

Installing Required Libraries

```
In [2]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
print("All imported")

All imported

In [3]: df = pd.read_csv("Expanded_data_with_more_features.csv")
print(df.head())

Unnamed: 0  Gender  EthnicGroup  ParentEduc  LunchType  TestPrep  \
0          0  female         NaN  bachelor's degree  standard    none
1          1  female    group C    some college  standard    NaN
2          2  female    group B  master's degree  standard    none
3          3  male     group A  associate's degree  free/reduced  none
4          4  male     group C    some college  standard    none

ParentMaritalStatus  PracticeSport  IsFirstChild  NrSiblings  TransportMeans  \
0      married      regularly      yes          3.0      school_bus
1      married      sometimes     yes          0.0      NaN
2      single       sometimes     yes          4.0      school_bus
3      married      never         no           1.0      NaN
4      married      sometimes     yes          0.0      school_bus

WklyStudyHours  MathScore  ReadingScore  WritingScore
0              < 5        71           71           74
1      5 - 10        69           90           88
2              < 5        87           93           91
3      5 - 10        45           56           42
4      5 - 10        76           78           75

In [8]: df.describe()

Out[8]:
```

	Unnamed: 0	NrSiblings	MathScore	ReadingScore	WritingScore
count	30641.000000	29069.000000	30641.000000	30641.000000	30641.000000
mean	499.556607	2.145894	66.558402	69.377533	68.418622
std	288.747894	1.458242	15.361616	14.758952	15.443525
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 [9]: df.info()

<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         28801 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       30810 non-null  object
8   IsFirstChild        29737 non-null  object
9   NrSiblings          29069 non-null  float64
10  TransportMeans       27507 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 [ ]:

In [10]: df.isnull().sum()

Out[10]:
```

	Unnamed: 0	Gender	EthnicGroup	ParentEduc	LunchType	TestPrep
Gender	0	0	1840	1845	0	1830
ParentEduc	0	1190	1190	931	994	1572
TransportMeans	0	3134	955	0	0	0
WklyStudyHours	0	0	0	0	0	0
MathScore	0	0	0	0	0	0
ReadingScore	0	0	0	0	0	0
WritingScore	0	0	0	0	0	0
dtype:	int64	object	object	object	object	object

Dropping unnamed column

```
In [4]: df = df.drop("Unnamed: 0", axis=1)
print(df.head())

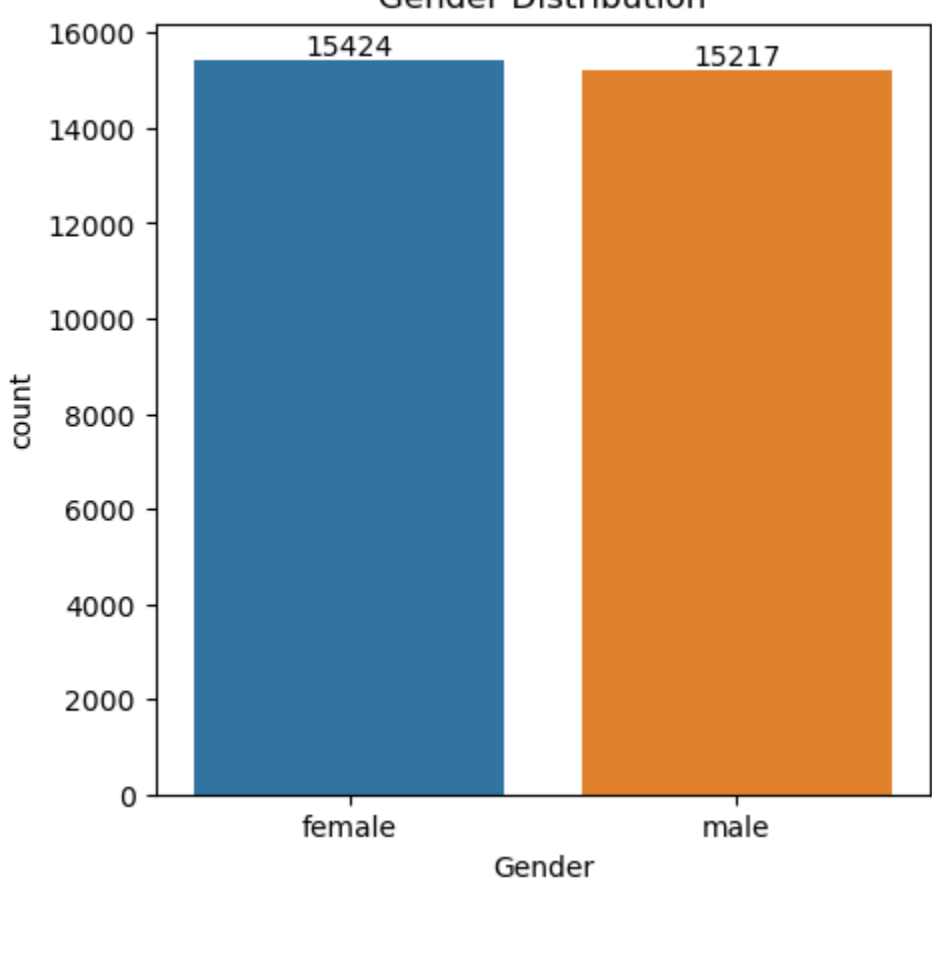
Gender  EthnicGroup  ParentEduc  LunchType  TestPrep  \
0  female         NaN  bachelor's degree  standard    none
1  female    group C    some college  standard    NaN
2  female    group B  master's degree  standard    none
3  male     group A  associate's degree  free/reduced  none
4  male     group C    some college  standard    none

