

Title: Factor affecting student performance

Group 1

Group Members:

1. Lim Chin Woon [Leader] (s2124307)
2. Muhammad Amirul Daniel bin Badrul Hisham (s2115750)
3. Nurul Shahirah binti Sha'ari (s2120876)
4. Qiyi Liu (s2106632)
5. Yeoh Li Tian (s2120306)

Overview and Motivation

With the popularity of the Internet, sensors and various digital terminal devices, an interconnected world is taking shape. At the same time, with the explosive exponential growth of data, data analysis has penetrated into all aspects of our daily life and pushed us to an era of profound change. The education field has no exception. Through the mining and analysis of a large number of students' education process and organization management data collected by modern schools, we can get the characteristics with strong data causality with students' final scores, and then make a reliable early prediction of students' future performance. The prediction of students' academic performance not only helps the school to timely intervene in the curriculum, enrich the classroom content and improve the training plan, but also helps to analyze students' academic performance, understand students' learning trends, encourage excellent students to achieve better results, and encourage students with relatively weak foundation to establish confidence and make continuous progress.

Related Work

1. Relationship of Parental Support on Healthy Habits, School Motivations and Academic Performance in Adolescents

<https://doi.org/10.3390/ijerph17030882>

-The objective of the study was to analyze how parental support relates to the physical activity practice, satisfaction with sports, level of physical activity, academic performance and alcohol consumption. Adolescents with little parental support show ($p < 0.001$) more boredom, less fun, worse academic performance and higher alcohol consumption. Gender shows differences ($p < 0.001$) experiencing girls more boredom, less fun, less PA practice and higher academic performance than boys. Age establishes ($p < 0.01$) that older adolescents (15–16 years old) experience more boredom, less fun, less PA practice, lower academic performance and higher alcohol consumption than young boys and girls (12–14 years old).

2. The Influence of Parents Educational Level on Secondary School Students Academic Achievements in District Rajanpur.

<https://files.eric.ed.gov/fulltext/EJ1079955.pdf>

-Impact of parents educational level on students academic achievement at secondary level of education. After analysis of the data the research finds significant positive relationship between parents education level and academic achievements of students.

3. Gender Gaps in Student Academic Achievement and Inequality

https://www.michaelsmith.cz/wp-content/uploads/2018/02/Gender_Gaps_Tsai_Smith_Hause_r_RSE.pdf

-PISA surveys, suggest that boys do better in math and science, whereas girls do better in reading but in overall, gender differences in family and school influences do not account for gender differences in academic achievement in any of the six countries.

4. Student Performance Prediction by Using Data Mining Classification Algorithms

https://www.researchgate.net/profile/Dorina-Kabakchieva/publication/272178031_Student_Performance_Prediction_by_Using_Data_Mining_Classification_Algorithms/links/5a151ffda6fdc6d697bc01ea/Student-Performance-Prediction-by-Using-Data-Mining-Classification-Algorithms.pdf

-Predicting student performance using 4 algorithms which are a Rule Learner, a Decision tree, a Neural network, and K-Nearest Neighbour

5. Predicting Student Performance in a Portuguese Secondary Institution

https://www.leonshpaner.com/projects/post/student_performance_models/Student_Performance_Models.pdf

-Study on grades in on Portugal High school. Uses 6 different machine learning model and compare 8 different metric scores

Research Question

We managed to obtain a relevant dataset here: <https://www.kaggle.com/ishandutta/student-performance-data-set>.

Based on the dataset, we will be answering these specific questions:

1. Studying how the variables can influence individual student performance in Portugal language study.
2. Predict student performance according to the variables after training the model using different machine learning algorithm.
3. Evaluate model accuracy and performance.

Data Cleaning/Data Preprocessing

Load library

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3 import pandas as pd
4 import seaborn as sns
```

Import and read dataset

The initial dataset is separated with semicolon(;). The data is converted into table form.

```

1 dataseturl = 'https://drive.google.com/uc?id=18gjZIQ060KZqxbk-
  BD2MpaEKZqJqUWwG'
2 pd.set_option('display.max_columns', None)
3 rawdf = pd.read_csv(dataseturl, sep=';')
4 print(rawdf.head())

```

```

1  school sex  age address famsize Pstatus  Medu  Fedu  Mjob  Fjob \
2  0      GP  F   18      U    GT3      A    4    4  at_home teacher
3  1      GP  F   17      U    GT3      T    1    1  at_home  other
4  2      GP  F   15      U    LE3      T    1    1  at_home  other
5  3      GP  F   15      U    GT3      T    4    2  health  services
6  4      GP  F   16      U    GT3      T    3    3   other   other
7
8  reason guardian  traveltime  studytime  failures  schoolsup  famsup paid \
9  0  course  mother      2      2.0      0      yes    no    no
10 1  course  father      1      2.0      0      no    yes    no
11 2  other  mother      1      2.0      0      yes    no    no
12 3  home  mother      1      3.0      0      no    yes    no
13 4  home  father      1      2.0      0      no    yes    no
14
15  activities nursery higher internet romantic  famrel  freetime  goout  Dalc
16 \
17 0      no    yes    yes      no      no      4      3      4      1
18 1      no    no    yes    yes      no      5      3      3      1
19 2      no    yes    yes    yes      no      4      3      2      2
20 3      yes    yes    yes    yes    yes      3      2      2      1
21 4      no    yes    yes      no      no      4      3      2      1
22
23 walc  health  absences  G1  G2  G3
24 0      1      3      4    0  11  11
25 1      1      3      2    9  11  11
26 2      3      3      6   12  13  12
27 3      1      5      0   14  14  14
28 4      2      5      0   11  13  13

```

Attributes for the variables in the dataset.

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

10. Fjob - father's job (nominal: "teacher", "health" care related, civil "services" (e.g. administrative or police), "at_home" or "other")
11. reason - reason to choose this school (nominal: close to "home", school "reputation", "course" preference or "other")
12. guardian - student's guardian (nominal: "mother", "father" or "other")
13. traveltime - home to school travel time (numeric: 1 - <15 min., 2 - 15 to 30 min., 3 - 30 min. to 1 hour, or 4 - >1 hour)
14. studytime - weekly study time (numeric: 1 - <2 hours, 2 - 2 to 5 hours, 3 - 5 to 10 hours, or 4 - >10 hours)
15. failures - number of past class failures (numeric: n if 1<=n<3, else 4)
16. schoolsup - extra educational support (binary: yes or no)
17. famsup - family educational support (binary: yes or no)
18. paid - extra paid classes within the course subject (Math or Portuguese) (binary: yes or no)
19. activities - extra-curricular activities (binary: yes or no)
20. nursery - attended nursery school (binary: yes or no)
21. higher - wants to take higher education (binary: yes or no)
22. internet - Internet access at home (binary: yes or no)
23. romantic - with a romantic relationship (binary: yes or no)
24. famrel - quality of family relationships (numeric: from 1 - very bad to 5 - excellent)
25. freetime - free time after school (numeric: from 1 - very low to 5 - very high)
26. goout - going out with friends (numeric: from 1 - very low to 5 - very high)
27. Dalc - workday alcohol consumption (numeric: from 1 - very low to 5 - very high)
28. Walc - weekend alcohol consumption (numeric: from 1 - very low to 5 - very high)
29. health - current health status (numeric: from 1 - very bad to 5 - very good)
30. absences - number of school absences (numeric: from 0 to 93)

These grades are related with the Portuguese subject:

1. G1 - first period grade (numeric: from 0 to 20)
2. G2 - second period grade (numeric: from 0 to 20)
3. G3 - final grade (numeric: from 0 to 20, output target)

Checking any non-availability (NA) inside the dataset

```
1 | rawdf.isna().sum()
```

```
1 | school      8
2 | sex         3
3 | age         0
4 | address     0
5 | famsize     0
6 | Pstatus     0
7 | Medu        0
8 | Fedu        0
9 | Mjob        0
10 | Fjob        0
11 | reason      0
12 | guardian    0
13 | traveltime  0
```

```

14 | studytime      2
15 | failures       0
16 | schoolsup      0
17 | famsup         0
18 | paid           0
19 | activities     0
20 | nursery        0
21 | higher         0
22 | internet       0
23 | romantic       0
24 | famrel         0
25 | freetime       0
26 | goout          0
27 | dalc           0
28 | walc           0
29 | health         0
30 | absences       0
31 | G1             0
32 | G2             0
33 | G3             0
34 | dtype: int64

```

We can see there are variables that contain NA value. We will handle accordingly.

1. student's school (school)
2. student's sex (sex)
3. weekly study time(studytime)

Imputation

We know that school and sex are in nominal pattern, hence We will handle using imputation by mode.

```

1 | #Impute missing value in school and sex using mode.
2 | cols = ['school', 'sex']
3 |
4 | for col in cols:
5 |     rawdf[col] = rawdf[col].fillna(rawdf[col].mode()[0])
6 |     print(rawdf['school'].isna().sum())

```

```

1 | 0
2 | 0

```

Weekly study time is a numeric data where the dataset has categorized into a range of 1, 2, 3, 4. In this case, we will handle using imputation by median.

```

1 | rawdf['studytime'] = rawdf['studytime'].fillna(rawdf['studytime'].median())
2 | rawdf['studytime'].isna().sum()

```

```

1 | 0

```

Check values pattern consistency

```
1 for nom_var in rawdf.columns[rawdf.dtypes.eq('object')]:
2     print(rawdf[nom_var].value_counts().sort_index())
3     print('\n')
```

```
1 GP      423
2 MS      226
3 Name: school, dtype: int64
```

```
1 F       383
2 M       266
3 Name: sex, dtype: int64
```

```
1 R       197
2 U       452
3 Name: address, dtype: int64
```

```
1 GT3     457
2 LE3     192
3 Name: famsize, dtype: int64
```

```
1 A        80
2 T       569
3 Name: Pstatus, dtype: int64
```

```
1 at_home    135
2 health     48
3 other      258
4 services   136
5 teacher    72
6 Name: Mjob, dtype: int64
```

```
1 | at_home      42
2 | health       23
3 | other        367
4 | services     181
5 | teacher      36
6 | Name: Fjob, dtype: int64
```

```
1 | course       285
2 | home         149
3 | other        72
4 | reputation   143
5 | Name: reason, dtype: int64
```

```
1 | father      153
2 | mother      455
3 | other       41
4 | Name: guardian, dtype: int64
```

