

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

Importing Most Frequently used Libraries for Data Analysis

```
In [2]: df= pd.read_csv(r"C:\Users\Apek\Desktop\Data Technical Proj\Expanded_data_with_more_features.csv")

In [3]: df
```

Out[3]:	Unnamed: 0	Gender	EthnicGroup	ParentEduc	LunchType	TestPrep	ParentMaritalStatus	PracticeSport	IsFirstChild	NrSiblings	TransportMeans	WklyStudyHours	MathScore	ReadingScore	WritingScore
0	0	female	NaN	bachelor's degree	standard	none	married	regularly	yes	3.0	school_bus	< 5	71		
1	1	female	group C	some college	standard	NaN	married	sometimes	yes	0.0	NaN	5 - 10	69		
2	2	female	group B	master's degree	standard	none	single	sometimes	yes	4.0	school_bus	< 5	87		
3	3	male	group A	associate's degree	free/reduced	none	married	never	no	1.0	NaN	5 - 10	45		
4	4	male	group C	some college	standard	none	married	sometimes	yes	0.0	school_bus	5 - 10	76		
...
30636	816	female	group D	high school	standard	none	single	sometimes	no	2.0	school_bus	5 - 10	59		
30637	890	male	group E	high school	standard	none	single	regularly	no	1.0	private	5 - 10	58		
30638	911	female	NaN	high school	free/reduced	completed	married	sometimes	no	1.0	private	5 - 10	61		
30639	934	female	group D	associate's degree	standard	completed	married	regularly	no	3.0	school_bus	5 - 10	82		
30640	960	male	group B	some college	standard	none	married	never	no	1.0	school_bus	5 - 10	64		

30641 rows × 15 columns

Before Analysis Summarize or Understand Data

```
In [4]: print(df.head())

Out[5]: Index(['Unnamed: 0', 'Gender', 'EthnicGroup', 'ParentEduc', 'LunchType', 'TestPrep', 'ParentMaritalStatus', 'PracticeSport', 'IsFirstChild', 'NrSiblings', 'TransportMeans', 'WklyStudyHours', 'MathScore', 'ReadingScore', 'WritingScore'], dtype='object')

In [6]: df.describe() ## It gives the Central Tendency of Numerical Fields Present in Dataset

Out[6]:
```

	Unnamed: 0	NrSiblings	MathScore	ReadingScore	WritingScore
count	30641.000000	29089.000000	30641.000000	30641.000000	30641.000000
mean	499.556607	2.145894	66.558402	69.377533	68.418622
std	288.747894	1.456242	15.361616	14.758952	15.442525
min	0.000000	0.000000	0.000000	10.000000	4.000000
25%	249.000000	1.000000	56.000000	59.000000	58.000000
50%	500.000000	2.000000	67.000000	70.000000	69.000000
75%	750.000000	3.000000	78.000000	80.000000	79.000000
max	999.000000	7.000000	100.000000	100.000000	100.000000

```
In [7]: df.info() ## Gives the information of each column Data type and Not null count

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30641 entries, 0 to 30640
Data columns (total 15 columns):
# Column Non-Null Count Dtype
---  ---
0 Unnamed: 0 30641 non-null int64
1 Gender 30641 non-null object
2 EthnicGroup 28891 non-null object
3 ParentEduc 28796 non-null object
4 LunchType 30641 non-null object
5 TestPrep 28811 non-null object
6 ParentMaritalStatus 29451 non-null object
7 PracticeSport 30639 non-null object
8 IsFirstChild 29737 non-null object
9 NrSiblings 29069 non-null float64
10 TransportMeans 27587 non-null object
11 WklyStudyHours 29686 non-null object
12 MathScore 30641 non-null int64
13 ReadingScore 30641 non-null int64
14 WritingScore 30641 non-null int64
dtypes: float64(1), int64(4), object(10)
memory usage: 3.5+ MB

In [8]: ## Counting Null Values in DataSet
df.isna().sum()

Out[8]:
```