ParentMaritalStatus  PracticeSport  IsFirstChild  NrSiblings  TransportMeans  \
0      married      regularly      yes          3.0      school_bus
1      married      sometimes     yes          0.0      NaN
2      single       sometimes     yes          4.0      school_bus
3      married      never         no           1.0      NaN
4      married      sometimes     yes          0.0      school_bus

WklyStudyHours  MathScore  ReadingScore  WritingScore
0              < 5        71           71           74
1      5 - 10        69           90           88
2              < 5        87           93           91
3      5 - 10        45           56           42
4      5 - 10        76           78           75
```

Gender Distribution

```
In [34]: plt.figure(figsize=(5,5))
ax = sns.countplot(data=df, x = "Gender")
ax.bar_label(ax.containers[0])
plt.title("Gender Distribution")
plt.show()
```

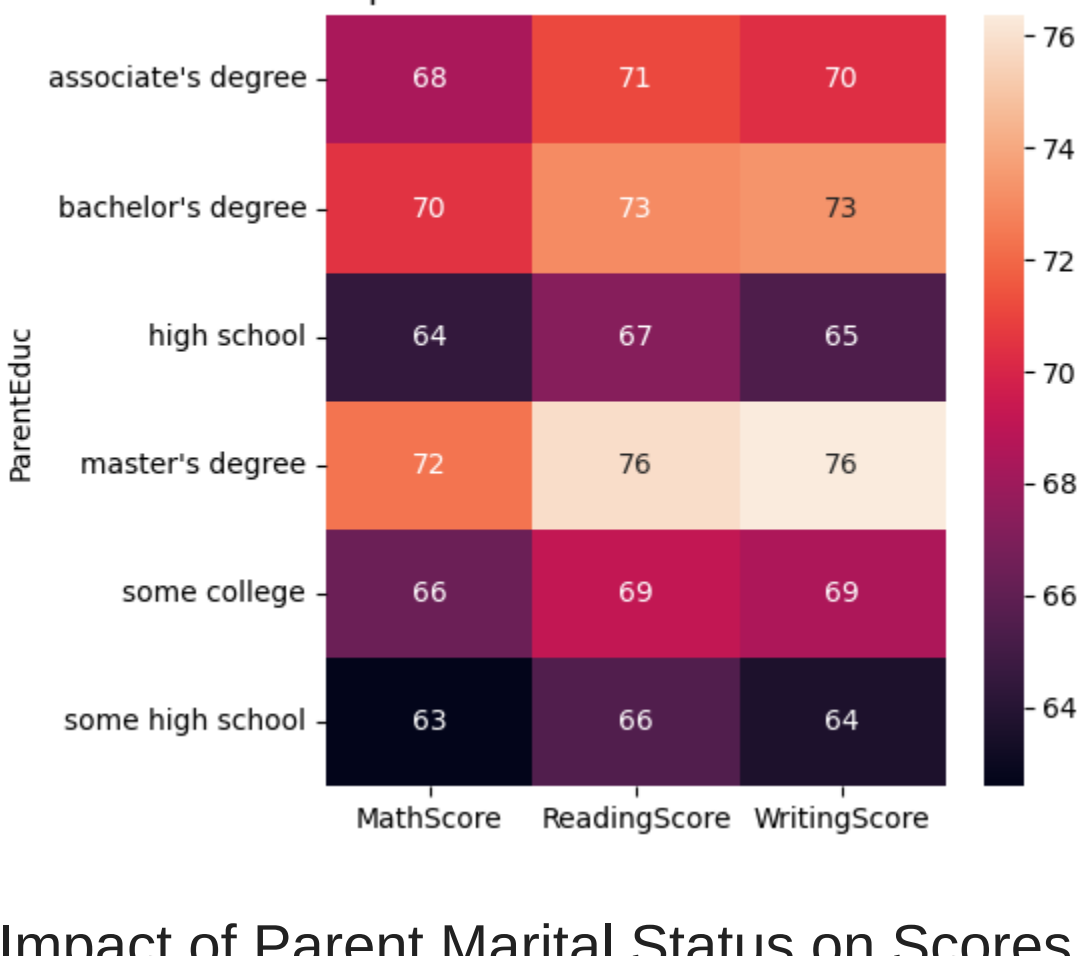


Impact of Parent Education on scores

```
In [14]: gb = df.groupby("ParentEduc").agg({'MathScore': 'mean', 'ReadingScore': 'mean', 'WritingScore': "mean"})
print(gb)
```

ParentEduc	MathScore	ReadingScore	WritingScore
associate's degree	68.365586	71.124324	70.299099
bachelor's degree	70.466627	73.062020	73.331069
high school	64.435781	67.213997	65.421136
master's degree	72.336134	75.832921	76.356896
some college	66.390472	69.179708	68.501432
some high school	62.584013	65.510785	63.632409