```
1 | no      581
2 | yes     68
3 | Name: schoolsup, dtype: int64
```

```
1 | no      251
2 | yes     398
3 | Name: famsup, dtype: int64
```

```
1 | no      610
2 | yes     39
3 | Name: paid, dtype: int64
```

```
1 | no      334
2 | yes     315
3 | Name: activities, dtype: int64
```

```
1 | no      128
2 | yes     521
3 | Name: nursery, dtype: int64
```

```
1 no      69
2 yes     580
3 Name: higher, dtype: int64
```

```
1 no      151
2 yes     498
3 Name: internet, dtype: int64
```

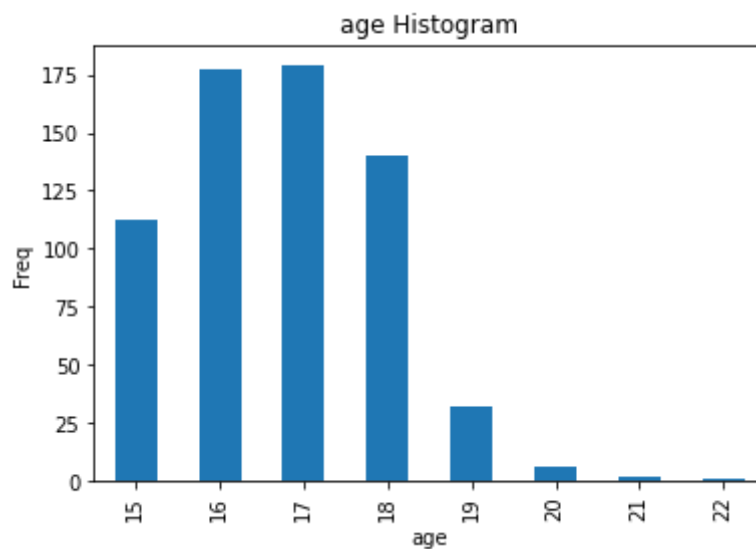
```
1 no      410
2 yes     239
3 Name: romantic, dtype: int64
```

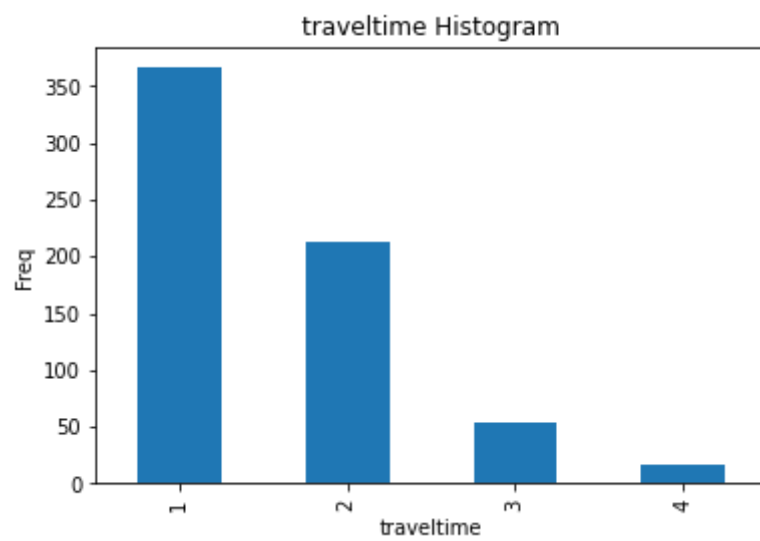
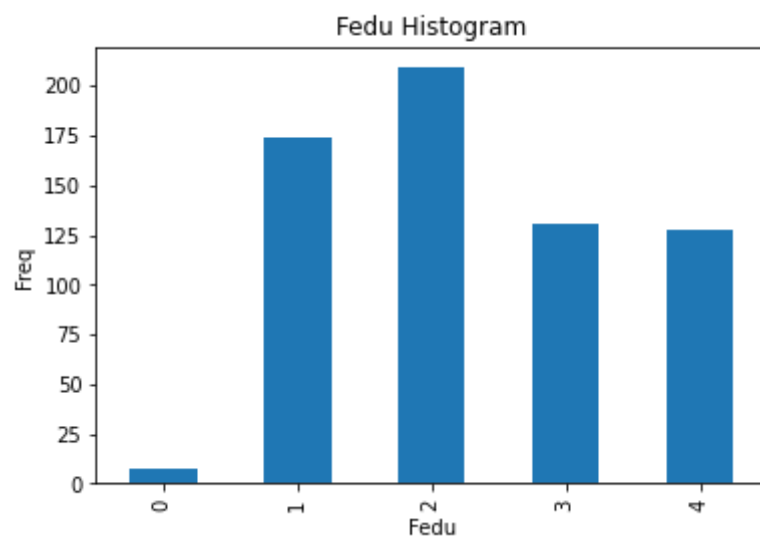
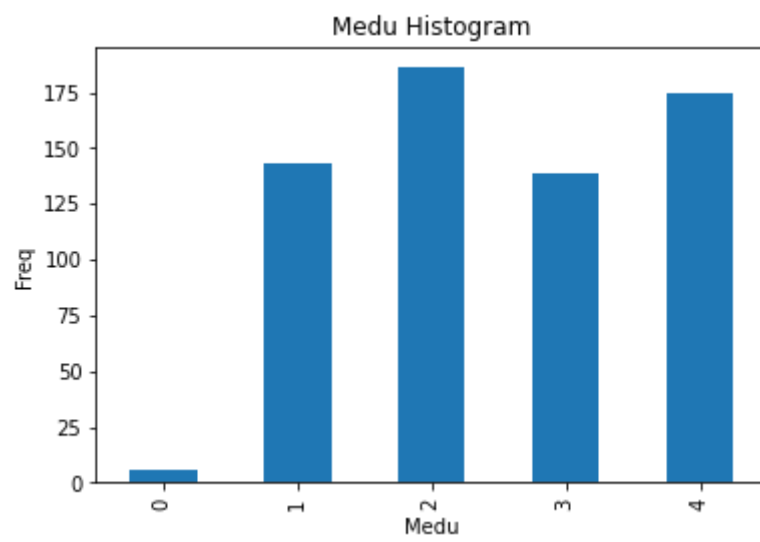
We can see that there are no inconsistent naming of values.

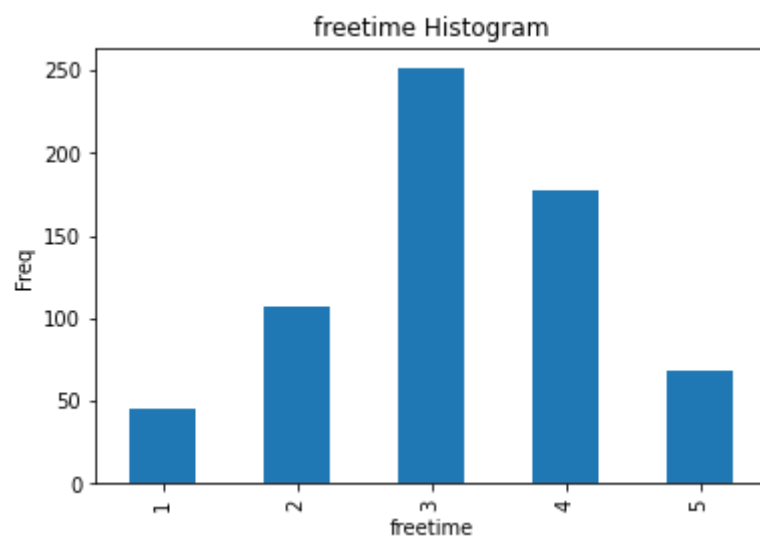
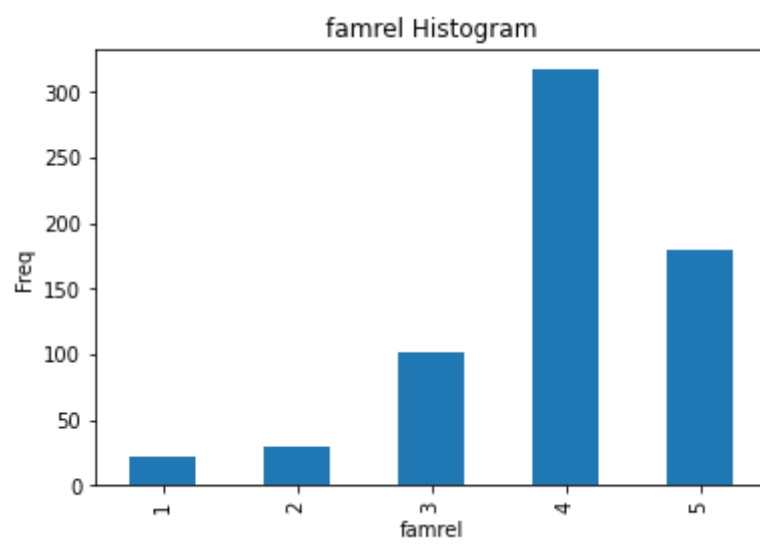
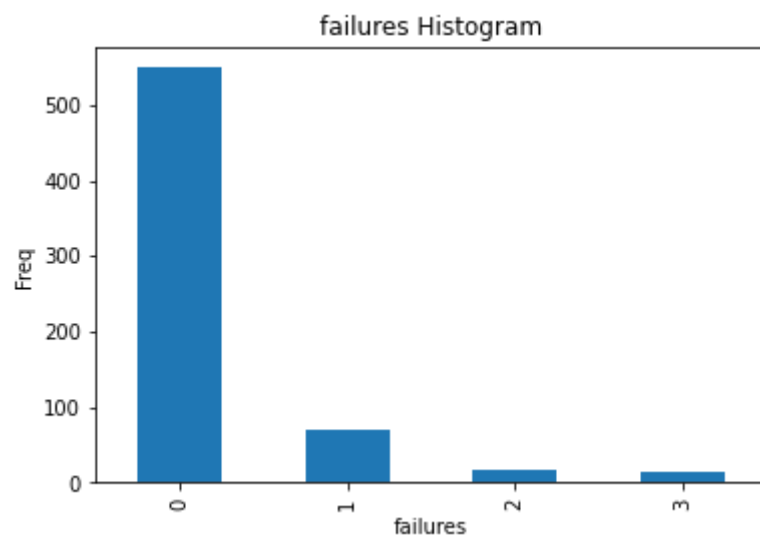
Visualising distribution of numerical variables:

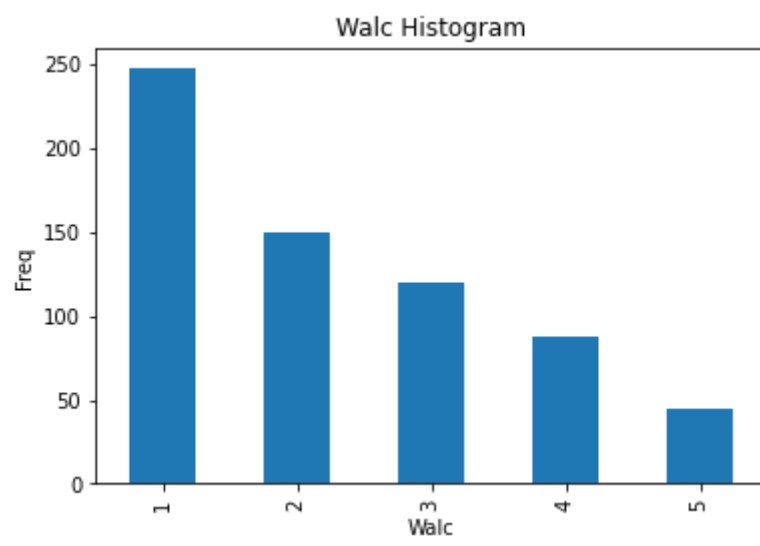
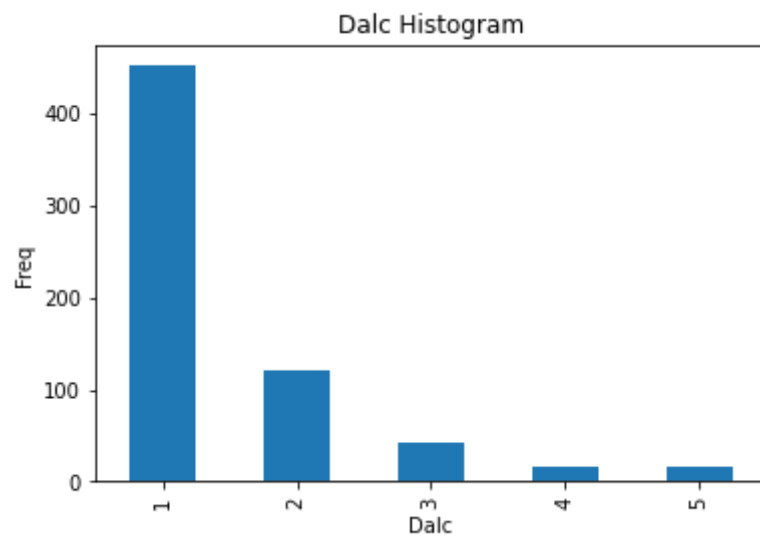
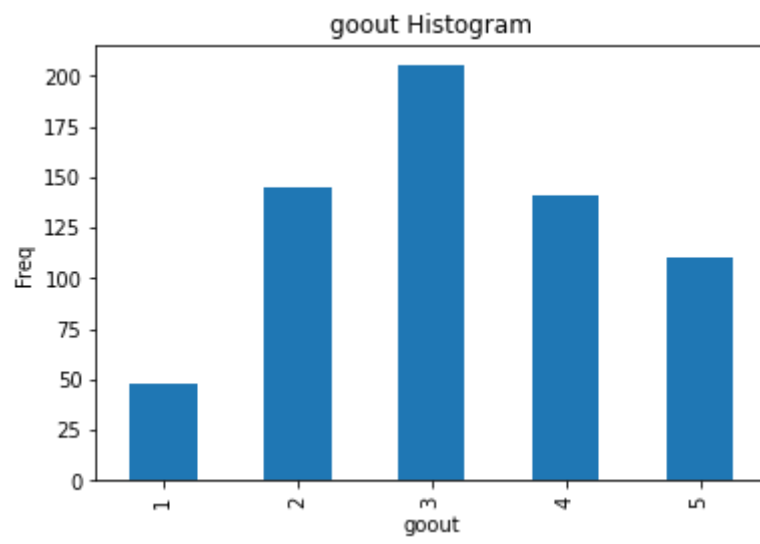
