

## Minimum Viable Product (MVP)

The goal of this project is do a deep analysis of what could be possible factor on whether a student is likely to get a high score or a low score. The data obtained does not contain that much information, but we will use Random Forest Regressor algorithm to predict what score a student will get in the foreseen feature by analysing the features

This jupyter\_notebook contains :

- importing data with panda for the file ("student-mat.csv")
- Exploring the data
- Dealing with missing value and drop the empty columns
- describe the data
- use Plotly for distribution plot for visualizing in three grades .

Extract the Data and Gather General Information of the Dataset

```
In [2]: import pandas as pd
import numpy as np
import plotly
from plotly import tools
```

```
df = pd.read_csv("student-mat.csv")
df.head()
```

Out[2]:

	school	sex	age	address	famsize	Pstatus	Medu	Fedu	Mjob	Fjob	...	famrel	freetime	goout	Dalc	Walc	health	absences	G1	G2	G3
0	GP	F	18	U	GT3	A	4	4	at_home	teacher	...	4	3	4	1	1	3	6	5	6	6
1	GP	F	17	U	GT3	T	1	1	at_home	other	...	5	3	3	1	1	3	4	5	5	6
2	GP	F	15	U	LE3	T	1	1	at_home	other	...	4	3	2	2	3	3	10	7	8	10
3	GP	F	15	U	GT3	T	4	2	health	services	...	3	2	2	1	1	5	2	15	14	15
4	GP	F	16	U	GT3	T	3	3	other	other	...	4	3	2	1	2	5	4	6	10	10

5 rows × 33 columns

Information about the Variables

```
In [4]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 395 entries, 0 to 394
Data columns (total 33 columns):
#   Column             Non-Null Count  Dtype  
---  --
0   school              395 non-null    object  
1   sex                 395 non-null    object  
2   age                 395 non-null    int64   
3   address             395 non-null    object  
4   famsize             395 non-null    object  
5   Pstatus             395 non-null    object  
6   Medu                395 non-null    int64   
7   Fedu                395 non-null    int64   
8   Mjob                395 non-null    object  
9   Fjob                395 non-null    object  
10  reason              395 non-null    object  
11  guardian            395 non-null    object  
12  traveltime          395 non-null    int64   
13  studytime           395 non-null    int64   
14  failures            395 non-null    int64   
15  schoolsup           395 non-null    object  
16  famsup              395 non-null    object  
17  paid                395 non-null    object  
18  activities          395 non-null    object  
19  nursery             395 non-null    object  
20  higher              395 non-null    object  
21  internet            395 non-null    object  
22  romantic            395 non-null    object  
23  famrel              395 non-null    int64   
24  freetime            395 non-null    int64   
25  goout               395 non-null    int64   
26  Dalc                395 non-null    int64   
27  Walc                395 non-null    int64   
28  health              395 non-null    int64   
29  absences            395 non-null    int64   
30  G1                  395 non-null    int64   
31  G2                  395 non-null    int64   
32  G3                  395 non-null    int64
```

```
File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 C
In [5]: df.describe()

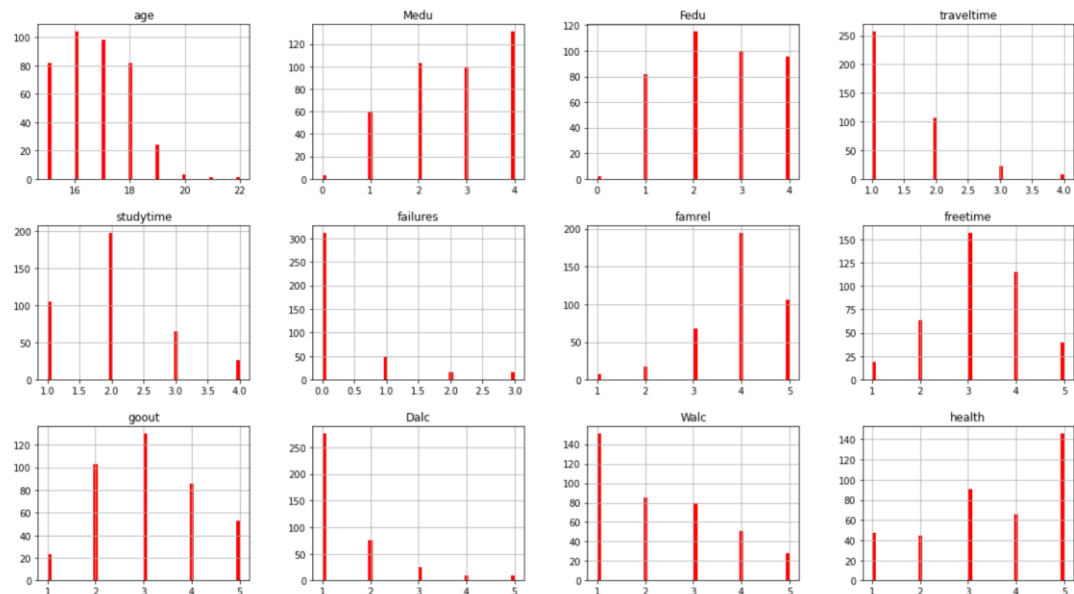
Out[5]:
```

	age	Medu	Fedu	traveltime	studytime	failures	famrel	freetime	goout	Dalc	Walc	health
count	395.000000	395.000000	395.000000	395.000000	395.000000	395.000000	395.000000	395.000000	395.000000	395.000000	395.000000	395.000000
mean	16.696203	2.749367	2.521519	1.448101	2.035443	0.334177	3.944304	3.235443	3.108861	1.481013	2.291139	3.554430
std	1.276043	1.094735	1.088201	0.697505	0.839240	0.743651	0.896659	0.998862	1.113278	0.890741	1.287897	1.390303
min	15.000000	0.000000	0.000000	1.000000	1.000000	0.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000
25%	16.000000	2.000000	2.000000	1.000000	1.000000	0.000000	4.000000	3.000000	2.000000	1.000000	1.000000	3.000000
50%	17.000000	3.000000	2.000000	1.000000	2.000000	0.000000	4.000000	3.000000	3.000000	1.000000	2.000000	4.000000
75%	18.000000	4.000000	3.000000	2.000000	2.000000	0.000000	5.000000	4.000000	4.000000	2.000000	3.000000	5.000000
max	22.000000	4.000000	4.000000	4.000000	4.000000	3.000000	5.000000	5.000000	5.000000	5.000000	5.000000	5.000000

Then we see how our Data is Distributed

```
In [6]: import matplotlib.pyplot as plt

df.hist(bins=50, figsize=(20,15), color='r')
plt.show()
```



Then distribution of students grades

Distribution of Student's Grades

