending=
2
3 plt.figure(figsize=(10, 8))
4 sns.barplot(x=correlations.values, y=correlations.index)
5 plt.title('Correlation of Features with G3')
6 plt.xlabel('Correlation Coefficient')
7 plt.ylabel('Features')
8 plt.show()
```



## Let's Calculate the VIF index to understand the multicollinearity among the variables.

It is recommended to drop any variables which have VIF higher than 5 or use it with caution as they tend to show multi collinearity and have been influenced by other variables.

```
In [20]: 1 from statsmodels.stats.outliers_influence import variance_inflation_
2
3 # X is the design matrix (independent variables)
4 X = df_norm.drop(['G3'], axis =1)
5 print(X.columns)
6
7 # Calculate VIF for each variable
8 vif_data = pd.DataFrame()
9 vif_data["Variable"] = X.columns
10 vif_data["VIF"] = [variance_inflation_factor(X.values, i) for i in range(X.shape[1])]
11
12 print(vif_data.sort_values(by='VIF',ascending=False))
13
```

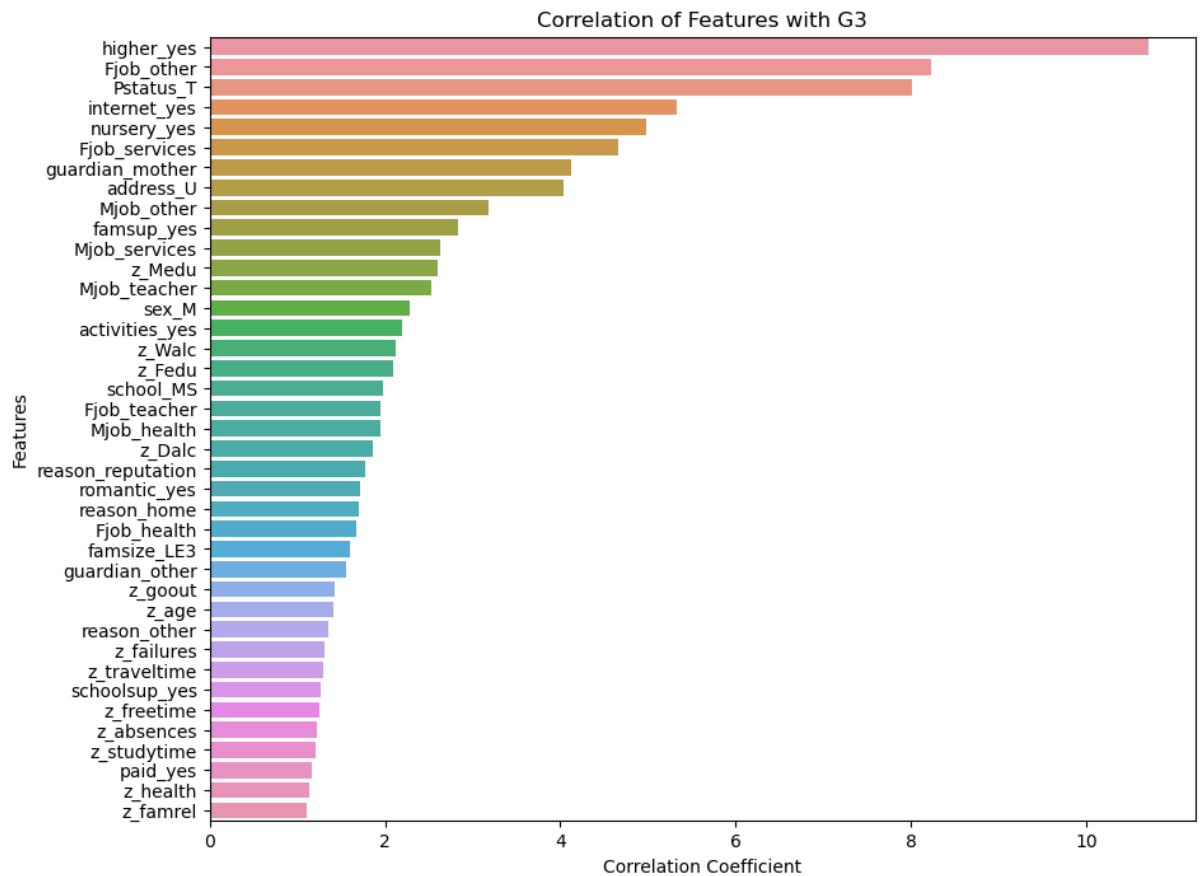
```

Index(['z_age', 'z_Medu', 'z_Fedu', 'z_traveltime', 'z_studytime',
      'z_failures', 'z_famrel', 'z_freetime', 'z_goout', 'z_Dalc', 'z_
Walc',
      'z_health', 'z_absences', 'school_MS', 'sex_M', 'address_U',
      'famsize_LE3', 'Pstatus_T', 'Mjob_health', 'Mjob_other',
      'Mjob_services', 'Mjob_teacher', 'Fjob_health', 'Fjob_other',
      'Fjob_services', 'Fjob_teacher', 'reason_home', 'reason_other',
      'reason_reputation', 'guardian_mother', 'guardian_other',
      'schoolsup_yes', 'famsup_yes', 'paid_yes', 'activities_yes',
      'nursery_yes', 'higher_yes', 'internet_yes', 'romantic_yes'],
      dtype='object')

```

	Variable	VIF
36	higher_yes	10.707842
23	Fjob_other	8.231278
17	Pstatus_T	8.011992
37	internet_yes	5.334340
35	nursery_yes	4.984757
24	Fjob_services	4.655083
29	guardian_mother	4.130585
15	address_U	4.044509
19	Mjob_other	3.181300
32	famsup_yes	2.830123
20	Mjob_services	2.633578
1	z_Medu	2.608084
21	Mjob_teacher	2.522472
14	sex_M	2.276130
34	activities_yes	2.198683
10	z_Walc	2.125139
2	z_Fedu	2.087526
13	school_MS	1.976517
25	Fjob_teacher	1.954323
18	Mjob_health	1.946154
9	z_Dalc	1.857122
28	reason_reputation	1.775952
38	romantic_yes	1.718174
26	reason_home	1.705876
22	Fjob_health	1.669730
16	famsize_LE3	1.600328
30	guardian_other	1.555884
8	z_goout	1.430683
0	z_age	1.408341
27	reason_other	1.349793
5	z_failures	1.306807
3	z_traveltime	1.291682
31	schoolsup_yes	1.270769
7	z_freetime	1.258458
12	z_absences	1.221183
4	z_studytime	1.214419
33	paid_yes	1.170930
11	z_health	1.140784
6	z_famrel	1.114400