We can use histogram to visualise the distributions of numerical variables:

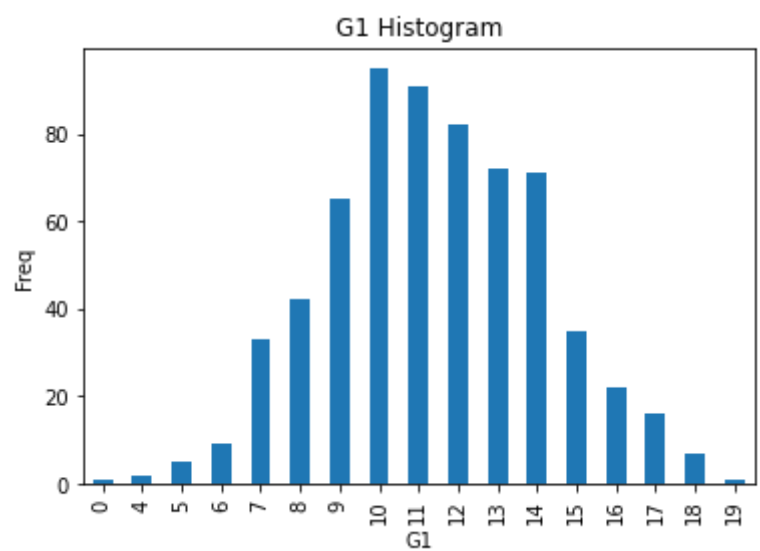
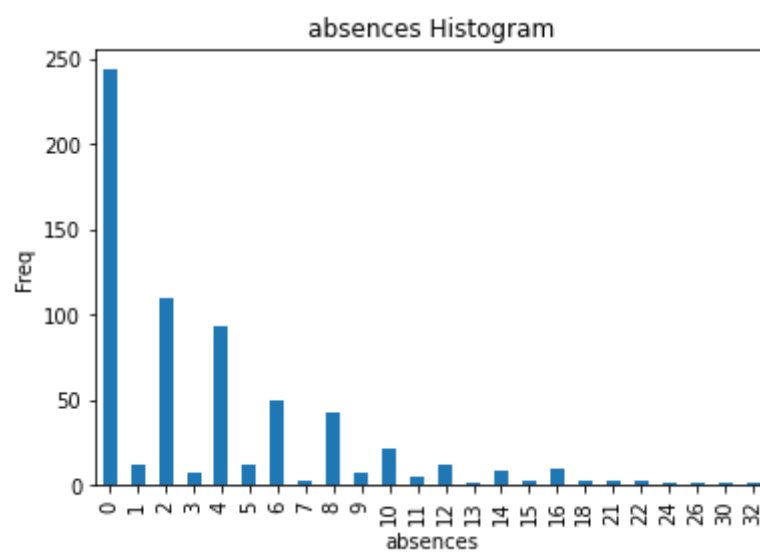
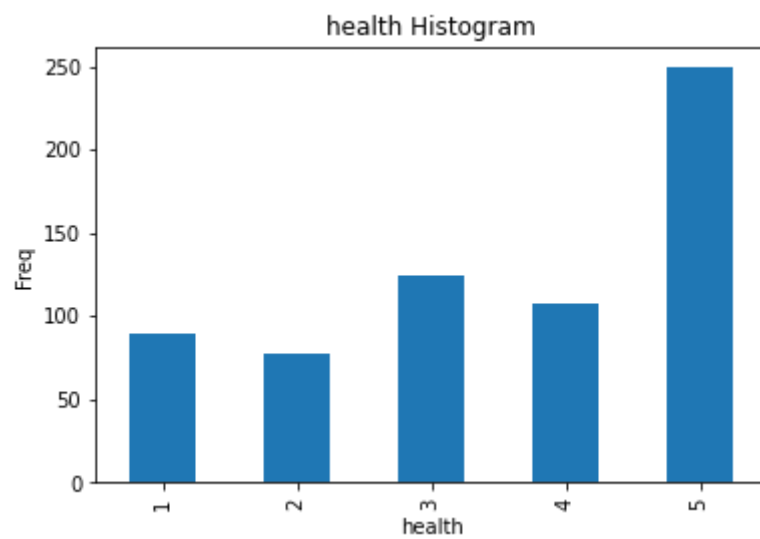
```
1 for numvar in rawdf.columns[rawdf.dtypes.eq('int64')]:
2     plt.title(numvar + " Histogram")
3     plt.xlabel(numvar)
4     plt.ylabel('Freq')
5     rawdf[numvar].value_counts().sort_index().plot(kind="bar")
6     plt.show()
7
```

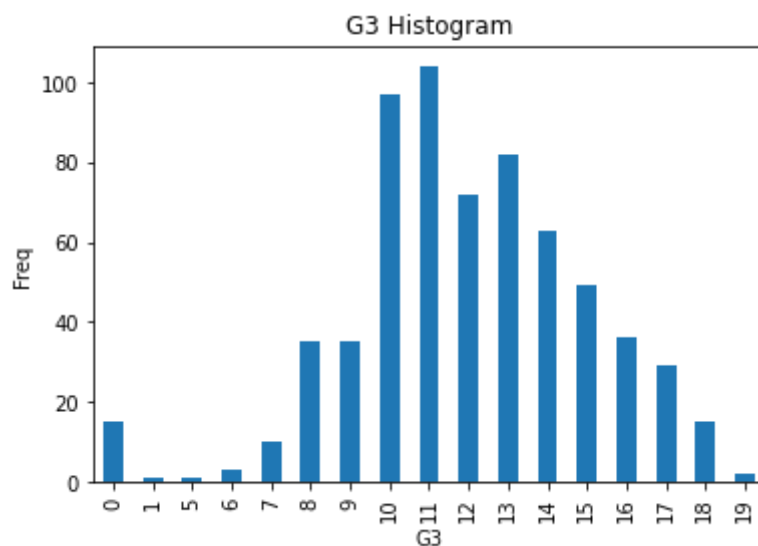
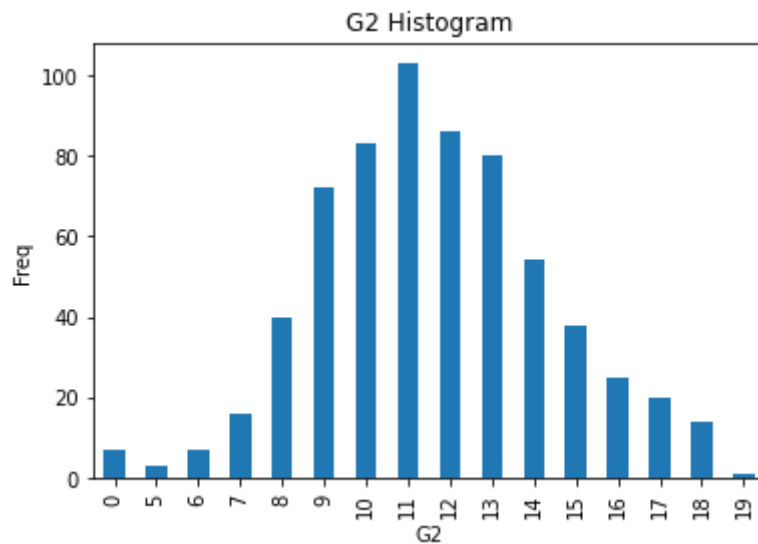












We see that there are very small number of students with Medu=0 or Fedu=0. On the other hand, the number of students with G3=0 is unexpectedly high.

Feature selection

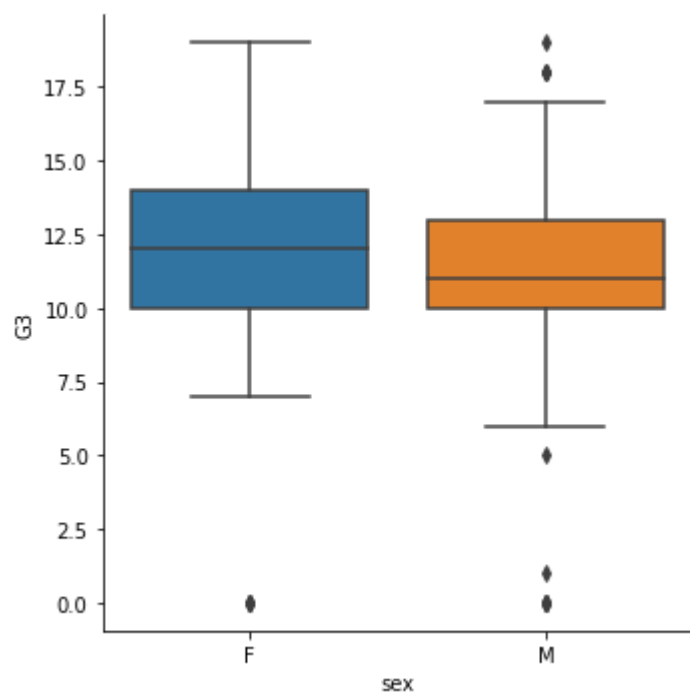
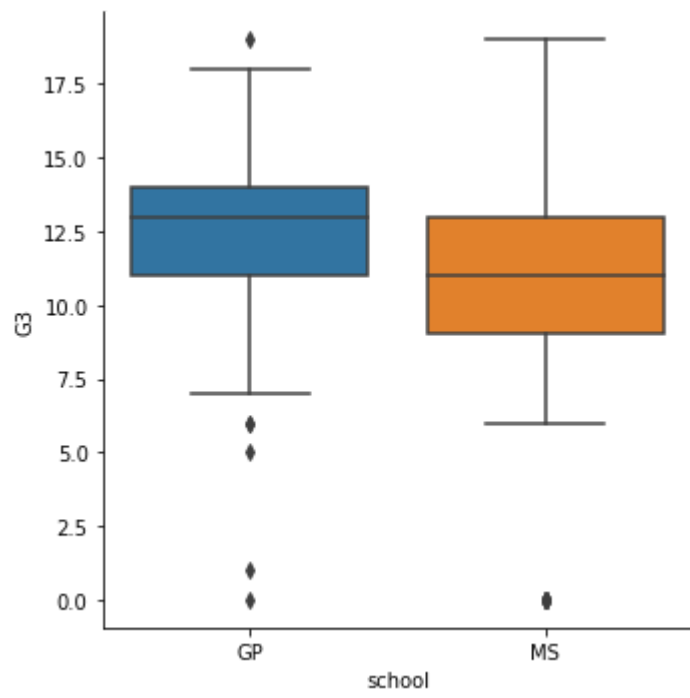
For categorical variables:

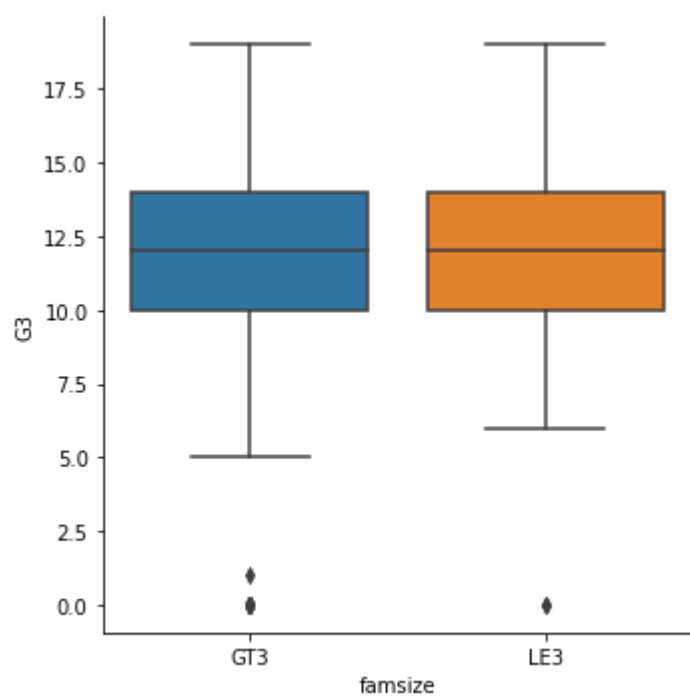
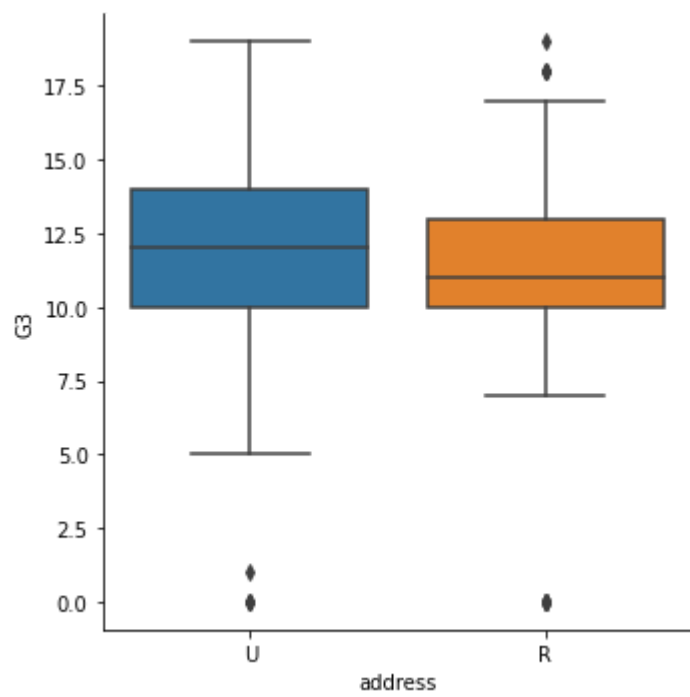
We use a boxplot to visualise the distribution of target G3 for different group of students

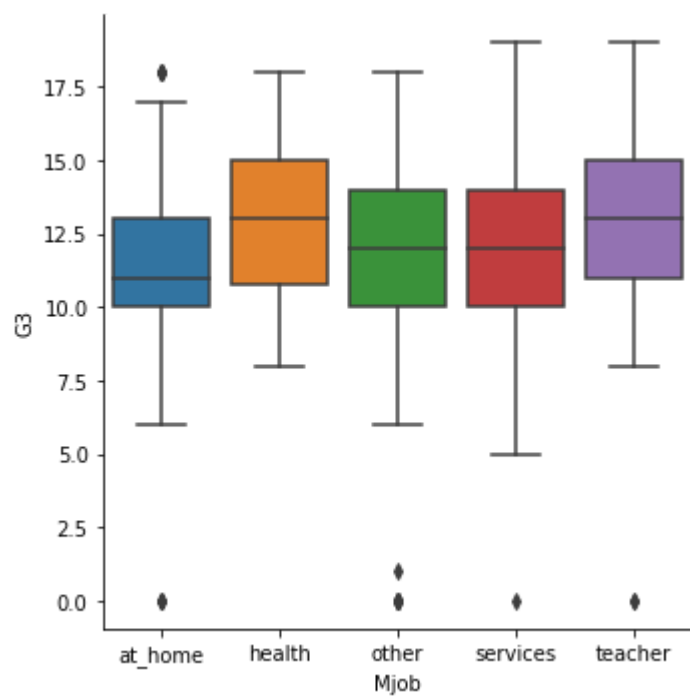
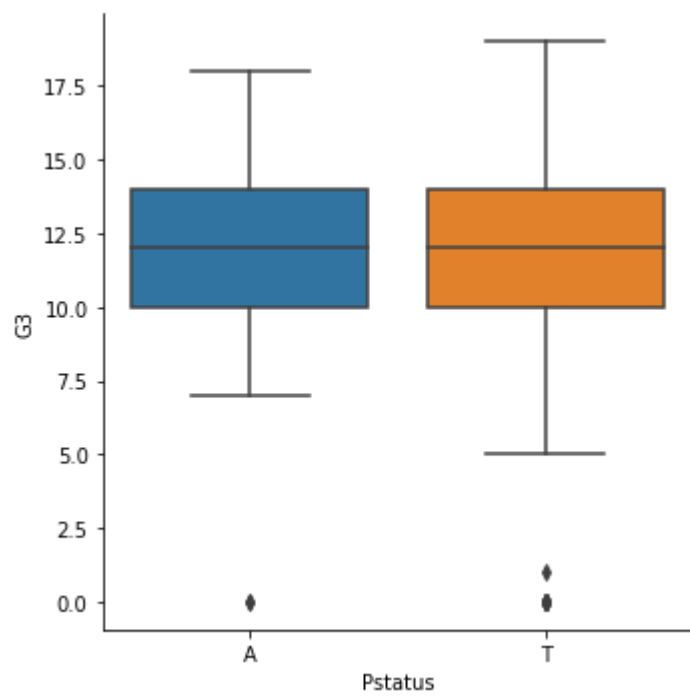
```