```
In [35]: plt.figure(figsize=(5,5))
sns.heatmap(gb, annot = True)
plt.title("Impact of Parent Education on scores")
plt.show()
```

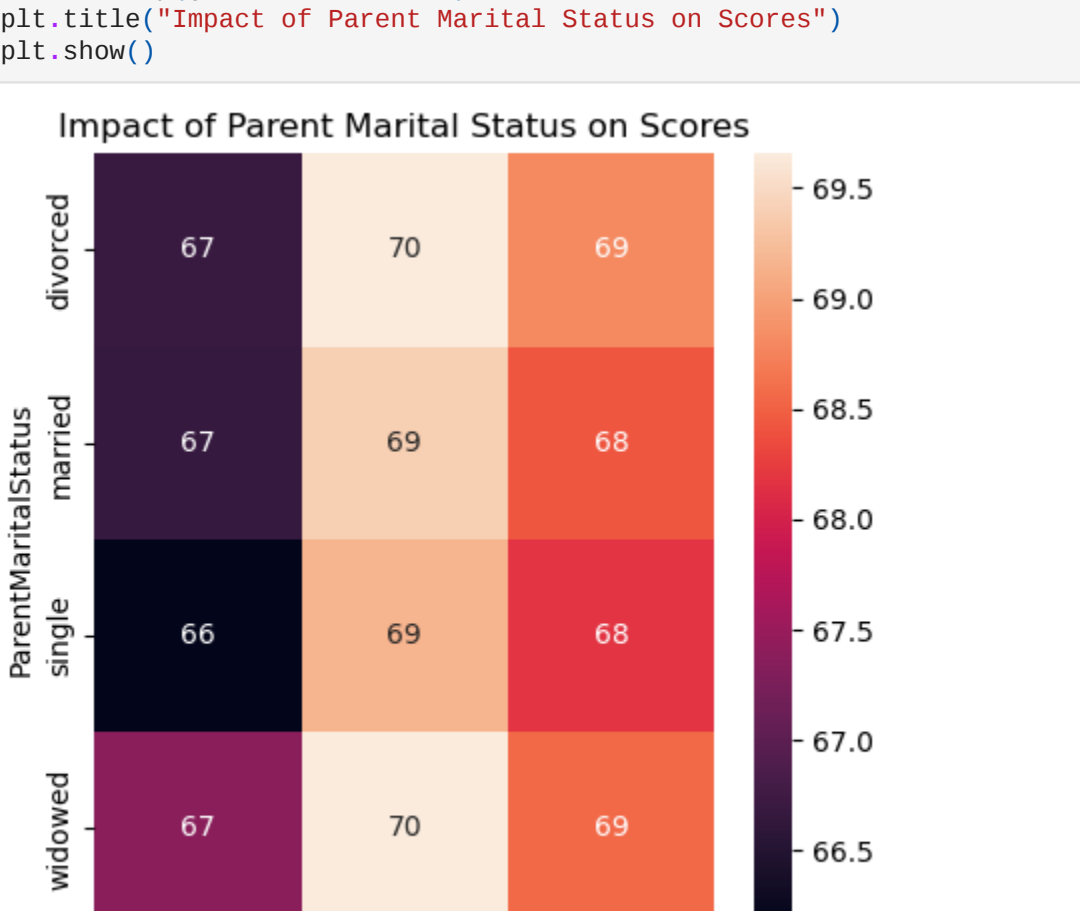


Impact of Parent Marital Status on Scores