```
In [21]: 1 # Plot the VIF factor in descending order.
2 vif_data_sorted = vif_data.sort_values(by='VIF', ascending=False)
3 plt.figure(figsize=(10, 8))
4 sns.barplot(x=vif_data_sorted.VIF, y=vif_data_sorted.Variable)
5 plt.title('Correlation of Features with G3')
6 plt.xlabel('Correlation Coefficient')
7 plt.ylabel('Features')
8 plt.show()
```



```
In [22]: 1 df_norm.columns
```

```
Out[22]: Index(['z_age', 'z_Medu', 'z_Fedu', 'z_traveltime', 'z_studytime',
               'z_failures', 'z_famrel', 'z_freetime', 'z_gooout', 'z_Dalc', 'z_
               Walc',
               'z_health', 'z_absences', 'G3', 'school_MS', 'sex_M', 'address_
               U',
               'famsize_LE3', 'Pstatus_T', 'Mjob_health', 'Mjob_other',
               'Mjob_services', 'Mjob_teacher', 'Fjob_health', 'Fjob_other',
               'Fjob_services', 'Fjob_teacher', 'reason_home', 'reason_other',
               'reason_reputation', 'guardian_mother', 'guardian_other',
               'schoolsup_yes', 'famsup_yes', 'paid_yes', 'activities_yes',
               'nursery_yes', 'higher_yes', 'internet_yes', 'romantic_yes'],
              dtype='object')
```

```

In [23]: 1 # Split the data in Training Data and Validation Data trainData, valid
2
3 num_bins = 4
4
5 # Create a new categorical column based on equal-frequency bins
6 df_norm['G3_bin'] = pd.qcut(df_norm['G3'], q=num_bins, labels=False,
7 print((df_norm.columns))
8 print(df_norm['G3_bin'].value_counts())
9 print(pd.DataFrame(df_norm.groupby('G3_bin')['G3'].agg(['min', 'max',
10
11 predictors = ['z_age', 'z_Medu', 'z_Fedu', 'z_traveltime', 'z_studyt
12               'z_freetime', 'z_goout', 'z_Dalc', 'z_Walc', 'z_health'
13               'famsup_yes', 'paid_yes', 'activities_yes', 'nursery_ye
14               'sex_M', 'address_U', 'famsize_LE3', 'Pstatus_T', 'Mjob
15               'Mjob_teacher', 'Fjob_health', 'Fjob_other', 'Fjob_serv
16               'reason_other', 'reason_reputation', 'guardian_mother',
17 outcome = 'G3'
18 label_predictor = 'G3_bin'
19 X = df_norm[predictors]
20 y = df_norm[label_predictor]
21 train_X, valid_X, train_y, valid_y = train_test_split(X, y, test_size=
22 print('Training : ', train_X.shape)
23 print('Validation : ', valid_X.shape)
24

```

```

Index(['z_age', 'z_Medu', 'z_Fedu', 'z_traveltime', 'z_studytime',
      'z_failures', 'z_famrel', 'z_freetime', 'z_goout', 'z_Dalc', 'z_
Walc',
      'z_health', 'z_absences', 'G3', 'school_MS', 'sex_M', 'address_
U',
      'famsize_LE3', 'Pstatus_T', 'Mjob_health', 'Mjob_other',
      'Mjob_services', 'Mjob_teacher', 'Fjob_health', 'Fjob_other',
      'Fjob_services', 'Fjob_teacher', 'reason_home', 'reason_other',
      'reason_reputation', 'guardian_mother', 'guardian_other',
      'schoolsup_yes', 'famsup_yes', 'paid_yes', 'activities_yes',
      'nursery_yes', 'higher_yes', 'internet_yes', 'romantic_yes', 'G3
_bin'],
      dtype='object')
0    181
1    176
2    145
3    131
Name: G3_bin, dtype: int64
      min max count
G3_bin
0         5  10    181
1        11  12    176
2        13  14    145
3        15  19    131
Training : (411, 39)
Validation : (222, 39)