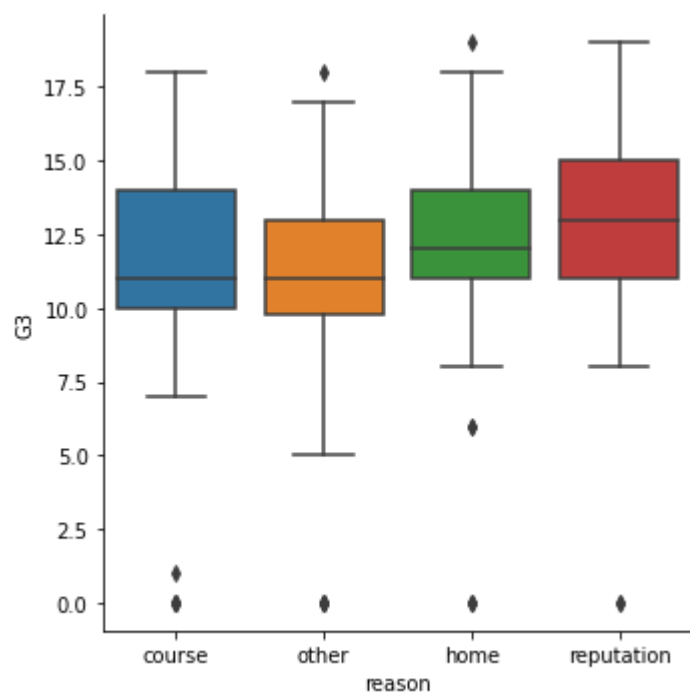
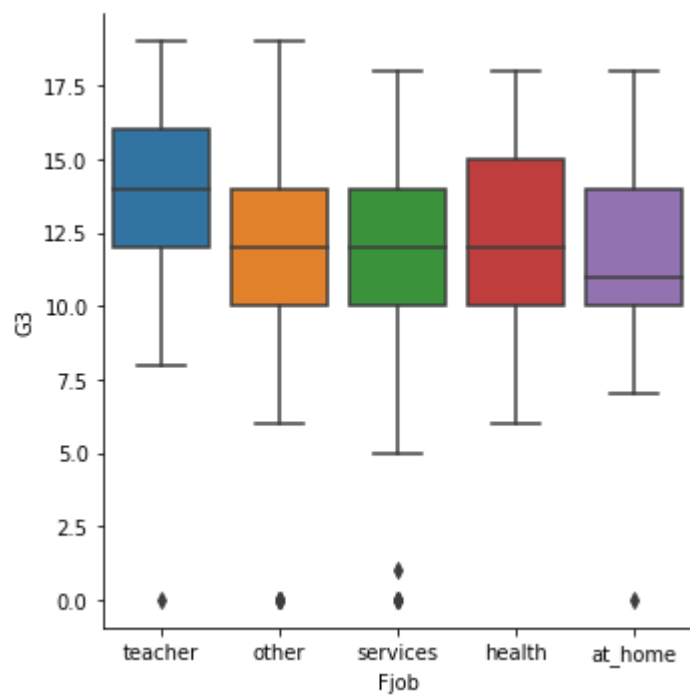
1 #filter out categorical variables
2 catvar = rawdf.columns[rawdf.dtypes.eq('object')]
3
4 # plot box plot
5 for cvar in catvar:
6     sns.catplot(x=cvar, y="G3", kind="box", data=rawdf)
7     plt.show()