Unnamed: 0	Gender	EthnicGroup	ParentEduc	LunchType	TestPrep	ParentMaritalStatus	PracticeSport	IsFirstChild	NrSiblings	TransportMeans	WklyStudyHours	MathScore	ReadingScore	WritingScore	
0	0	NaN	bachelor's degree	standard	none	married	regularly	yes	3.0	school_bus	< 5	71			
1	1	female	group C	some college	standard	NaN	married	sometimes	yes	0.0	NaN	5 - 10	69		
2	2	female	group B	master's degree	standard	none	single	sometimes	yes	4.0	school_bus	< 5	87		
3	3	male	group A	associate's degree	free/reduced	none	married	never	no	1.0	NaN	5 - 10	45		
4	4	male	group C	some college	standard	none	married	sometimes	yes	0.0	school_bus	5 - 10	76		

Drop unused Fields from DataFrame

```
In [9]: df=df.drop('Unnamed: 0',axis=1)

In [10]: df.head()

Out[10]:
```

Gender	EthnicGroup	ParentEduc	LunchType	TestPrep	ParentMaritalStatus	PracticeSport	IsFirstChild	NrSiblings	TransportMeans	WklyStudyHours	MathScore	ReadingScore	WritingScore
female	NaN	bachelor's degree	standard	none	married	regularly	yes	3.0	school_bus	< 5	71		
1 female	group C	some college	standard	NaN	married	sometimes	yes	0.0	NaN	5 - 10	69		
2 female	group B	master's degree	standard	none	single	sometimes	yes	4.0	school_bus	< 5	87		
3 male	group A	associate's degree	free/reduced	none	married	never	no	1.0	NaN	5 - 10	45		
4 male	group C	some college	standard	none	married	sometimes	yes	0.0	school_bus	5 - 10	76		

As minimum value in Maths subject is zero filtering out their Gender, ParentEduc, ParentMaritalStatus and PracticeSport

```
In [11]: df[df['MathScore']==0][['Gender', 'ParentEduc', 'ParentMaritalStatus', 'PracticeSport']]

Out[11]:
```

Gender	ParentEduc	ParentMaritalStatus	PracticeSport
55 female	some high school	single	regularly

Evaluate the Records of group A ,group B category from Ethnic Group and count the total records

```
In [12]: eg=df[df['EthnicGroup'].isin(['group A','group B'])]
eg

Out[12]:
```

Gender	EthnicGroup	ParentEduc	LunchType	TestPrep	ParentMaritalStatus	PracticeSport	IsFirstChild	NrSiblings	TransportMeans	WklyStudyHours	MathScore	ReadingScore
2 female	group B	master's degree	standard	none	single	sometimes	yes	4.0	school_bus	< 5	87	
3 male	group A	associate's degree	free/reduced	none	married	never	no	1.0	NaN	5 - 10	45	
5 female	group B	associate's degree	standard	none	married	regularly	yes	1.0	school_bus	5 - 10	73	
6 female	group B	some college	standard	completed	widowed	never	no	1.0	private	5 - 10	85	
7 male	group B	some college	free/reduced	none	married	sometimes	yes	1.0	private	> 10	41	
...
30627 female	group A	high school	standard	completed	married	never	no	NaN	school_bus	> 10	58	
30628 female	group B	NaN	free/reduced	none	single	sometimes	no	1.0	school_bus	5 - 10	55	
30630 male	group B	associate's degree	free/reduced	none	married	sometimes	no	4.0	private	5 - 10	43	
30634 male	group A	associate's degree	free/reduced	completed	NaN	sometimes	no	2.0	school_bus	5 - 10	65	
30640 male	group B	some college	standard	none	married	never	no	1.0	school_bus	5 - 10	64	

8045 rows × 14 columns

```
In [13]: ## Isin Function can be used to filter rows based on multiple values

In [14]: eg.shape[0]

Out[14]: 8845
```

Calculate the students records who have scored greater than 70 marks in each subject

```
In [37]: sr=df[(df['MathScore']>70) & (df['ReadingScore']>70) & (df['WritingScore']>70)]

In [38]: sr

Out[38]:
```