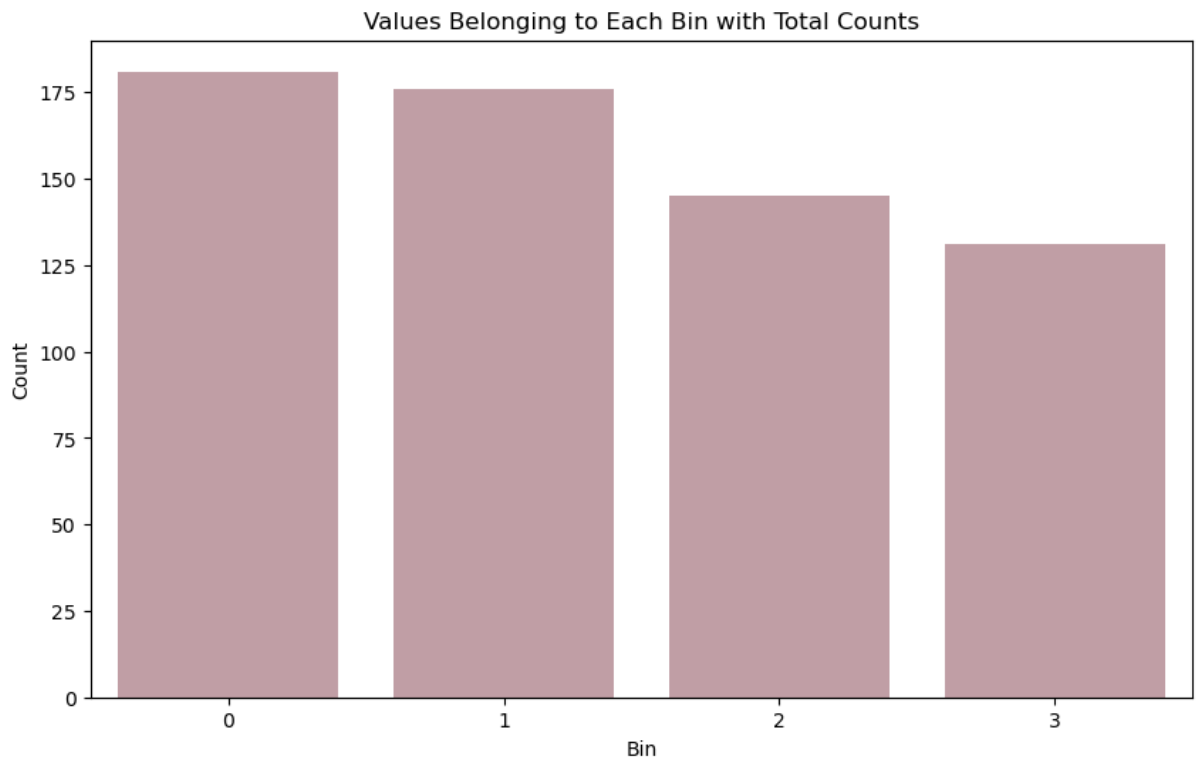
```



```
In [24]: 1 # plot bins information to understand its distribution
2 count_per_bin = df_norm['G3_bin'].value_counts().sort_index()
3 print(count_per_bin)
4 # Plot the information using seaborn
5 plt.figure(figsize=(10, 6))
6 # Grouped bar plot showing which values belong to which bin
7 sns.barplot(x='G3_bin', y='G3', data=df_norm, color='red', errorbar=None)
8 # Stacked bar plot showing the total counts per bin
9 sns.barplot(x=count_per_bin.index, y=count_per_bin.values, color='lightcoral')
10
11 # Adding labels and title
12 plt.title('Values Belonging to Each Bin with Total Counts')
13 plt.xlabel('Bin')
14 plt.ylabel('Count')
15
16 # Display the plot
17 plt.show()
```

```
0    181
1    176
2    145
3    131
```

Name: G3\_bin, dtype: int64



```
In [25]: 1 def train_model(variables):
2         model = LinearRegression()
3         model.fit(train_X[variables], train_y)
4         return model
5
6 def score_model(model, variables):
7         return AIC_score(train_y, model.predict(train_X[variables]), model)
8
9 bestBE_model, bestbe_variables = backward_elimination(train_X.columns, train_y)
10
11
12 print(bestbe_variables)
13
```

Variables: z\_age, z\_Medu, z\_Fedu, z\_traveltime, z\_studytime, z\_failures, z\_famrel, z\_freetime, z\_goout, z\_Dalc, z\_Walc, z\_health, z\_absences, schoolsup\_yes, famsup\_yes, paid\_yes, activities\_yes, nursery\_yes, higher\_yes, internet\_yes, romantic\_yes, school\_MS, sex\_M, address\_U, famsize\_LE3, Pstatus\_T, Mjob\_health, Mjob\_other, Mjob\_services, Mjob\_teacher, Fjob\_health, Fjob\_other, Fjob\_services, Fjob\_teacher, reason\_home, reason\_other, reason\_reputation, guardian\_mother, guardian\_other

Start: score=1118.85

Step: score=1116.85, remove Fjob\_teacher

Step: score=1114.86, remove Mjob\_other

Step: score=1112.87, remove romantic\_yes

Step: score=1110.91, remove guardian\_other

Step: score=1108.96, remove z\_freetime

Step: score=1107.03, remove Mjob\_services

Step: score=1105.18, remove famsup\_yes

Step: score=1103.39, remove reason\_other

Step: score=1101.64, remove address\_U

Step: score=1099.90, remove activities\_yes

Step: score=1098.14, remove nursery\_yes

Step: score=1096.58, remove Mjob\_health

Step: score=1095.02, remove z\_Medu

Step: score=1093.48, remove famsize\_LE3

Step: score=1092.05, remove paid\_yes

Step: score=1090.88, remove z\_traveltime

Step: score=1089.88, remove z\_Walc

Step: score=1088.75, remove z\_Dalc

Step: score=1087.81, remove z\_age

Step: score=1087.01, remove reason\_home

Step: score=1086.33, remove z\_goout

Step: score=1085.86, remove z\_Fedu

Step: score=1085.74, remove internet\_yes

Step: score=1085.67, remove Pstatus\_T

Step: score=1085.55, remove guardian\_mother

Step: score=1085.55, remove None

['z\_studytime', 'z\_failures', 'z\_famrel', 'z\_health', 'z\_absences', 'schoolsup\_yes', 'higher\_yes', 'school\_MS', 'sex\_M', 'Mjob\_teacher', 'Fjob\_health', 'Fjob\_other', 'Fjob\_services', 'reason\_reputation']

In [26]: 1 regressionSummary(valid\_y, bestBE\_model.predict(valid\_X[bestbe\_varial

Regression statistics

Mean Error (ME) : 0.0212  
 Root Mean Squared Error (RMSE) : 0.9816  
 Mean Absolute Error (MAE) : 0.8139

```
In [27]: 1 # BEst Model is
2 col_list = bestbe_variables
3 #outcome = 'G3'
4 grade_lm = LinearRegression()
5 grade_lm.fit(train_X[col_list], train_y)
6
7 # print coefficients
8 print('intercept ', grade_lm.intercept_)
9 print(pd.DataFrame({'Predictor': col_list, 'coefficient': grade_lm.co
10
11 # print performance measures
12 print(regressionSummary(train_y, grade_lm.predict(train_X[col_list])
13
14 pred_y = grade_lm.predict(train_X[col_list])
15
16 print('adjusted r2 : ', adjusted_r2_score(train_y, pred_y, grade_lm)
17 print('AIC : ', AIC_score(train_y, pred_y, grade_lm))
18 print('BIC : ', BIC_score(train_y, pred_y, grade_lm))
19
```