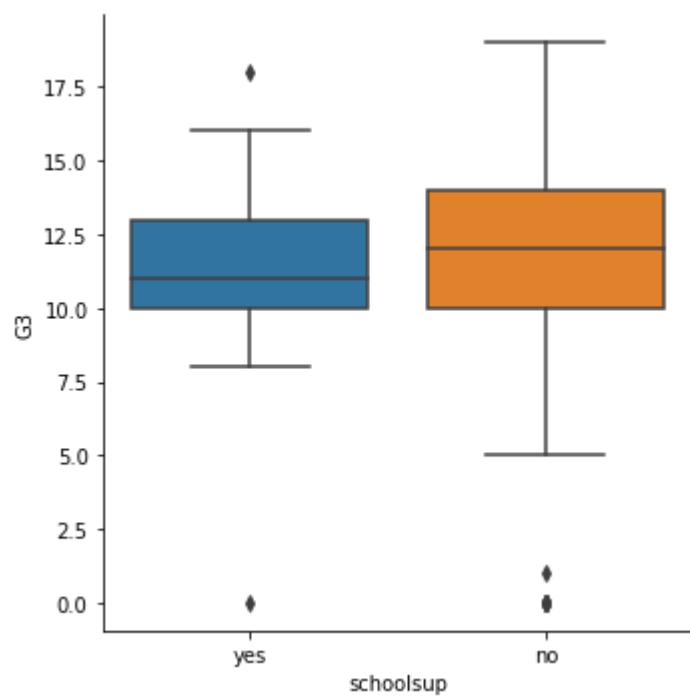
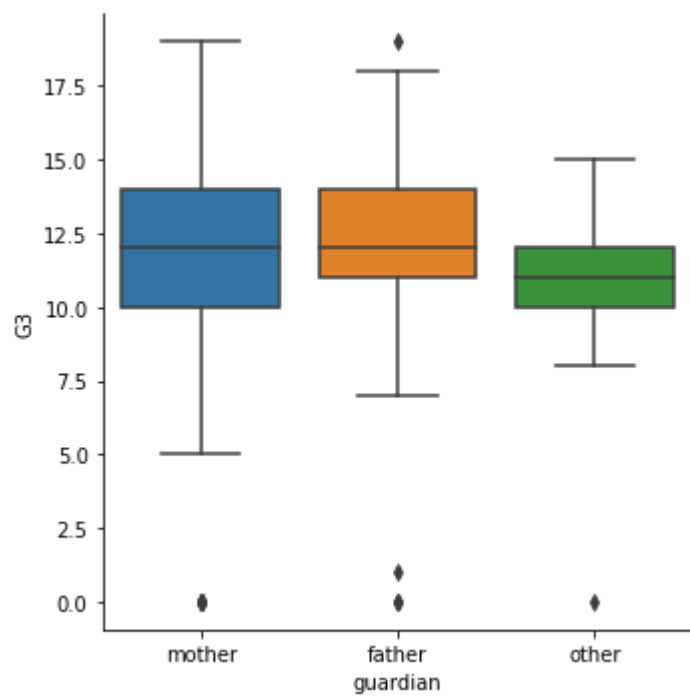
```

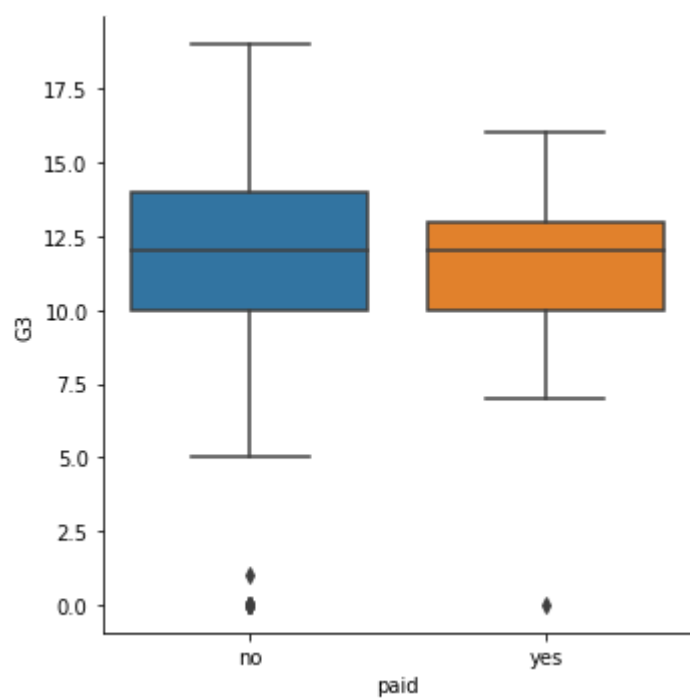
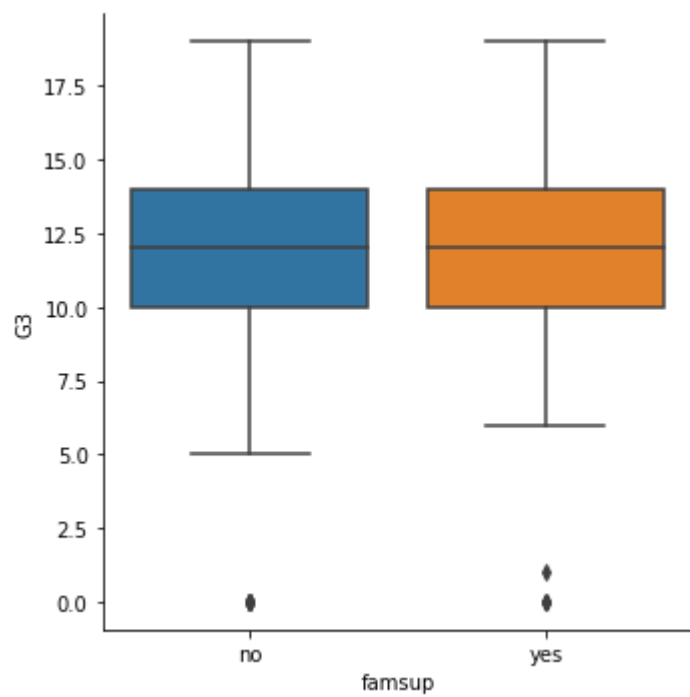


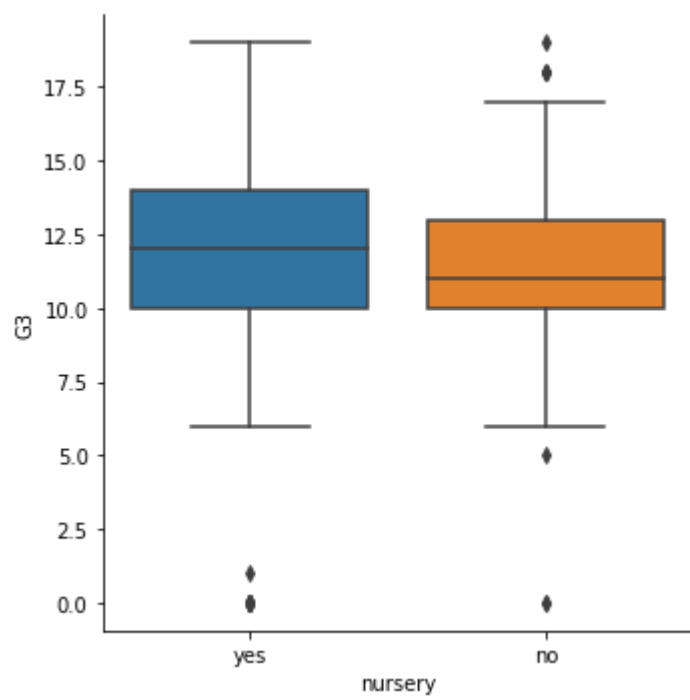
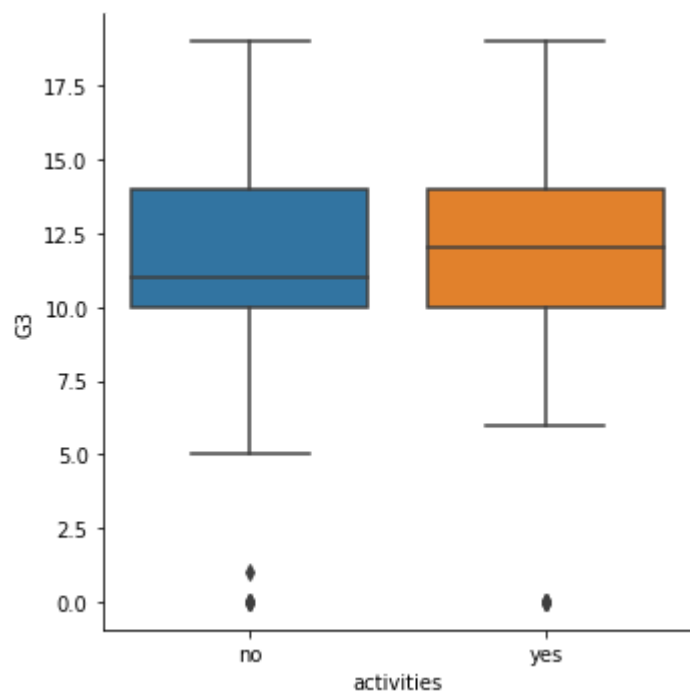


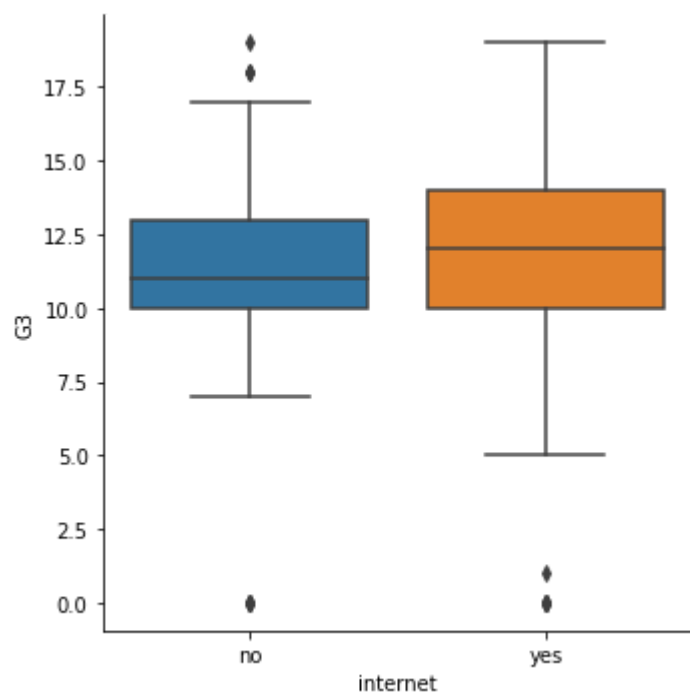
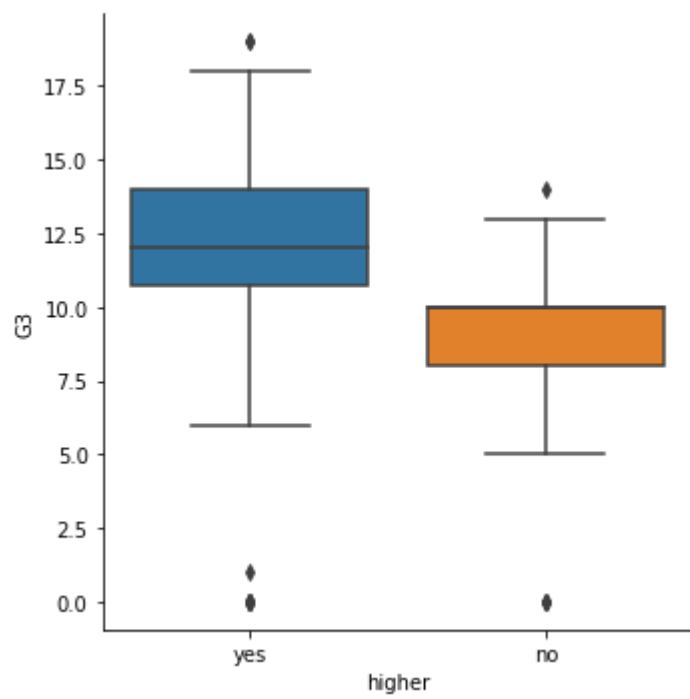


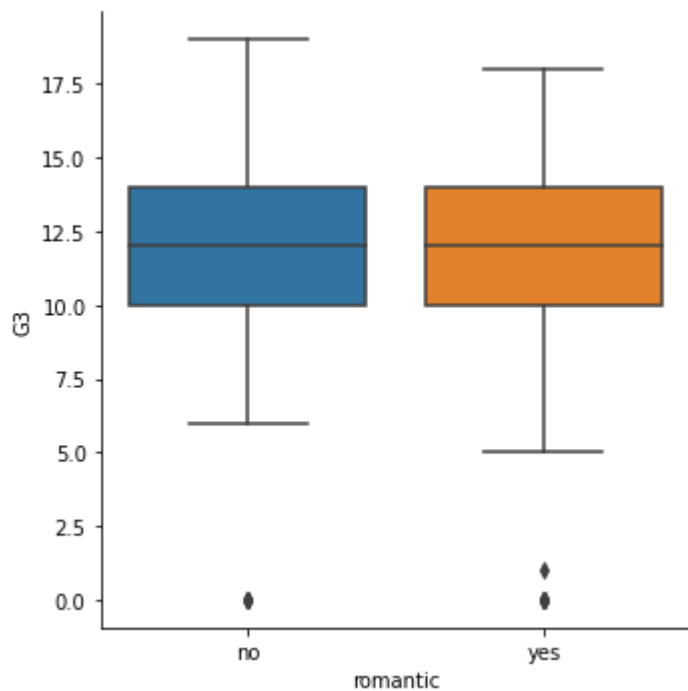












From the figures, we can see that there is not much variation across different groups in the following variables:

1. family size (famsize),
2. parent's cohabitation status (Pstatus),
3. family educational support (famsup),
4. extra paid classes within the course subject (paid),
5. with a romantic relationship (romantic).

Hence we will consider to drop these variables.

For numerical variables:

We use a heatmap and scatterplot to investigate on the correlation of each variables with G3

```
1 plt.figure(figsize=(1, 16))
2 corrcol = rawdf.corr()[['G3']]
3 sns.heatmap(corrcol.sort_values(by='G3', ascending=False), annot=True,
  cbar=False)
```