Gender	EthnicGroup	ParentEduc	LunchType	TestPrep	ParentMaritalStatus	PracticeSport	IsFirstChild	NrSiblings	TransportMeans	WklyStudyHours	MathScore	ReadingScore
0 female	NaN	bachelor's degree	standard	none	married	regularly	yes	3.0	school_bus	< 5	71	
2 female	group B	master's degree	standard	none	single	sometimes	yes	4.0	school_bus	< 5	87	
4 male	group C	some college	standard	none	married	sometimes	yes	0.0	school_bus	5 - 10	76	
5 female	group B	associate's degree	standard	none	married	regularly	yes	1.0	school_bus	5 - 10	73	
6 female	group B	some college	standard	completed	widowed	never	no	1.0	private	5 - 10	85	
...
30620 male	group C	bachelor's degree	standard	none	married	NaN	no	4.0	private	5 - 10	91	
30623 female	group B	NaN	standard	completed	married	regularly	no	2.0	NaN	5 - 10	75	
30632 female	group D	some college	standard	none	married	regularly	no	3.0	private	5 - 10	82	
30633 female	group C	master's degree	standard	completed	married	never	no	2.0	school_bus	5 - 10	84	
30639 female	group D	associate's degree	standard	completed	married	regularly	no	3.0	school_bus	5 - 10	82	

9747 rows × 14 columns

Gender Distribution

```
In [15]: plt.figure(figsize=(4,5))