```
intercept 1.2228448518811723
Predictor coefficient
0 z_studytime 0.154383
1 z_failures -0.252755
2 z_famrel 0.064940
3 z_health -0.075270
4 z_absences -0.173099
5 schoolsup_yes -0.601194
6 higher_yes 0.719596
7 school_MS -0.513364
8 sex_M -0.278409
9 Mjob_teacher 0.469715
10 Fjob_health -0.527434
11 Fjob_other -0.257072
12 Fjob_services -0.452015
13 reason_reputation 0.298566
```

Regression statistics

Mean Error (ME) : -0.0000  
 Root Mean Squared Error (RMSE) : 0.8718  
 Mean Absolute Error (MAE) : 0.7144  
 None  
 adjusted r2 : 0.35280277889614864  
 AIC : 1085.54631800051  
 BIC : 1149.8438094324497

```
In [28]: 1 # Validation Data
2
3 # Use predict() to make predictions on a new set
4 grade_lm_pred = grade_lm.predict(valid_X[col_list])
5
6 result = pd.DataFrame({'Predicted': grade_lm_pred, 'Actual': valid_y
7                        'Residual': valid_y - grade_lm_pred})
8 print(result.head(20))
9
10 # Compute common accuracy measures
11 regressionSummary(valid_y, grade_lm_pred)
```

	Predicted	Actual	Residual
329	1.926291	1	-0.926291
247	2.284399	1	-1.284399
390	2.010154	3	0.989846
145	1.520071	0	-1.520071
497	-0.448724	0	0.448724
513	2.052429	3	0.947571
165	0.756830	1	0.243170
77	2.659055	2	-0.659055
534	1.451858	1	-0.451858
163	0.730479	0	-0.730479
271	1.667949	1	-0.667949
31	1.605020	3	1.394980
55	2.004507	1	-1.004507
90	2.017232	1	-1.017232
576	0.229733	0	-0.229733
76	2.555632	1	-1.555632
2	1.133995	1	-0.133995
256	1.878179	3	1.121821
311	1.883826	2	0.116174
334	1.971111	3	1.028889

Regression statistics

```
Mean Error (ME) : 0.0212
Root Mean Squared Error (RMSE) : 0.9816
Mean Absolute Error (MAE) : 0.8139
```

Let us remove the high VIF variables to see the impact on the model Given that the regression model has these attributes ['z\_studytime', 'z\_failures', 'z\_famrel', 'z\_health', 'z\_absences', 'schoolsup\_yes', 'higher\_yes', 'school\_MS', 'sex\_M', 'Mjob\_teacher', 'Fjob\_health', 'Fjob\_other', 'Fjob\_services', 'reason\_reputation']

It makes sense to remove the high VIF index variables to reduce the multicollinearity.

However, we want keep higher\_yes as it means the student have expectations for further studying and we believe it should influence the approach towards their education and reflect seriousness towards education.

But there is debatable counter argument such that it could just be a wrong data as its accuracy could not be predicted and directly associated to the student as its just a future event. Anyone can say they have higher education aspirations and how many follow through with it is

debatable. So we are inclined to take such attributes out of our model which are not current facts and something in future.

Also we are inclined to create one more variable if either of the parent is working more than one job, then what would be the impact. This will mean that student see less of the parent. In order

```
In [29]: 1 columns_to_sum = ['Fjob_health', 'Fjob_other', 'Fjob_services', 'Fjob_
2 temp_data = pd.DataFrame()
3 temp_data['sum_of_Fjob'] = df_norm[columns_to_sum].sum(axis=1)
4
5 df_norm.loc[temp_data['sum_of_Fjob']>=2].head()
```

```
Out[29]:
```

z_age	z_Medu	z_Fedu	z_traveltime	z_studytime	z_failures	z_famrel	z_freetime	z_goout	z_D
-------	--------	--------	--------------	-------------	------------	----------	------------	---------	-----

It appears that there are no data points which indicate that father is working two jobs and hence the impact would be more driven from the fact that which kind of job sector is in. So we will keep fathers job sectors in analysis.

So the updated list becomes - ['z\_studytime', 'z\_failures', 'z\_famrel', 'z\_health', 'z\_absences', 'schoolsup\_yes', 'school\_MS', 'sex\_M', 'Mjob\_teacher', 'Fjob\_health', 'Fjob\_services', 'reason\_reputation']

So we gonna train the model with these and see what the results we get.

```
In [30]: 1 updated_predictors = ['z_studytime', 'z_failures', 'z_famrel', 'z_hes  
2             'schoolsup_yes', 'school_MS', 'sex_M', 'Mjob_t  
3             'Fjob_services', 'reason_reputation']  
4 # Best Model is  
5 col_list = updated_predictors  
6 outcome = 'G3'  
7 vif_grade_lm = LinearRegression()  
8 vif_grade_lm.fit(train_X[col_list], train_y)  
9  
10 # print coefficients  
11 print('intercept ', vif_grade_lm.intercept_)  
12 print(pd.DataFrame({'Predictor': col_list, 'coefficient': vif_grade_  
13  
14 # print performance measures  
15 print(regressionSummary(train_y, vif_grade_lm.predict(train_X[col_li  
16  
17 pred_y = vif_grade_lm.predict(train_X[col_list])  
18  
19 print('adjusted r2 : ', adjusted_r2_score(train_y, pred_y, vif_grade_  
20 print('AIC : ', AIC_score(train_y, pred_y, vif_grade_lm))  
21 print('BIC : ', BIC_score(train_y, pred_y, vif_grade_lm))  
22  
23 # Validation Data  
24  
25 # Use predict() to make predictions on a new set  
26 vif_grade_lm_pred = vif_grade_lm.predict(valid_X[col_list])  
27  
28 result = pd.DataFrame({'Predicted': vif_grade_lm_pred, 'Actual': val  
29                        'Residual': valid_y - vif_grade_lm_pred})  
30 print(result.head(20))  
31  
32 # Compute common accuracy measures  
33 regressionSummary(valid_y, vif_grade_lm_pred)  
34
```