```
In [28]: gb1 = df.groupby("ParentMaritalStatus").agg({'MathScore': 'mean', 'ReadingScore': 'mean', 'WritingScore': "mean"})
print(gb1)
```

ParentEduc	MathScore	ReadingScore	WritingScore
associate's degree	68.365586	71.124324	70.299099
bachelor's degree	70.466627	73.062020	73.331069
high school	64.435781	67.213997	65.421136
master's degree	72.336134	75.832921	76.356896
some college	66.390472	69.179708	68.501432
some high school	62.584013	65.510785	63.632409

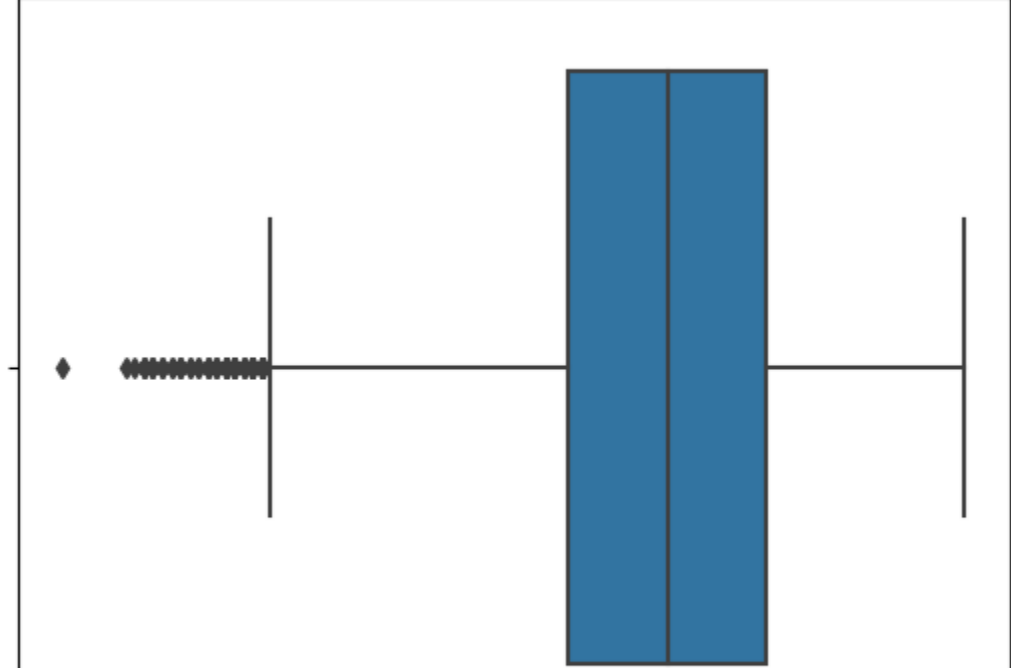
```
In [37]: plt.figure(figsize=(5,5))
sns.heatmap(gb1, annot = True)
plt.title("Impact of Parent Marital Status on Scores")
plt.show()
```



Checking for Outliers

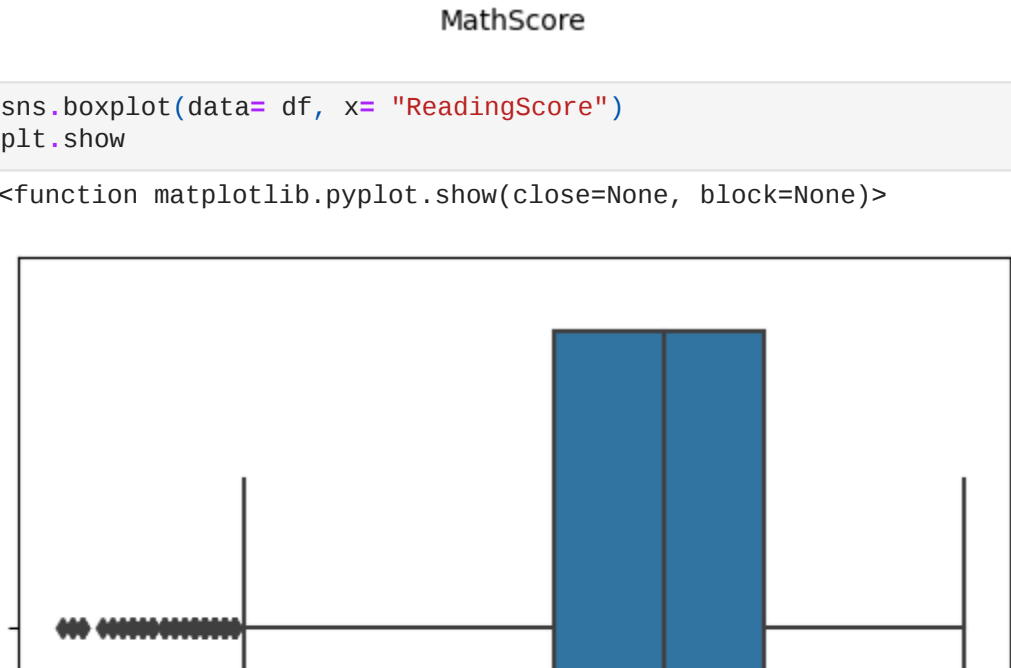
```
In [39]: sns.boxplot(data= df, x= "MathScore")
plt.show()