plt.title('Gender Distribution')
ax=sns.countplot(data=df, x='Gender')
ax.bar_label(ax.containers[0]) ## Adding Labels to know exact count
plt.show()
```



from the above chart we have analysed that female count is more than male count in the Dataset

```
In [ ]:
```

Analysing What Parent's Education is Impacting on the Scores of the Students

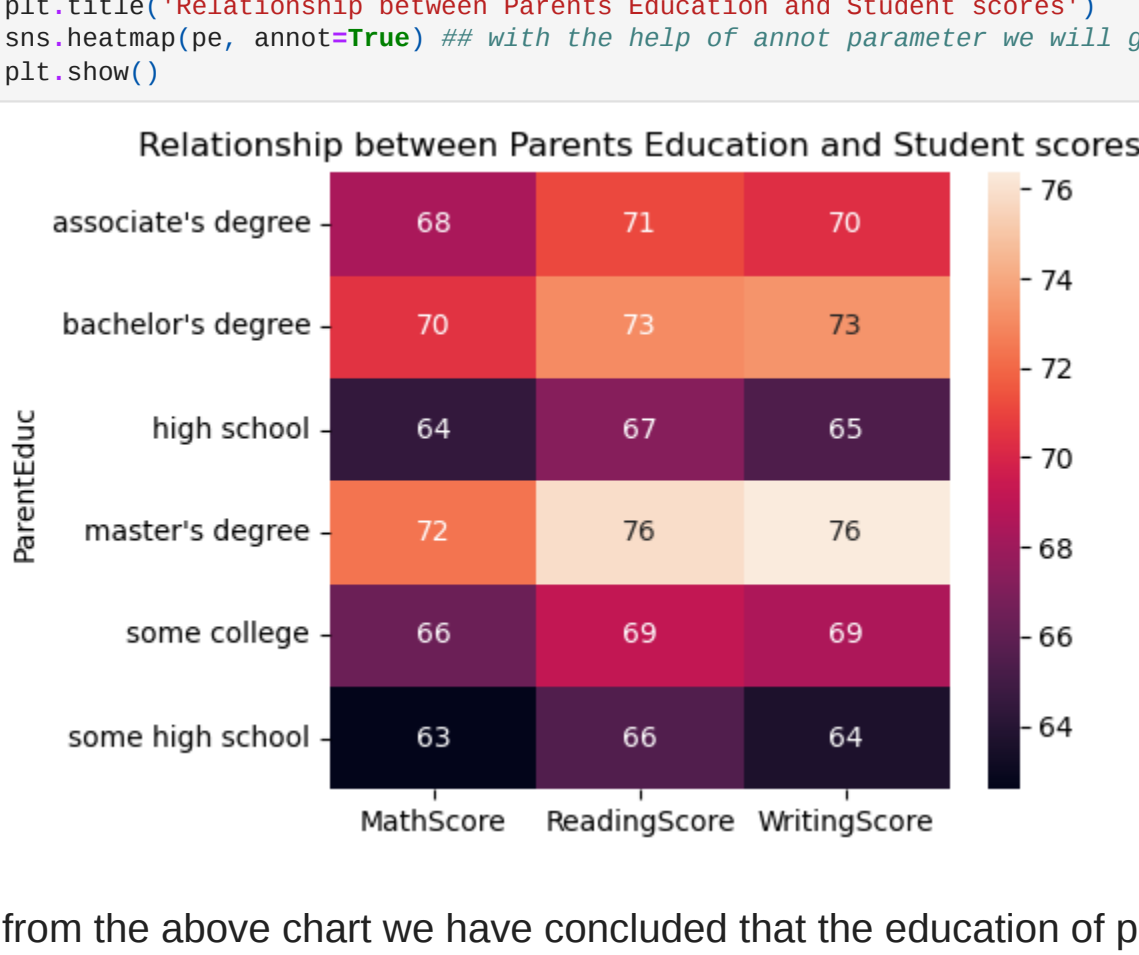
```
In [16]: pe=df.groupby('ParentEduc').agg({'MathScore':'mean','ReadingScore':'mean','WritingScore':'mean'})

In [17]: print(pe)

ParentEduc  MathScore  ReadingScore  WritingScore
associate's degree  68.365586  71.124324  70.299899
bachelor's degree  70.466627  73.062620  73.331069
high school  64.425721  67.221997  65.421136
master's degree  72.336134  75.832921  76.356896
some college  66.398472  69.179708  68.591432
some high school  62.584613  65.518785  63.632409

In [18]: ## Visualizing the above Data with the Heatmap

In [19]: plt.figure(figsize=(5,4))
plt.title('Relationship between Parents Education and Student scores')
sns.heatmap(pe, annot=True) ## with the help of annot parameter we will get the values in to the heatmap
plt.show()
```