```

intercept  1.6259639843896965
Predictor  coefficient
0          z_studytime      0.186416
1          z_failures      -0.290868
2          z_famrel         0.053646
3          z_health        -0.072550
4          z_absences      -0.197613
5          schoolsup_yes    -0.508348
6          school_MS       -0.543775
7          sex_M           -0.263777
8          Mjob_teacher     0.575553
9          Fjob_health     -0.273077
10         Fjob_services   -0.216441
11         reason_reputation 0.325079

```

Regression statistics

```

Mean Error (ME) : -0.0000
Root Mean Squared Error (RMSE) : 0.8979
Mean Absolute Error (MAE) : 0.7379
None
adjusted r2 : 0.3168205826874245
AIC : 1105.8545994699402
BIC : 1162.1149044728875

```

	Predicted	Actual	Residual
329	1.943735	1	-0.943735
247	2.324615	1	-1.324615
390	1.982694	3	1.017306
145	1.438725	0	-1.438725
497	0.046313	0	-0.046313
513	1.763998	3	1.236002
165	1.140237	1	-0.140237
77	2.719729	2	-0.719729
534	1.427572	1	-0.427572
163	0.587365	0	-0.587365
271	1.631975	1	-0.631975
31	1.592320	3	1.407680
55	1.954580	1	-0.954580
90	2.021479	1	-1.021479
576	0.894450	0	-0.894450
76	2.735652	1	-1.735652
2	1.169686	1	-0.169686
256	1.819606	3	1.180394
311	1.847720	2	0.152280
334	2.059235	3	0.940765

Regression statistics

```

Mean Error (ME) : 0.0222
Root Mean Squared Error (RMSE) : 1.0030
Mean Absolute Error (MAE) : 0.8292

```

## Comparing updated Model wrt to New model after dropping Multicollinear columns

After comparing both the models it appears that model's efficiency has not reduced drastically and we could see that it faired similar. So it makes sense to drop the columns with higher VIF >5 and focus on only remaining columns.

The final list of variables is -

```
['z_studytime', 'z_failures', 'z_famrel', 'z_health', 'z_absences', 'schoolsup_yes', 'school_MS', 'sex_M', 'Mjob_teacher', 'Fjob_health', 'Fjob_services', 'reason_reputation']
```

Among all the variables it could be noted that for some variables we as a society cannot do much and its comes with a package, however we can try to focus on the variables which could influence the behaviour and outcome.

For example, we notice that student health, student absences, and study time are key contributing factors which we could influence as well.

Analysis of Key Factors which could be influenced and affecting Student performance

### **1. Health**

The health of students is a crucial aspect that significantly impacts their academic performance. Ensuring good health is not only vital from an educational perspective but also contributes to the overall well-being of students. Education departments and school authorities can play a pivotal role by focusing on initiatives such as health insurance coverage and regular health check-ups. This proactive approach can address health-related challenges that students might face at home.

### **2. Student Absences:**

Managing and reducing student absences is essential for maintaining a consistent learning environment. Communication with parents and implementing strict attendance policies may positively influence attendance rates. A deeper analysis could explore additional factors contributing to high absences, such as alcohol consumption, late-night activities on weekdays, and involvement in extracurricular activities. However, for the current analysis, the focus will be limited to understanding the impact of absences on academic grades.

### **3. Study Time:**

The amount of time students dedicate to studying is influenced by various factors, including the role of authorities, parents, and the broader societal context. Conducting a comprehensive study to identify reasons behind variations in study time—whether low or high—can provide insights into effective strategies. Understanding and addressing these factors can contribute to the development of initiatives that promote optimal study habits among students.

## **kNN for Predicting**

With the multiple regression model we were able to gain decent insights and could predict the grade ranges with decent accuracy. However, we can also try to explore the other modeling techniques to establish and identify if we could establish any new model which can give us any better results.



We could use kNN and see if it could which neighbors share the grades and what traights are common and what we could do influence them as an exercise.

If this model gives us more insights and better results we can observe and take actions based on it.

As a base exercise what we did is we ran a grid search on the column list we figured in multiple linear regression model above and see if it could fit better here?

After grid search it results revealed that its is giving a prediction accuracy of 34% which is slightly above the Naive model accuracy of 28.59% is not a significant gain.

So we will not invest much in it and move to another modeling technique of classification Trees

```
In [31]: 1 # Naive Model accuracy Calculation
          2 naive_df = df_norm['G3_bin'].value_counts()
          3 print('Naive Model\'s accuracy is : ',round((naive_df.max()/naive_df
Naive Model's accuracy is : 28.59 %
```

```
In [32]: 1 all_cols = ['z_age', 'z_Medu', 'z_Fedu', 'z_traveltime', 'z_studytime',
          2             'z_freetime', 'z_goout', 'z_Dalc', 'z_Walc', 'z_health',
          3             'famsup_yes', 'paid_yes', 'activities_yes', 'nursery_yes',
          4             'sex_M', 'address_U', 'famsize_LE3', 'Pstatus_T', 'Mjob_
          5             'Mjob_teacher', 'Fjob_health', 'Fjob_other', 'Fjob_serv
          6             'reason_other', 'reason_reputation', 'guardian_mother',
          7             print(all_cols)

['z_age', 'z_Medu', 'z_Fedu', 'z_traveltime', 'z_studytime', 'z_failure
s', 'z_famrel', 'z_freetime', 'z_goout', 'z_Dalc', 'z_Walc', 'z_health',
h', 'z_absences', 'schoolsup_yes', 'famsup_yes', 'paid_yes', 'activities
s_yes', 'nursery_yes', 'higher_yes', 'internet_yes', 'romantic_yes', 's
chool_MS', 'sex_M', 'address_U', 'famsize_LE3', 'Pstatus_T', 'Mjob_heal
th', 'Mjob_other', 'Mjob_services', 'Mjob_teacher', 'Fjob_health', 'Fjo
b_other', 'Fjob_services', 'Fjob_teacher', 'reason_home', 'reason_oth
er', 'reason_reputation', 'guardian_mother', 'guardian_other']
```

## Knn with all columns