<function matplotlib.pyplot.show(close=None, block=None)>
```



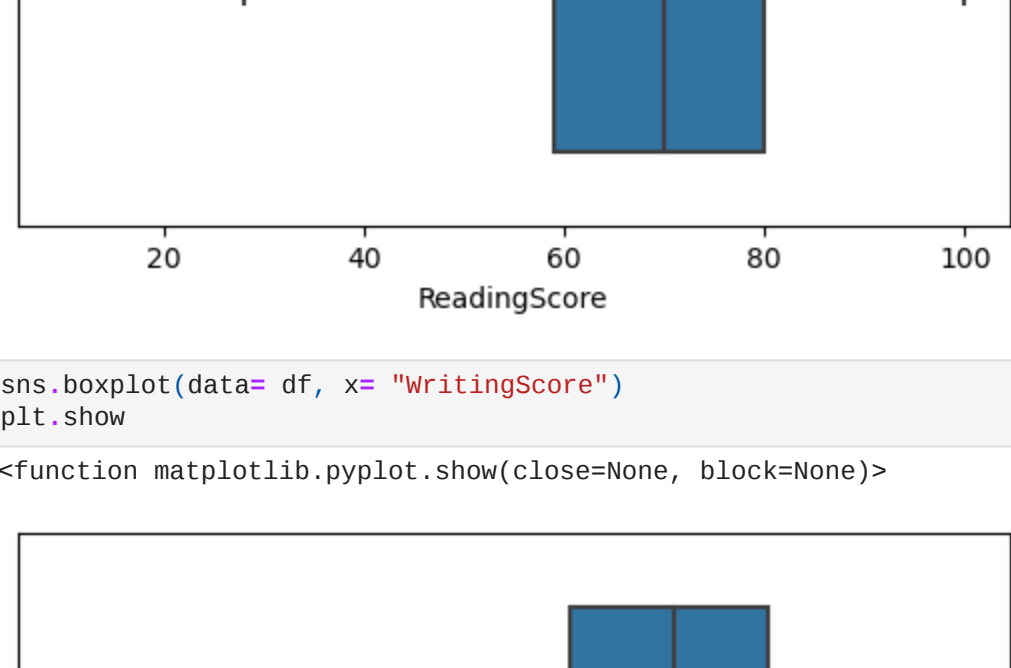
```
In [40]: sns.boxplot(data= df, x= "ReadingScore")
plt.show()

<function matplotlib.pyplot.show(close=None, block=None)>
```



```
In [41]: sns.boxplot(data= df, x= "WritingScore")
plt.show()

<function matplotlib.pyplot.show(close=None, block=None)>
```



Distribution Of Ethnic Groups

```
In [42]: print(df["EthnicGroup"].unique())

[nan 'group C' 'group B' 'group A' 'group D' 'group E']
```

```
In [60]: groupA = df.loc[(df['EthnicGroup'] == "group A")].count()
group = df.loc[(df['EthnicGroup'] == "group B")].count()
groupC = df.loc[(df['EthnicGroup'] == "group C")].count()
groupD = df.loc[(df['EthnicGroup'] == "group D")].count()
groupE = df.loc[(df['EthnicGroup'] == "group E")].count()
```

```
m1 = [groupA['EthnicGroup'], groupB['EthnicGroup'], groupC['EthnicGroup'], groupD['EthnicGroup'], groupE['EthnicGroup']]
l1 = ['Group A', 'Group B', 'Group C', 'Group D', 'Group E']
plt.title("Distribution Of Ethnic Groups")
plt.pie(m1, labels=l1, autopct='%1.0f%%')
```

```
Out[60]: ([<matplotlib.patches.Wedge at 0x21cfd4159d0>,
<matplotlib.patches.Wedge at 0x21cfd423130>,
<matplotlib.patches.Wedge at 0x21cfd423850>,
<matplotlib.patches.Wedge at 0x21cfd423f70>,
<matplotlib.patches.Wedge at 0x21cfd4316d0>],
[Text(1.067934312184989, 0.26365943230411354, 'Group A'),
Text(0.479656756487504, 0.98913832591447, 'Group B'),
Text(-1.0208489310562472, 0.40971631644507106, 'Group C'),
Text(-0.14172979830625182, -1.09083118046381, 'Group D'),
Text(0.9948566425203025, -0.4693189329584214, 'Group E')],
[Text(0.582509627028272, 0.14381423589224373, '6%'),
Text(0.26163095808408303, 0.5390529995953347, '20%'),
Text(-0.5568266896670430, 0.2234816271518569, '32%'),
Text(-0.07730716271250097, -0.5949988257075327, '26%'),
Text(0.5426490777383467, -0.255992145250048, '14%')])
```



```
In [63]: ax = sns.countplot(data=df, x= "EthnicGroup")
ax.bar_label(ax.containers[0])
```

```
Out[63]: [Text(0, 0, '9212'),
Text(0, 0, '5826'),
Text(0, 0, '2219'),
Text(0, 0, '7503'),
Text(0, 0, '4041')]
```