```
1 <matplotlib.axes._subplots.AxesSubplot at 0x7f46ea6a2050>
```



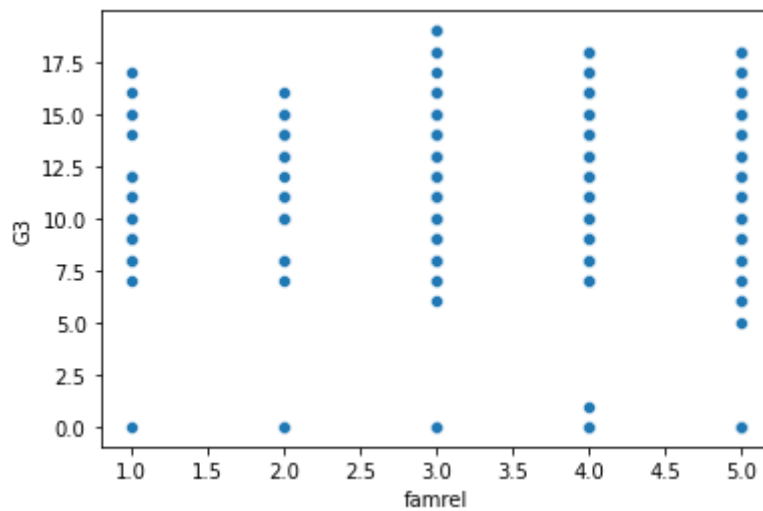
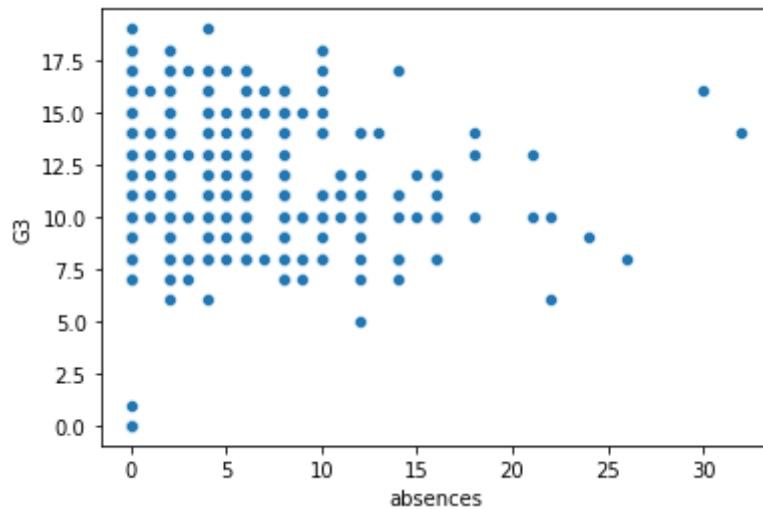
By observing the heatmap, we can identify those variables which show a very weak correlation ($|r| < 0.1$) to our output:

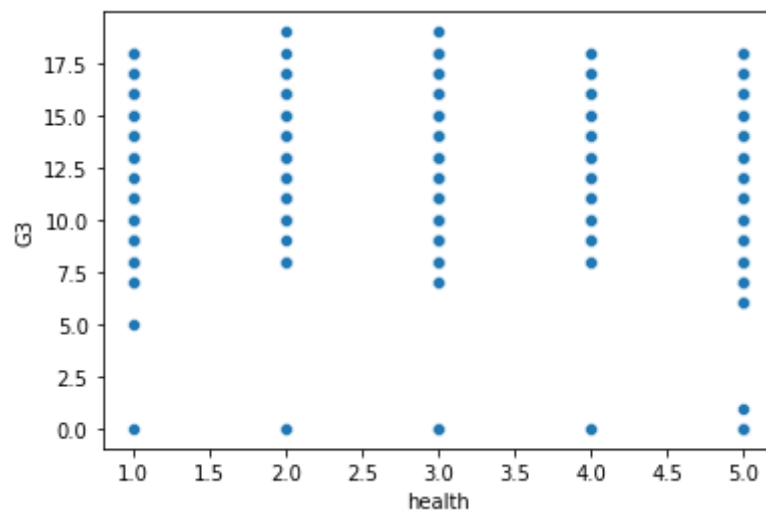
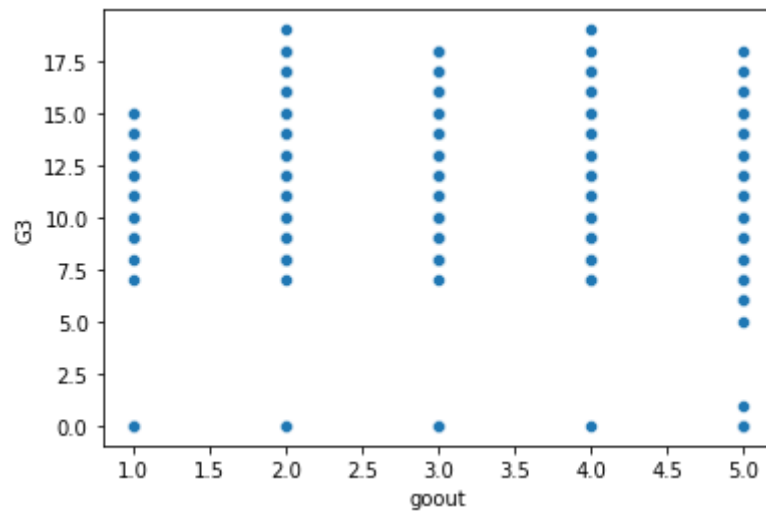
1. number of school absences (absences),
2. quality of family relationships (famrel),

3. going out with friends (goout),
4. current health status (health).

We will use scatterplot to identify if there are any abnormal values that affect the correlation.

```
1 #Plot G3 versus absences/famrel/goout/health
2 vars = ['absences', 'famrel', 'goout', 'health']
3
4 for var in vars:
5     sns.scatterplot(x=rawdf[var], y = rawdf.G3)
6     plt.show()
```





We can see that there is some extreme condition in absences variable. But it is still worth to keep as it has a negative correlation pattern against G3.

For the remaining variables, data is evenly distributed. Hence, we will keep the absences variable for further analysis and consider to drop the other three.

Chi2 Test

To support our previous finding, we do a Chi-Squared Test to find features relevancy score to our target.

```

1  from sklearn.feature_selection import SelectKBest, chi2
2  from sklearn.preprocessing import LabelEncoder
3
4  chi2_df = rawdf.copy()
5  label_encoder = LabelEncoder()
6
7  for nom_var in chi2_df.columns[chi2_df.dtypes.eq('object')]:
8      chi2_df[nom_var] = label_encoder.fit_transform(chi2_df[nom_var])
9
10 x = chi2_df.drop('G3',axis=1)
11 y = chi2_df['G3']
12

```

```
13 # Select two features with highest chi-squared statistics
14 chi2_selector = SelectKBest(chi2, k='all')
15 chi2_selector.fit(x, y)
16
17 # Look at scores returned from the selector for each feature
18 chi2_scores = pd.DataFrame(list(zip(chi2_df.columns,
19                                     chi2_selector.scores_)), columns=['features', 'score'])
19 chi2_scores.sort_values(by='score', ascending=False)
```

```
1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }
```

	features	score
31	G2	418.452113
30	G1	333.615973
29	absences	294.362769
14	failures	255.779332
0	school	60.035797
26	Dalc	35.491003
8	Mjob	34.888847
10	reason	33.617500
6	Medu	33.567952
7	Fedu	32.368762
27	Walc	25.772026
15	schoolsup	21.756390
13	studytime	19.574536
28	health	14.518698
17	paid	13.252095
1	sex	12.930063
12	traveltime	12.919294
20	higher	11.343756
18	activities	11.131398
3	address	11.106601
22	romantic	10.741687
25	goout	10.105568
24	freetime	9.030872
4	famsize	8.130862
16	famsup	7.533753
9	Fjob	7.277847
11	guardian	6.306938
21	internet	5.694293
2	age	5.223356
19	nursery	3.187132

	features	score
23	famrel	3.039747
5	Pstatus	0.487380

```

1 cleandf = rawdf.drop(['famsize', 'Pstatus', 'famsup', 'paid', 'romantic',
2   'famrel', 'goout', 'health'], axis=1)
3 print(cleandf.columns)

```

```

1 Index(['school', 'sex', 'age', 'address', 'Medu', 'Fedu', 'Mjob', 'Fjob',
2   'reason', 'guardian', 'traveltime', 'studytime', 'failures',
3   'schoolsup', 'activities', 'nursery', 'higher', 'internet', 'freetime',
4   'Dalc', 'Walc', 'absences', 'G1', 'G2', 'G3'],
5   dtype='object')

```

Dropping of Variables

We will be dropping variables based on the findings from the 3 previous sections. First we use boxplot, correlation and scatter plot to filter out a total of 8 candidates, then we make the final decision with the help of Chi-Squared Test. A critical value of 15 is set, where we find that all the 8 variables have chi2 value below it. Hence, they are all dropped:

1. family size (famsize),
2. parent's cohabitation status (Pstatus),
3. family educational support (famsup),
4. extra paid classes within the course subject (paid),
5. with a romantic relationship (romantic),
6. quality of family relationships (famrel),
7. going out with friends (goout),
8. current health status (health).

```

1 cleandf = rawdf.drop(['famsize', 'Pstatus', 'famsup', 'paid', 'romantic',
2   'famrel', 'goout', 'health'], axis=1)
3 print(cleandf.columns)

```

```

1 Index(['school', 'sex', 'age', 'address', 'Medu', 'Fedu', 'Mjob', 'Fjob',
2   'reason', 'guardian', 'traveltime', 'studytime', 'failures',
3   'schoolsup', 'activities', 'nursery', 'higher', 'internet', 'freetime',
4   'Dalc', 'Walc', 'absences', 'G1', 'G2', 'G3'],
5   dtype='object')

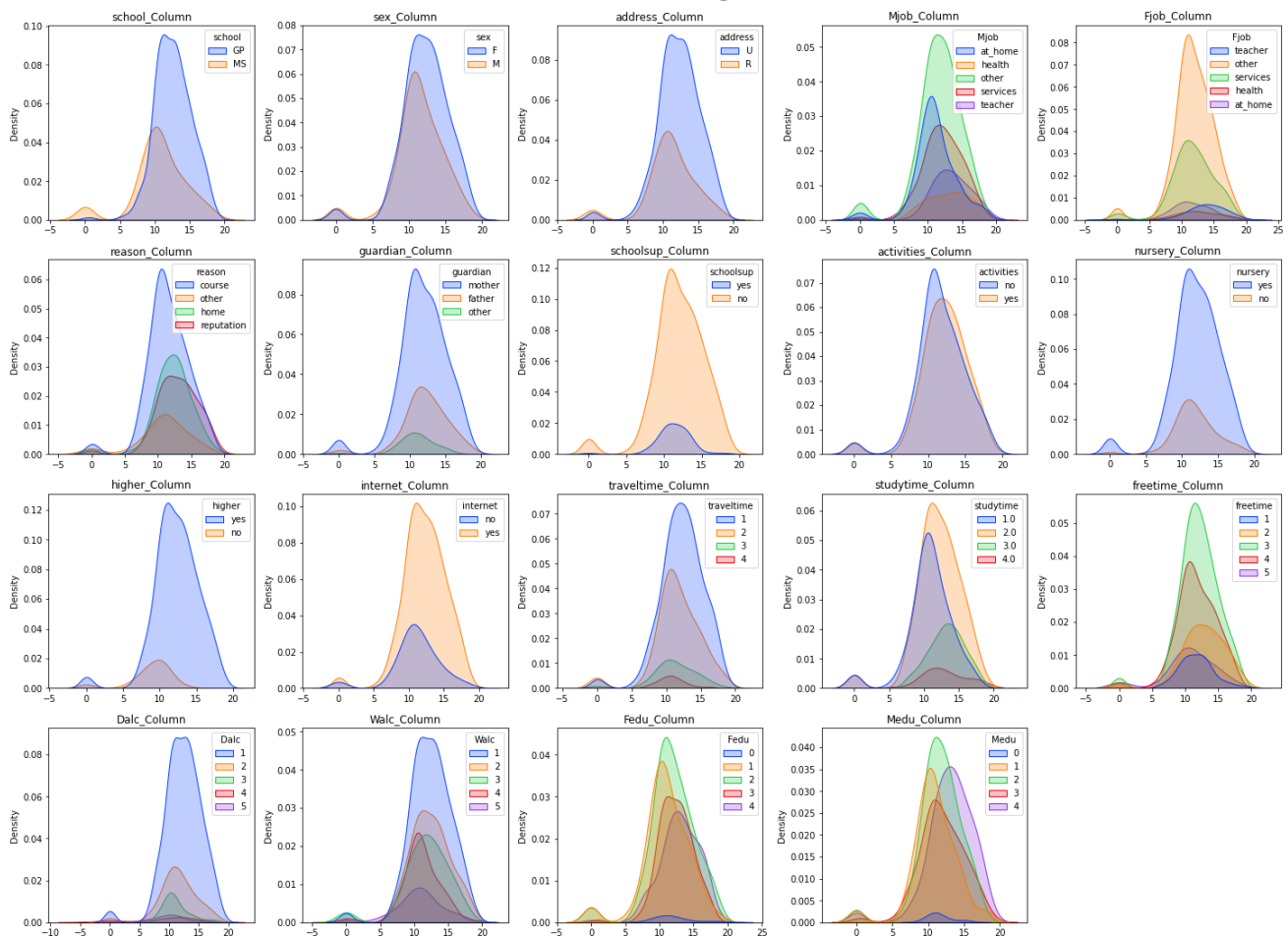
```

Exploratory Data Analysis

Categorical Feature