from the above chart we have concluded that the education of parent's have a good impact on student score's

```
In [ ]:
```

Analysing Does Parent's Marital Status Has Impact on student's score

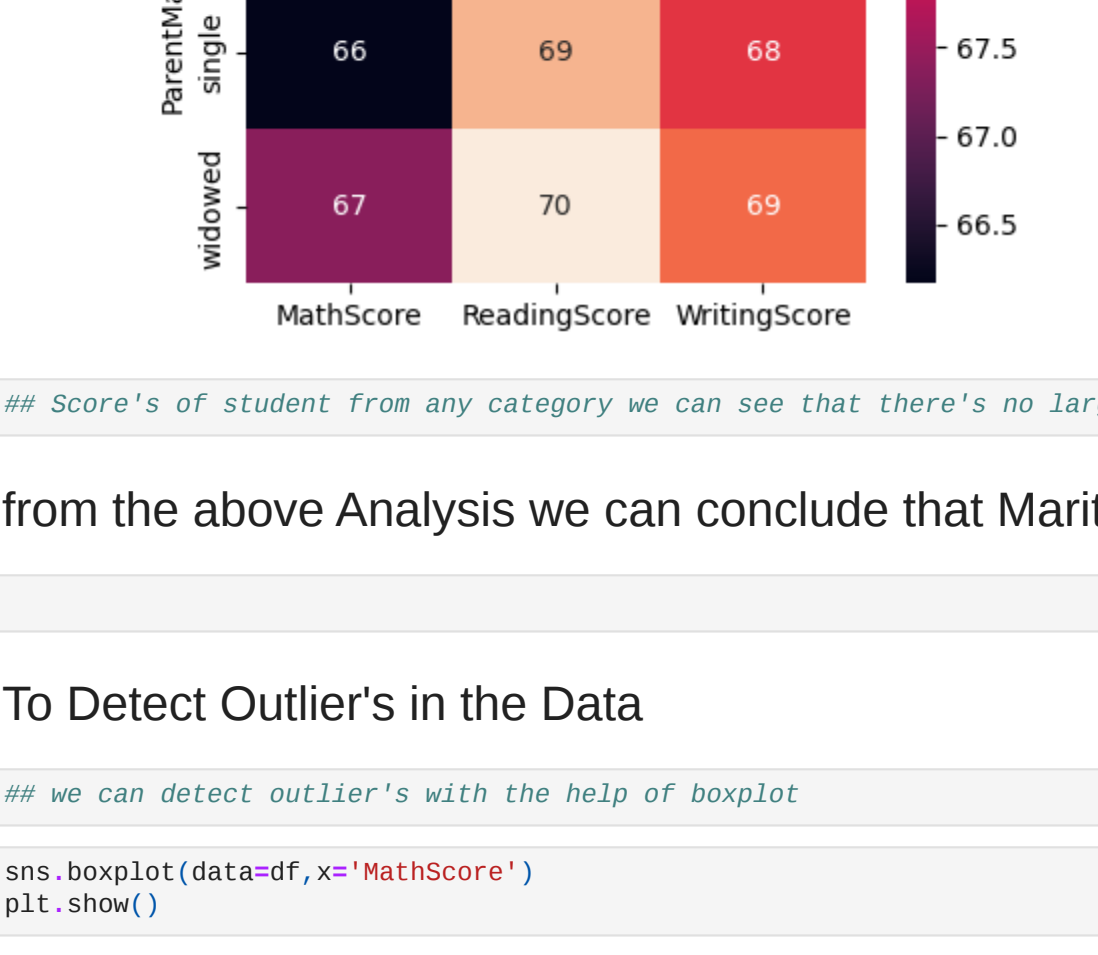
```
In [20]: pm=df.groupby('ParentMaritalStatus').agg({'MathScore':'mean','ReadingScore':'mean','WritingScore':'mean'})

In [21]: pm

Out[21]:
```

ParentMaritalStatus	MathScore	ReadingScore	WritingScore
divorced	66.601197	69.655011	68.799146
married	66.657326	69.389575	68.420981
single	66.105704	69.157250	68.174440
widowed	67.368866	69.651438	68.563452

```
In [22]: plt.figure(figsize=(5,4))
plt.title('Relationship between Parent's Marital Status and Student score's')
sns.heatmap(pm, annot=True)
plt.show()
```



Score's of student from any category we can see that there's no large difference

from the above Analysis we can conclude that Marital Status has no impact on student score's

```
In [ ]:
```

To Detect Outlier's in the Data

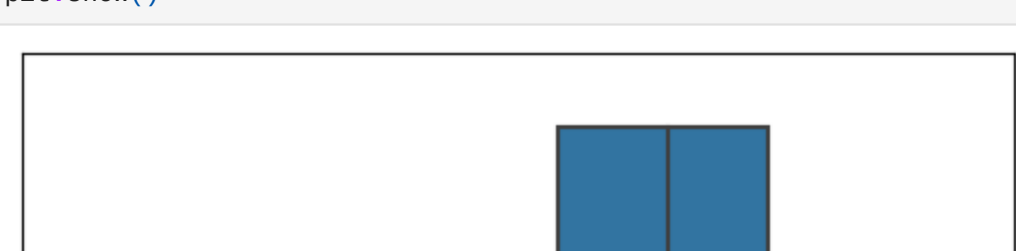
```
In [24]: ## we can detect outlier's with the help of boxplot
```

```
In [25]: sns.boxplot(data=df, x='MathScore')
plt.show()
```



```
In [26]: ## Above visual give's the idea that the boxplot is left skewed
## Outlier's are present in the range of left hand side (that is datapoints are lying in the range of 0-20 values)
```

```
In [27]: sns.boxplot(data=df, x='ReadingScore')
plt.show()
```



```
In [ ]:
```

```
In [28]: sns.boxplot(data=df, x='WritingScore')
plt.show()
```



From the above Visual we can conclude that one of the datapoint is lying on zero value in MathScore. Hence student's are having difficulties to score zero's in Maths Subject

```
In [ ]:
```