```
In [33]: 1 # Knn with all columns
          2 # Identify the best kNN and see its performance
          3
          4 knn_classifier = KNeighborsClassifier()
          5 # Define the hyperparameter grid for grid search
          6 param_grid = {
          7     'n_neighbors': list(range(1,15)), # Adjust the values based on
          8 }
          9
         10 # Perform grid search with cross-validation
         11 grid_search = GridSearchCV(knn_classifier, param_grid, scoring='accu
         12 grid_search.fit(train_X[all_cols], train_y)
         13
         14 # Get the best parameters and model
         15 best_params = grid_search.best_params_
         16 best_model = grid_search.best_estimator_
         17
         18 # Make predictions on the test set
         19 pred_y = best_model.predict(valid_X[all_cols])
         20
         21 # Evaluate the best model
         22 accuracy = accuracy_score(valid_y, pred_y)
         23 print(f'Best Model - Accuracy: {accuracy:.4f}')
         24 print('Best Model - Best Parameters:', best_params)
         25
         26 result = pd.DataFrame({'Predicted': pred_y, 'Actual': valid_y,
         27                       'Residual': valid_y - pred_y})
         28 print(result.head())
         29
```

Best Model - Accuracy: 0.3604

Best Model - Best Parameters: {'n\_neighbors': 2}

	Predicted	Actual	Residual
329	1	1	0
247	2	1	-1
390	2	3	1
145	2	0	-2
497	1	0	-1

## KNN with only columns from final agreed multiple regression analysis

```
In [34]: 1 # Identify the best kNN and see its performance
2
3 knn_classifier = KNeighborsClassifier()
4 # Define the hyperparameter grid for grid search
5 param_grid = {
6     'n_neighbors': list(range(1,15)), # Adjust the values based on
7 }
8
9 # Perform grid search with cross-validation
10 grid_search = GridSearchCV(knn_classifier, param_grid, scoring='accu
11 grid_search.fit(train_X[col_list], train_y)
12
13 # Get the best parameters and model
14 best_params = grid_search.best_params_
15 best_model = grid_search.best_estimator_
16
17 # Make predictions on the test set
18 pred_y = best_model.predict(valid_X[col_list])
19
20 # Evaluate the best model
21 accuracy = accuracy_score(valid_y, pred_y)
22 print(f'Best Model - Accuracy: {accuracy:.4f}')
23 print('Best Model - Best Parameters:', best_params)
24
25 result = pd.DataFrame({'Predicted': pred_y, 'Actual': valid_y,
26                        'Residual': valid_y - pred_y})
27 print(result.head())
```

Best Model - Accuracy: 0.3423

Best Model - Best Parameters: {'n\_neighbors': 9}

	Predicted	Actual	Residual
329	3	1	-2
247	3	1	-2
390	2	3	1
145	0	0	0
497	0	0	0

## Classification Trees

In this method we gonna train the mdeol for the columns we identified and notice if we could improve the performance of the model using a new technique. This model is good in identifying student traights and see if we could do anything if the accuracy is good.

We ran a grid search on the model using param grid and noticed it attained a total accuracy of 36.94% which is pretty much the same for the kNN and hence we wil stop here and will not divulge more into its analysis further an move on to next technique as Random Forrest.

```
In [35]: 1 # Start with an all parameters
2 param_grid = {
3     'max_depth': [3,5,10, 20, 30, 40],
4     'min_samples_split': [10, 15, 20, 40, 60, 80, 100,150],
5     'min_impurity_decrease': [0, 0.0005, 0.001, 0.005, 0.01],
6 }
7 gridSearch = GridSearchCV(DecisionTreeClassifier(), param_grid, cv=5
8     #that the available computer memory (CPU) will be used to ma
9 gridSearch.fit(train_X[all_cols], train_y)
10 print('Initial best score: ', gridSearch.best_score_)
11 print('Initial Best parameters: ', gridSearch.best_params_)
```

Initial best score: 0.43076697032030564

Initial Best parameters: {'max\_depth': 10, 'min\_impurity\_decrease': 0, 'min\_samples\_split': 80}

```
In [36]: 1 # Start with an initial guess for parameters
2 param_grid = {
3     'max_depth': [3,5,10, 20, 30, 40],
4     'min_samples_split': [10, 15, 20, 40, 60, 80, 100,150],
5     'min_impurity_decrease': [0, 0.0005, 0.001, 0.005, 0.01],
6 }
7 gridSearch = GridSearchCV(DecisionTreeClassifier(), param_grid, cv=5
8     #that the available computer memory (CPU) will be used to ma
9 gridSearch.fit(train_X[col_list], train_y)
10 print('Initial best score: ', gridSearch.best_score_)
11 print('Initial Best parameters: ', gridSearch.best_params_)
```

Initial best score: 0.396620628856891

Initial Best parameters: {'max\_depth': 3, 'min\_impurity\_decrease': 0, 'min\_samples\_split': 100}