```
1 nominal_columns = ['school', 'sex', 'address', 'Mjob', 'Fjob', 'reason', 'guardian',
2                   'schoolsup', 'activities', 'nursery', 'higher', 'internet']
3 ordinal_columns = ['traveltime', 'studytime', 'freetime', 'Dalc', 'Walc', 'Fedu', 'Medu']
4
5 fig, axes = plt.subplots(4, 5, figsize=(20, 15), constrained_layout=True)
6 fig.suptitle('G3 Distribution for each categorical feature', size=20,
7             fontweight='bold', fontfamily='serif')
8 axes = axes.ravel()
9 axes[-1].remove() # removing non plotted axis
10
11 for i in range(len(nominal_columns + ordinal_columns)):
12     # Creating plotting data variables
13     ax = axes[i]
14     col = (nominal_columns + ordinal_columns)[i]
15     # plotting function
16     sns.kdeplot(x='G3', data=cleandf, hue=col, ax=ax, shade=True, palette='bright')
17     # Adjusting plot
18     ax.set_xlabel('')
19     ax.set_title(col + '_Column')
```

G3 Distribution for each categorical feature



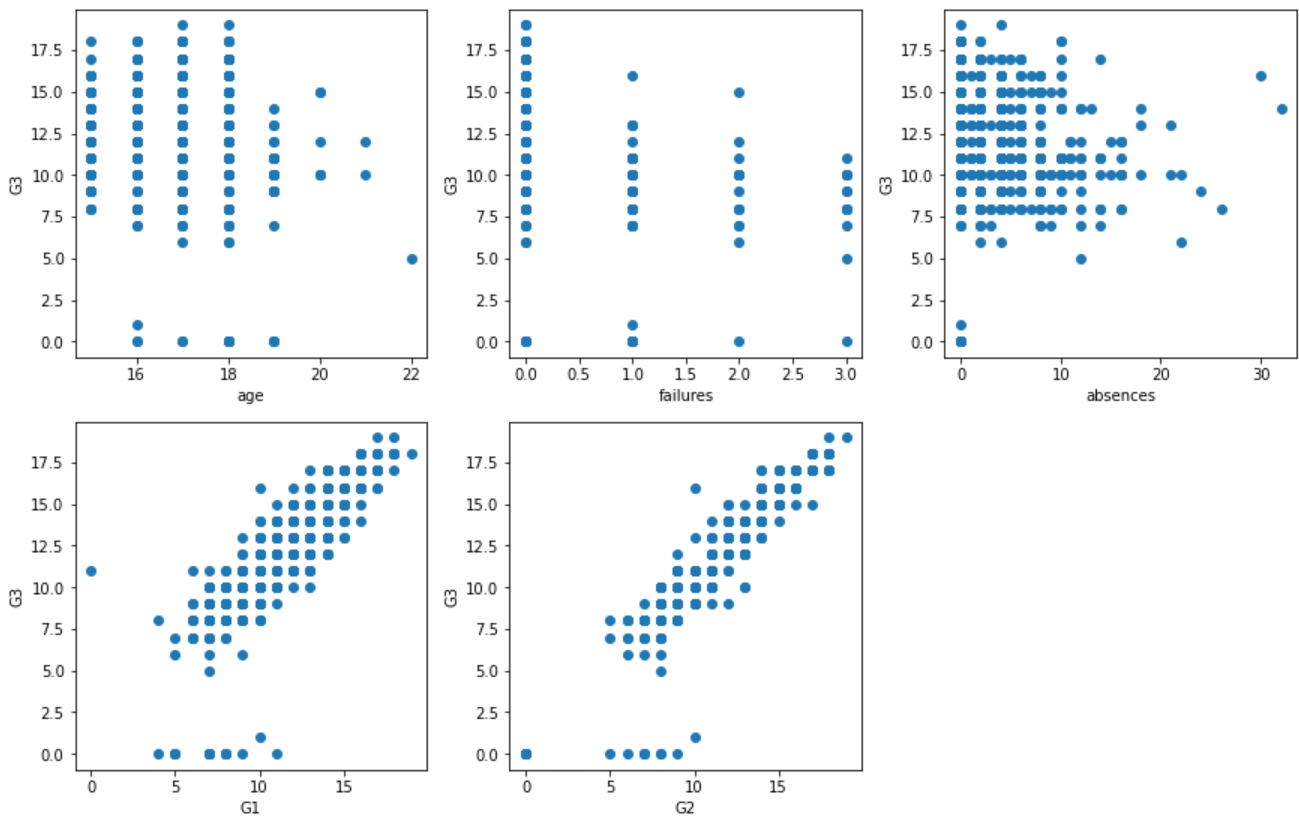
Continuous Feature

```

1  # 'age', 'failures', 'absences', 'G1', 'G2', 'G3'
2  import seaborn as sns
3
4  fig, axes = plt.subplots(2, 3, figsize=(12, 8), constrained_layout=True)
5  fig.suptitle('G3 Distribution for each continuous feature', size=20,
6               fontweight='bold', fontfamily='serif')
7
8  axes[0, 0].scatter(cleandf.age, cleandf.G3)
9  axes[0, 0].set_xlabel('age')
10 axes[0, 0].set_ylabel('G3')
11 axes[0, 1].scatter(cleandf.failures, cleandf.G3)
12 axes[0, 1].set_xlabel('failures')
13 axes[0, 1].set_ylabel('G3')
14 axes[0, 2].scatter(cleandf.absences, cleandf.G3)
15 axes[0, 2].set_xlabel('absences')
16 axes[0, 2].set_ylabel('G3')
17 axes[1, 0].scatter(cleandf.G1, cleandf.G3)
18 axes[1, 0].set_xlabel('G1')
19 axes[1, 0].set_ylabel('G3')
20 axes[1, 1].scatter(cleandf.G2, cleandf.G3)
21 axes[1, 1].set_xlabel('G2')
22 axes[1, 1].set_ylabel('G3')
23 fig.delaxes(axes[1, 2])
24 plt.show()
25

```

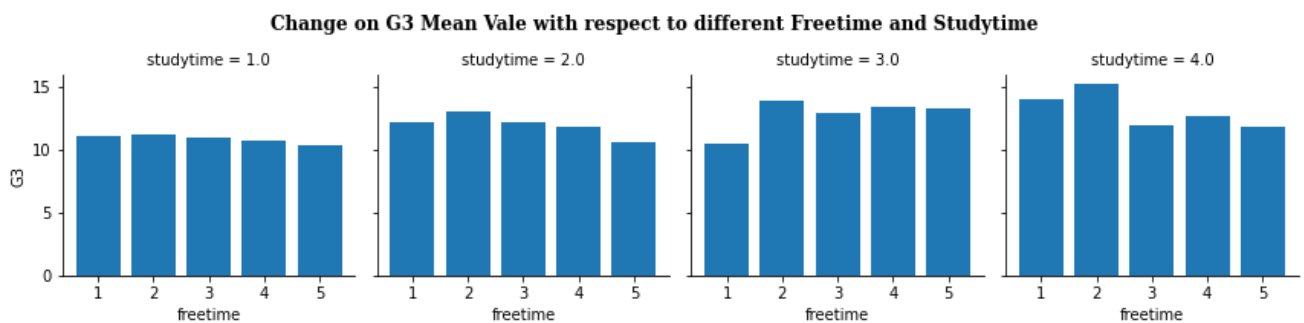
G3 Distribution for each continuous feature



Multivariant Relationship

Explore how student manage their time can affect their exam score.

```
1 def exam_mean_facetgrid(df,face_column,x_col,mean_col='G3',title=''):
2     plot_data = df.groupby([x_col,face_column]).mean()[mean_col].reset_index()
3     g=sns.FacetGrid(data=plot_data,col=face_column,margin_titles=True)
4     g.map(plt.bar,x_col,mean_col)
5     g.fig.subplots_adjust(top=0.8)
6     g.fig.suptitle(f'{title}',size=12, fontweight='bold', fontfamily='serif')
7
8 exam_mean_facetgrid(df=cleandf, face_column='studytime',
9 x_col='freetime',title='Change on G3 Mean Value with respect to different Freetime and
10 Studytime')
```



Obviously more study time will have better result.

We can also spot that student with studytime = 4.0 (> 10 hours) and freetime = 1&2 perform better than those student with freetime = 3&4&5.

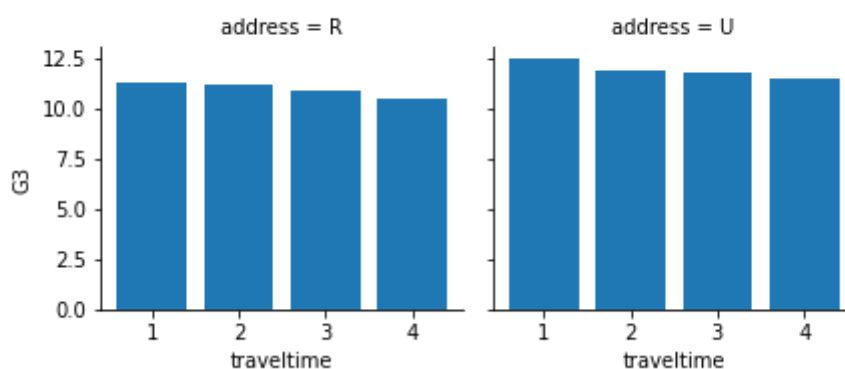
This may infere to following posibility:

1. Student study > 10 hours (probably 16 hours) and make him/her has less free time.
2. Student might has activity that can increase their performance. For example, exercise.

Explore how home address and travel time affect to exam result.

```
1 exam_mean_facetgrid(df=cleandf, face_column='address',
2 x_col='traveltime',title='Change on G3 Mean Value with respect to different Address
3 and Traveltime')
```

Change on G3 Mean Vale with respect to different Address and Traveltime



We can see when student stay in urban, they take less time to travel and perform better in exam.

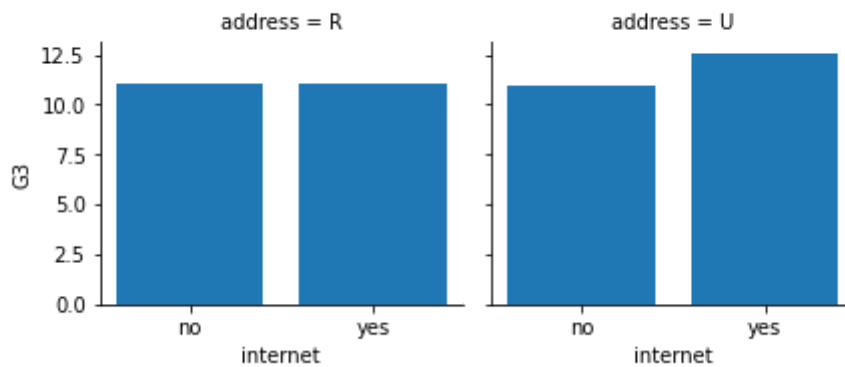
This may infer to following possibility:

1. Urban has a more convenient travel system.
2. Shorter travel time, more energetic to learn (no need to wake up so early).

Explore how home address and internet access affect exam performance.