How many group's are present in EthnicGroup

```
In [29]: print(df['EthnicGroup'].unique())
[nan 'group C' 'group B' 'group A' 'group D' 'group E']

In [30]: ## In the above output we can see the group's

In [31]: GroupA=df[(df['EthnicGroup']=='group A')].count()
GroupB=df[(df['EthnicGroup']=='group B')].count()
GroupC=df[(df['EthnicGroup']=='group C')].count()
GroupD=df[(df['EthnicGroup']=='group D')].count()
GroupE=df[(df['EthnicGroup']=='group E')].count()

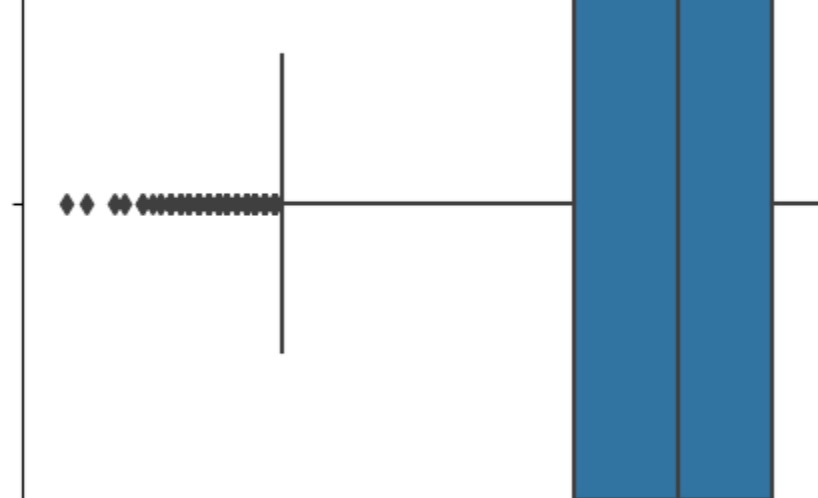
In [32]: print(GroupA)

Gender 2219
EthnicGroup 2219
ParentEduc 2078
LunchType 2219
TestPrep 2081
ParentMaritalStatus 2121
PracticeSport 2167
IsFirstChild 2168
NrSiblings 2096
TransportMeans 1999
WklyStudyHours 2146
MathScore 2219
ReadingScore 2219
WritingScore 2219
dtype: int64

In [33]: h1=[GroupA, 'GroupA', 'GroupC', 'GroupD', 'GroupE']
h1list=[GroupA['EthnicGroup'],GroupC['EthnicGroup'],GroupD['EthnicGroup'],GroupE['EthnicGroup']]

plt.pie(h1list, labels=h1, autopct='%1.2f%%')
plt.title('Distribution of Ethnic Groups')
plt.show()
```

Distribution of Ethnic Groups

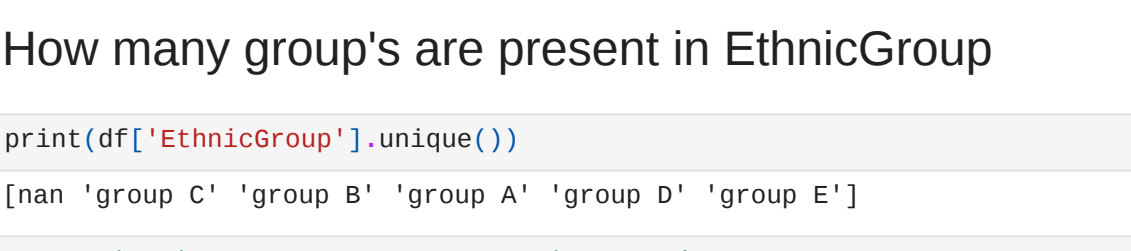


```
In [34]: print(h1list)

[2219, 5826, 9212, 7503, 4041]
```

Cross checking the Ethnic Group Values

```
In [35]: ab=sns.countplot(data=df, x='EthnicGroup')
ab.bar_label(ab.containers[0])
plt.show()
```



```
In [ ]:
```