```
In [37]: 1 # Adapt grid based on result from initial grid search
2 param_grid = {
3     'max_depth': list(range(2, 16)),
4     'min_samples_split': list(range(70, 90)),
5     'min_impurity_decrease': [0.0001,0.0009, 0.001, 0.0011],
6 }
7 gridSearch = GridSearchCV(DecisionTreeClassifier(), param_grid, cv=5
8 gridSearch.fit(train_X, train_y)
9 print('Improved score: ', gridSearch.best_score_)
10 print('Improved parameters: ', gridSearch.best_params_)
11
12 bestClassTree = gridSearch.best_estimator_
```

Improved score: 0.43564501910079345

Improved parameters: {'max\_depth': 8, 'min\_impurity\_decrease': 0.0001, 'min\_samples\_split': 82}

In [38]: 1 classificationSummary(valid\_y, bestClassTree.predict(valid\_X))

Confusion Matrix (Accuracy 0.3694)

	Prediction				
Actual	0	1	2	3	
0	44	8	4	3	
1	23	20	10	12	
2	12	14	6	17	
3	17	4	16	12	

## Random Forrest

In this methdo we ran a Random forrest on all the attributes and tried to indetify most important attributes and then run the model on the high important attributes to see if we have any improvement.

In [39]: 1 rf = RandomForestClassifier(n\_estimators=500, random\_state=1)  
2 rf.fit(train\_X, train\_y)

Out[39]:

▼

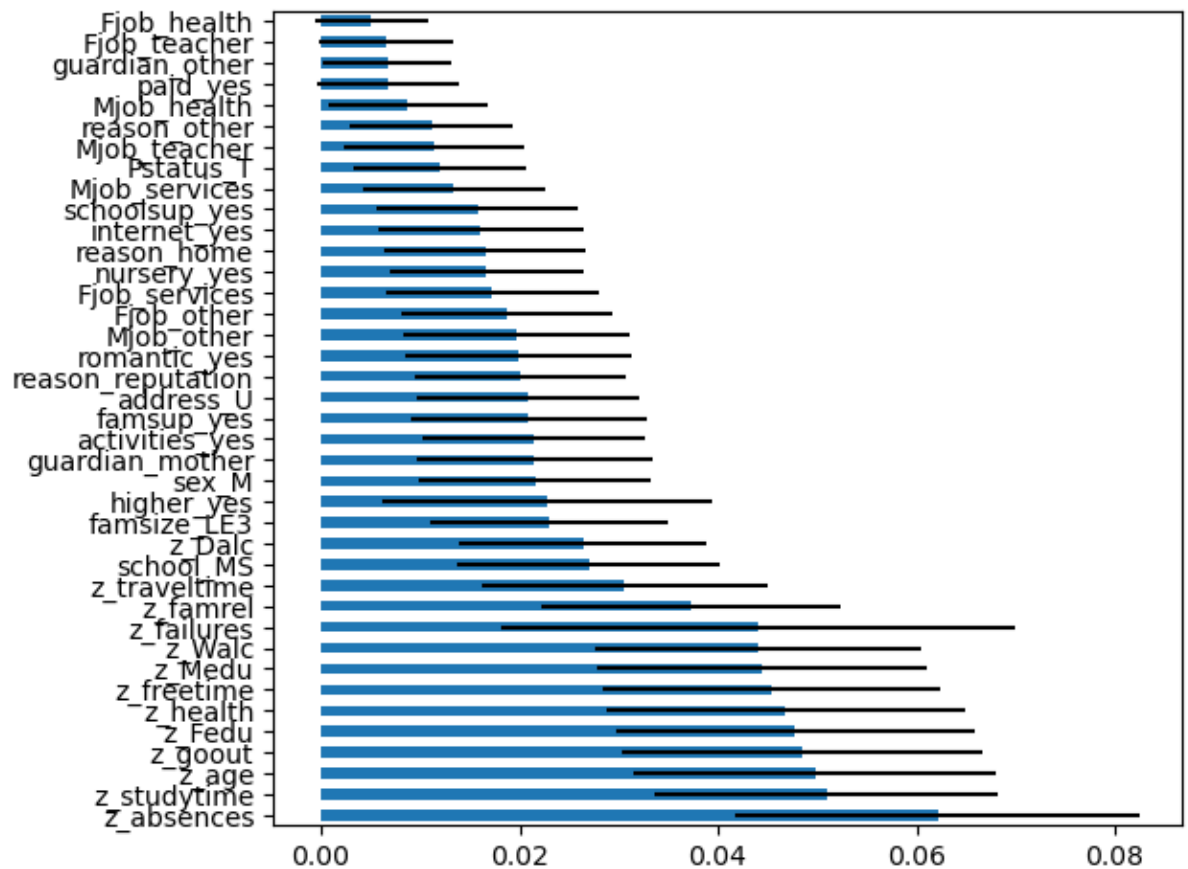
RandomForestClassifier

RandomForestClassifier(n\_estimators=500, random\_state=1)

```
In [40]: 1 importances = rf.feature_importances_
2 std = np.std([tree.feature_importances_ for tree in rf.estimators_],
3
4 df = pd.DataFrame({'feature': train_X.columns, 'importance': importances})
5 df = df.sort_values('importance', ascending=False)
6 print(df)
7 plt.figure(figsize=(15, 25))
8 ax = df.plot(kind='barh', xerr='std', x='feature', legend=False)
9 ax.set_ylabel('')
10
11 plt.tight_layout()
12 plt.show()
```

	feature	importance	std
12	z_absences	0.062206	0.020379
4	z_studytime	0.050927	0.017259
0	z_age	0.049788	0.018285
8	z_goout	0.048520	0.018143
2	z_Fedu	0.047774	0.018038
11	z_health	0.046831	0.018110
7	z_freetime	0.045443	0.016959
1	z_Medu	0.044439	0.016645
10	z_Walc	0.044092	0.016456
5	z_failures	0.044087	0.025917
6	z_famrel	0.037268	0.015042
3	z_traveltime	0.030566	0.014392
21	school_MS	0.027013	0.013215
9	z_Dalc	0.026381	0.012438
24	famsize_LE3	0.022955	0.011957
18	higher_yes	0.022817	0.016692
22	sex_M	0.021547	0.011707
37	guardian_mother	0.021497	0.011840
16	activities_yes	0.021388	0.011200
14	famsup_yes	0.020897	0.011855
23	address_U	0.020872	0.011169
36	reason_reputation	0.020041	0.010651
20	romantic_yes	0.019860	0.011369
27	Mjob_other	0.019719	0.011458
31	Fjob_other	0.018676	0.010640
32	Fjob_services	0.017291	0.010652
17	nursery_yes	0.016647	0.009774
34	reason_home	0.016536	0.010190
19	internet_yes	0.016130	0.010308
13	schoolsup_yes	0.015813	0.010172
28	Mjob_services	0.013369	0.009169
25	Pstatus_T	0.012052	0.008706
29	Mjob_teacher	0.011371	0.009147
35	reason_other	0.011143	0.008174
26	Mjob_health	0.008759	0.008031
15	paid_yes	0.006855	0.007151
38	guardian_other	0.006715	0.006498
33	Fjob_teacher	0.006618	0.006769
30	Fjob_health	0.005099	0.005746

<Figure size 1500x2500 with 0 Axes>



**Variables with high importance factors are -**

1. z\_absences
2. z\_studytime
3. z\_age
4. z\_goout
5. z\_Fedu
6. z\_health
7. z\_freetime
8. z\_Medu
9. z\_Walc
10. z\_failures
11. z\_famrel

We will do the boost with the above attributes and notice if there is any change?