```
1 exam_mean_facetgrid(df=cleandf, face_column='address', x_col='internet',title='Change on G3 Mean Value with respect to different Address and Internet Access')
```

Change on G3 Mean Vale with respect to different Address and Internet Access

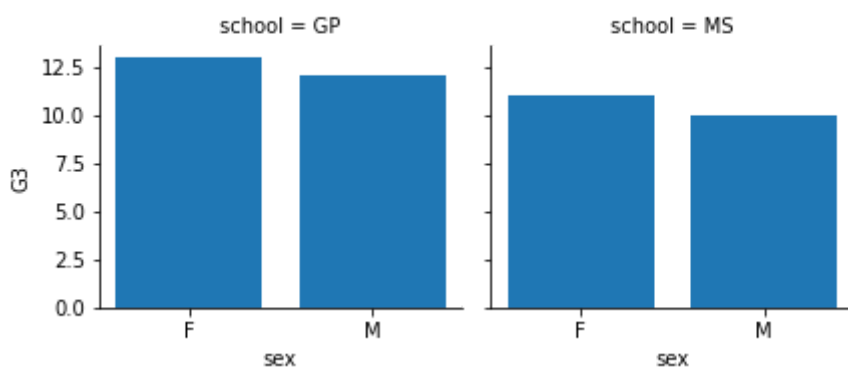


The students in urban more likely to have internet access and perform better in exam.(urban has more advance internet infrastructure-easier get access, easier to gain extra information for study)

From previous finding, female take higher median exam score in result. Explore how different school affect exam score for different sex.

```
1 exam_mean_facetgrid(df=cleandf, face_column='school', x_col='sex',title='Change on G3 Mean Value with respect to different school and sex')
```

Change on G3 Mean Vale with respect to different school and sex



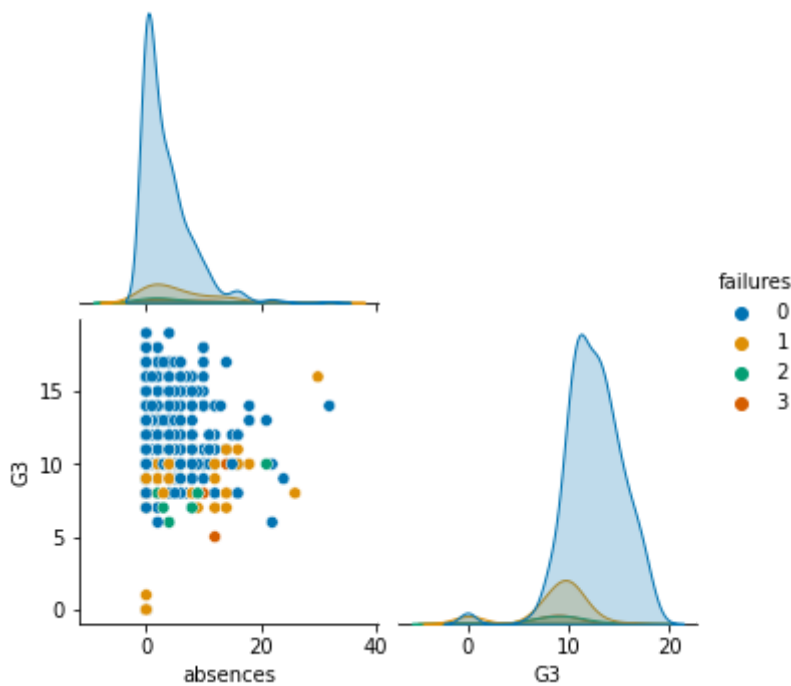
We can see female student contribute better result mo matter which school they belong to.

Explore how absences affect student failures and exam score.

```
1 pairplot_df = cleandf[['failures','absences','G3']]
2 sns.pairplot(pairplot_df, hue='failures', palette='colorblind', corner = True)
```



```
1 | <seaborn.axisgrid.PairGrid at 0x7f46e85cad10>
```



From the chart, we can observe that when the absences increase, they are tend to get more failure and lower exam score.

Data Modelling

Split into Train and Test Dataset

```
1 from sklearn.model_selection import train_test_split as tts
2 import time
3
4 # variables
5 x = cleandf.drop(['G3'], axis=1)
6
7 # target
8 y = cleandf[['G3']]
9
10 # Split training data and test data (70%/30%)
11 x_train, x_test, y_train, y_test = tts(x, y, test_size=0.3, random_state=40)
12
```

Feature Transformation and Standardization

```
1 from sklearn.preprocessing import OrdinalEncoder, StandardScaler
2 #encode categorical variable into ordinal
3 catvar = ['school', 'sex', 'address', 'Mjob', 'Fjob',
4           'reason', 'guardian', 'schoolsup', 'activities', 'nursery', 'higher', 'internet']
5 ordenc = OrdinalEncoder()
6 x_train[catvar] = ordenc.fit_transform(x_train[catvar])
```

```

6
7 #standardise x_train data
8 ss_x = StandardScaler()
9 x_train = pd.DataFrame(ss_x.fit_transform(x_train))
10 #x_train = pd.DataFrame(x_train)
11
12 #standardise y_train data
13 ss_y = StandardScaler()
14 y_train = ss_y.fit_transform(y_train)
15
16 #encode and standardise x_test and y_test
17 x_test[catvar] = ordenc.transform(x_test[catvar])
18 x_test = pd.DataFrame(ss_x.transform(x_test))
19 #y_test = ss_y.transform(y_test)

```

Model One: Random Forest Regressor

```

1 from sklearn.ensemble import RandomForestRegressor
2 from sklearn.model_selection import GridSearchCV
3
4 start_time = time.time()
5 #hyper parameter tuning
6 rfr_param_grid = {'n_estimators':range(50,101)}
7 rfrReg = GridSearchCV(RandomForestRegressor(random_state = 40), rfr_param_grid,
8 refit=True, verbose=0)
9 rfrReg.fit(x_train, y_train.ravel())
10
11 rfrtime = time.time()-start_time
12
13 print(rfrReg.best_estimator_)

```

```

1 | RandomForestRegressor(n_estimators=99, random_state=40)

```

Model Two: Support Vector Regressor

```

1 from sklearn import svm
2 from sklearn.model_selection import GridSearchCV
3
4 start_time = time.time()
5
6 #hyper parameter tuning
7 svr_param_grid = {'C': [0.1,1, 10, 100], 'gamma': [1,0.1,0.01,0.001], 'kernel':
8 ['rbf', 'poly', 'sigmoid']}
9 svrReg = GridSearchCV(svm.SVR(),svr_param_grid, refit=True, verbose=0)
10 svrReg.fit(x_train,y_train.ravel())
11
12 svrtime = time.time()-start_time
13
14 print(svrReg.best_estimator_)

```

```

1 | SVR(C=100, gamma=0.001, kernel='sigmoid')

```

Model Three: K Nearest Neighbors Regression Model

```

1 from sklearn.neighbors import KNeighborsRegressor
2 from sklearn.model_selection import GridSearchCV
3
4 start_time = time.time()
5
6 #hyper parameter tuning
7 knr_param_grid = {'n_neighbors': [3, 5, 7, 13, 15, 17, 19], 'weights':
  ['uniform', 'distance']}
8 knrReg = GridSearchCV(KNeighborsRegressor(), knr_param_grid, refit=True, verbose=0)
9 knrReg.fit(x_train, y_train.ravel())
10
11 knrtime = time.time()-start_time
12
13 print(knrReg.best_estimator_)

```

```

1 KNeighborsRegressor(n_neighbors=13, weights='distance')

```

Model Four: Decision Tree Regressor

```

1 from sklearn.tree import DecisionTreeRegressor
2
3 start_time = time.time()
4
5 DtReg = DecisionTreeRegressor(random_state = 40)
6 DtReg.fit(x_train, y_train)
7
8 dttime = time.time()-start_time

```

Model Five: Ridge Regressor

```

1 from sklearn.linear_model import Ridge
2
3 start_time = time.time()
4
5 rr_param_grid = {'alpha': [0, 0.2, 0.4, 0.6, 0.8, 1, 10, 100]}
6 ridgeReg = GridSearchCV(Ridge(), rr_param_grid, refit=True, verbose=0)
7 ridgeReg.fit(x_train, y_train)
8
9 ridgetime = time.time()-start_time
10
11 print(ridgeReg.best_estimator_)

```

```

1 Ridge(alpha=1)

```

Results Evalutaion

```

1 from sklearn.metrics import r2_score, mean_squared_error, mean_absolute_error
2
3 def modelacc(modelname, model, x_test, y_test, ss_y):
4     model_y_pred = model.predict(x_test)
5     model_y_pred = ss_y.inverse_transform(model_y_pred.reshape(-1,1))
6

```

```

7     modelr2 = r2_score(y_test, model_y_pred)
8     modelmse = mean_squared_error(y_test, model_y_pred)
9     modelmae = mean_absolute_error(y_test, model_y_pred)
10
11     reseval = pd.DataFrame({
12         'Model': [modelname],
13         'R2_score': [modelr2],
14         'MSE': [modelmse],
15         'MAE': [modelmae]
16     })
17
18     return(reseval)
19
20 #####
21 result = pd.DataFrame(columns= ["Model", "R2_score", "MSE", "MAE"])
22
23 #model one: rfr
24 rfrRegacc = modelacc('RFR', rfrReg, x_test, y_test, ss_y)
25 result = result.append(rfrRegacc, ignore_index = True)
26
27 #model two: svr
28 svrRegacc = modelacc('SVR', svrReg, x_test, y_test, ss_y)
29 result = result.append(svrRegacc, ignore_index = True)
30
31 #model three: KNR
32 neighacc = modelacc('KNR', knrReg, x_test, y_test, ss_y)
33 result = result.append(neighacc, ignore_index = True)
34
35 #model four: DT
36 DtRegacc = modelacc('DT', DtReg, x_test, y_test, ss_y)
37 result = result.append(DtRegacc, ignore_index = True)
38
39 #model five:
40 ridgeRegacc = modelacc('Ridge', ridgeReg, x_test, y_test, ss_y)
41 result = result.append(ridgeRegacc, ignore_index = True)
42
43 result['Time(s)'] = [rfrtime, svrtime, knrtime, dttime, ridgetime]
44 print(result)
45

```

	Model	R2_score	MSE	MAE	Time(s)
0	RFR	0.836988	1.967419	0.841388	46.833260
1	SVR	0.813252	2.253889	0.866753	4.999248
2	KNR	0.590293	4.944823	1.478549	0.671243
3	DT	0.620988	4.574359	1.117949	0.006966
4	Ridge	0.821261	2.157229	0.868879	0.333495

We see that random forest regressor (RFR) gives the best r2 score, but with a significantly longer time taken (especially when doing hyperparameter tuning). Ridge regressor comes close as second with a much shorter execution time. On the other hand, both k-nearest neighbour and decision tree regressor perform poorly on this set of data.

Summary

In summary, from 33 variables we have dropped 8 of it which are family size (famsize), parent's cohabitation status (Pstatus), family educational support (famsup), extra paid classes within the course subject (paid), with a romantic relationship (romantic), quality of family relationships (famrel), going out with friends (goout) and current health status (health) due to irrelevant score to the student grade. The criteria for dropping the variables is based on boxplot, scatterplot and pearson correlation coefficient and chi-square test. From exploratory data analysis, we can conclude that study time, travel time, student's home location, their absences and gender had affect their exam grade. Five machine learning-based regressor namely Random Forest Regressor, Support Vector Regressor, K Nearest Neighbors Regression Model, Decision Tree Regressor and Ridge Regressor were implemented in this project. As a result, Random Forest Regressor were significantly better than other regressor.

Limitation

1. One of the biggest limitation with K-NN is to choose the optimal number of neighbors to be consider while predicting the new data entry. Right number of neighbors is a must to produce the best prediction. Beside that, KNN also bad in large dataset. This is because, in large datasets, the cost of calculating the distance between the new point and each existing points is huge which degrades the performance of the algorithm.
2. Decision tree are inadequate to use when working with continuous variables as it will loses information when the variable are categorizes in different categories. As we can see that Decision Tree Regressor less performed on this dataset.

Member's Contribution

Member	Matric No.	Section
Lim Chin Woon (L)	s2124307	Feature Selection, EDA , Ridge Model.
Muhammad Amirul Daniel bin Badrul Hisham	s2115750	KNN Model, Limitation, Presentation.
Nurul Shahirah binti Sha'ari	s2120876	Related Work, Summary , Limitation, DT Model.
Qiyi Liu	s2106632	Overview and Motivation, Data Transformation, RFR Model.
Yeoh Li Tian	s2120306	Data Cleaning (except feature selection), Results Evaluation, SVR Model.

Research Question was by who again?
Maybe Data collection also?