```
In [41]: 1 rand_frst_col_list = ['z_absences', 'z_studytime', 'z_age', 'z_goout', 'z_
2          'z_freetime', 'z_Medu', 'z_Walc', 'z_failures', 'z_
3 boost = GradientBoostingClassifier()
4 boost.fit(train_X[rand_frst_col_list], train_y)
5 classificationSummary(valid_y, boost.predict(valid_X[rand_frst_col_l
```

Confusion Matrix (Accuracy 0.3649)

		Prediction			
Actual	0	1	2	3	
	0	28	16	7	8
	1	21	20	14	10
	2	9	15	12	13
	3	6	10	12	21

## lets run multiple regression on the above columns as well to notice if we have better results?

It seems the results are kind of tied up and hence we gonna stick with the old model only. However, it is debatable from the fact that the new model has more variables which could be influenced by external factors to influence the Grades. So a rational perspective is to include the variables and stick with a slight reduction in accuracy and more controlling power on the influencing factors.

So the updated list of variables is

'z\_absences', 'z\_studytime', 'z\_age', 'z\_goout', 'z\_Fedu', 'z\_health', 'z\_freetime', 'z\_Medu', 'z\_Walc', 'z\_

From the above list variables that we could influence externally are

1. Absence
2. Studytime
3. Goout
4. Health
5. Freetime
6. Weekend Alcohol Consumptions
7. Family Relations

These above variables may define the grades and students approach to the study. A small change or effort on the above list which changes the favors in better grade we could impact the grade and betterment of society as a whole.



```
In [42]: 1 # Multiple Regression after Random Forrest
2 col_list = rand_frst_col_list
3 rndm_grade_lm = LinearRegression()
4 rndm_grade_lm.fit(train_X[col_list], train_y)
5
6 # print coefficients
7 print('intercept ', rndm_grade_lm.intercept_)
8 print(pd.DataFrame({'Predictor': col_list, 'coefficient': rndm_grade_lm.coef_}))
9
10 # print performance measures
11 print(regressionSummary(train_y, rndm_grade_lm.predict(train_X[col_list])))
12
13 pred_y = rndm_grade_lm.predict(train_X[col_list])
14
15 print('adjusted r2 : ', adjusted_r2_score(train_y, pred_y, rndm_grade_lm))
16 print('AIC : ', AIC_score(train_y, pred_y, rndm_grade_lm))
17 print('BIC : ', BIC_score(train_y, pred_y, rndm_grade_lm))
18
19 # Validation Data
20
21 # Use predict() to make predictions on a new set
22 rndm_grade_lm_pred = rndm_grade_lm.predict(valid_X[col_list])
23
24 result = pd.DataFrame({'Predicted': rndm_grade_lm_pred, 'Actual': valid_y,
25                        'Residual': valid_y - rndm_grade_lm_pred})
26 print(result.head(20))
27
28 # Compute common accuracy measures
29 print('Random Forrest Generated Variable Model Summary')
30 regressionSummary(valid_y, rndm_grade_lm_pred)
31
32 print('\nVIF Fixed Model Summary ')
33 regressionSummary(valid_y, vif_grade_lm_pred)
34
```

```

intercept 1.3264315334451535
Predictor coefficient
0 z_absences -0.163113
1 z_studytime 0.230278
2 z_age 0.095586
3 z_goout -0.045332
4 z_Fedu 0.089760
5 z_health -0.077227
6 z_freetime -0.021675
7 z_Medu 0.150442
8 z_Walc -0.019077
9 z_failures -0.329948
10 z_famrel 0.052846

```

#### Regression statistics

```

Mean Error (ME) : 0.0000
Root Mean Squared Error (RMSE) : 0.9540
Mean Absolute Error (MAE) : 0.7936

```

None

```
adjusted r2 : 0.23066781438127193
```

```
AIC : 1153.698519088712
```

```
BIC : 1205.940230877163
```

	Predicted	Actual	Residual
329	1.297840	1	-0.297840
247	1.990478	1	-0.990478
390	1.775865	3	1.224135
145	1.231638	0	-1.231638
497	0.205498	0	-0.205498
513	1.953847	3	1.046153
165	1.250343	1	-0.250343
77	2.066557	2	-0.066557
534	1.586544	1	-0.586544
163	0.137119	0	-0.137119
271	1.644406	1	-0.644406
31	1.752808	3	1.247192
55	1.420302	1	-0.420302
90	1.839336	1	-0.839336
576	0.912990	0	-0.912990
76	2.108759	1	-1.108759
2	1.007666	1	-0.007666
256	1.647397	3	1.352603
311	1.354146	2	0.645854
334	1.742476	3	1.257524

#### Random Forrest Generated Variable Model Summary

#### Regression statistics

```

Mean Error (ME) : 0.0872
Root Mean Squared Error (RMSE) : 0.9563
Mean Absolute Error (MAE) : 0.7865

```

#### VIF Fixed Model Summary

#### Regression statistics

```
Mean Error (ME) : 0.0222
```

Root Mean Squared Error (RMSE) : 1.0030  
Mean Absolute Error (MAE) : 0.8292

## Summary

Yes the mean error has increased and RMSE has decreased we will stick to the new model as it gives more opportunity to course correct the grades and influence the behaviour.

## Let's Try Clustering as well !!

We can also try to cluster the data and see if we could get any more insights into the data and have the better results.

In [ ]:

1