

```
In [4]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

In [5]: 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 [7]: df.describe()

Out[7]:
```

	Unnamed: 0	N/Siblings	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	288.747094	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 [8]: 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   29403 non-null  object
 7  PracticeSport         30010 non-null  object
 8  IsFirstChild          29737 non-null  object
 9  Nrsiblings            29069 non-null  float64
10  TransportMeans        27507 non-null  object
11  WklyStudyHours        29666 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 [10]: df.isnull().sum()

Out[10]:
```

Unnamed: 0	Gender	EthnicGroup	ParentEduc	LunchType	TestPrep
0	0	0	1840	0	0
1	1	1840	1845	0	0
2	2	1845	1830	1190	631
3	3	1830	1190	631	904
4	4	1190	631	904	1572
5	5	631	904	1572	3134
6	6	904	1572	3134	955
7	7	1572	3134	955	0
8	8	3134	955	0	0
9	9	955	0	0	0
10	10	0	0	0	0
11	11	0	0	0	0
12	12	0	0	0	0
13	13	0	0	0	0
14	14	0	0	0	0

```
In [11]: df.drop('Unnamed: 0',axis = 1)
print(df.head())

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

#Drop unnamed column

In [11]: df = df.drop("Unnamed: 0",axis = 1)
print(df.head())

Gender EthnicGroup  ParentEduc  LunchType TestPrep  ParentMaritalStatus  PracticeSport  IsFirstChild  Nrsiblings  TransportMeans  WklyStudyHours  MathScore  ReadingScore  WritingScore
0  female         NaN  bachelor's degree      standard    none      married      regularly      yes          3.0      school_bus
1  female      group C      some college      standard    NaN      married      sometimes    yes          0.0         NaN
2  female      group B  master's degree      standard    none      single      sometimes    yes          4.0      school_bus
3  male        group A  associate's degree  free/reduced    none      married      never        no          1.0         NaN
4  male        group C      some college      standard    none      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
```

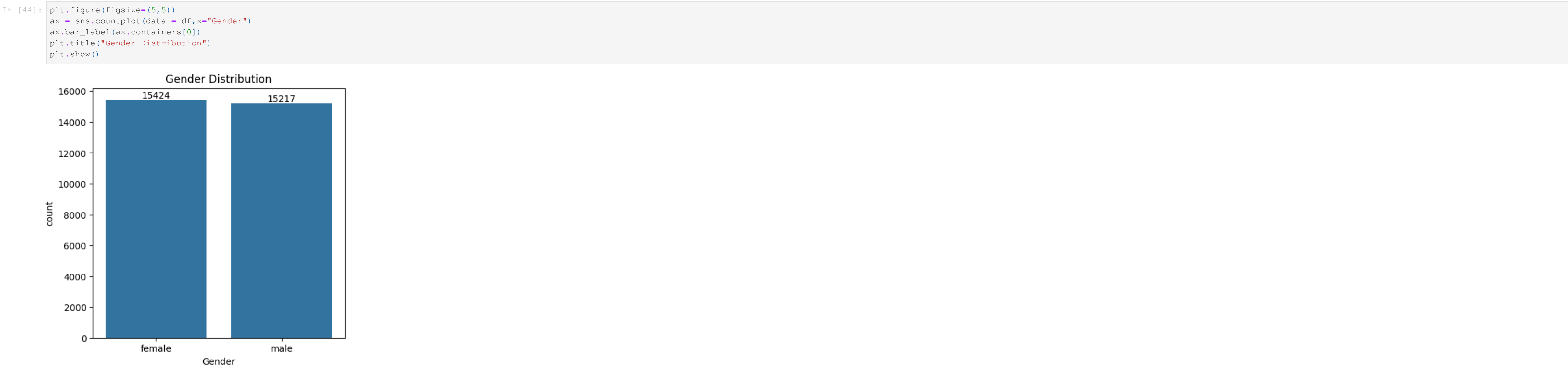
Change weekly study hours column

```
In [19]: df["WklyStudyHours"] = df["WklyStudyHours"].str.replace("5-Oct","5-Oct")
df.head()

Out[19]:
```

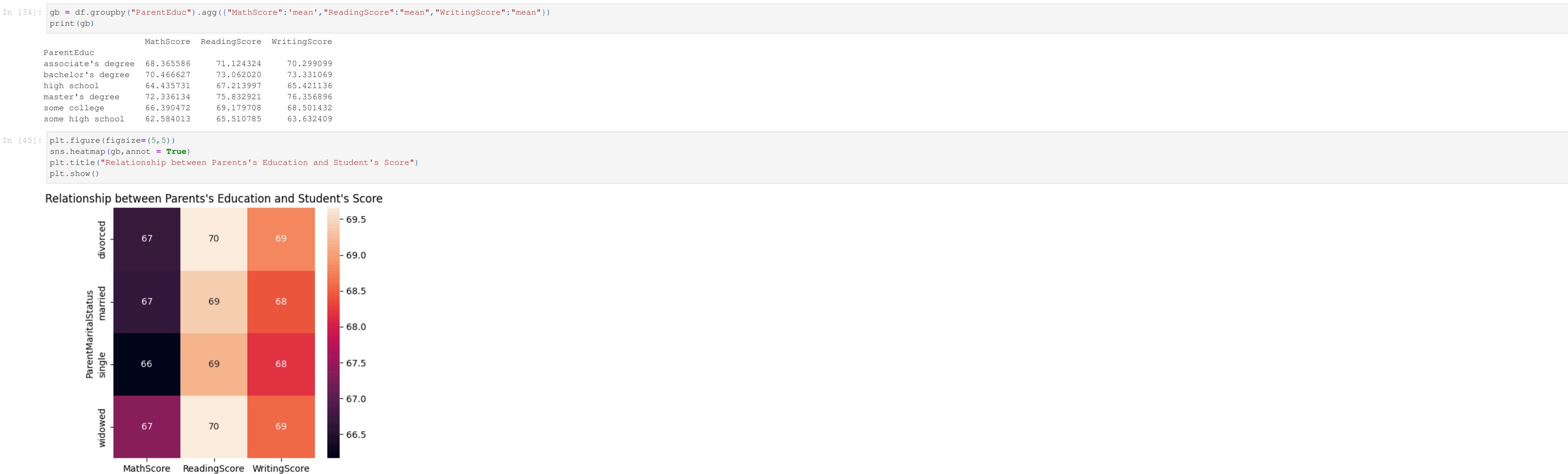
Gender	EthnicGroup	ParentEduc	LunchType	TestPrep	ParentMaritalStatus	PracticeSport	IsFirstChild	Nrsiblings	TransportMeans	WklyStudyHours	MathScore	ReadingScore	WritingScore	
0	female	NaN	bachelor's degree	standard	none	married	regularly	yes	3.0	school_bus	< 5	71	71	74
1	female	group C	some college	standard	NaN	married	sometimes	yes	0.0	NaN	5 - 10	69	90	88
2	female	group B	master's degree	standard	none	single	sometimes	yes	4.0	school_bus	< 5	87	93	91
3	male	group A	associate's degree	free/reduced	none	married	never	no	1.0	NaN	5 - 10	45	56	42
4	male	group C	some college	standard	none	married	sometimes	yes	0.0	school_bus	5 - 10	76	78	75

Gender distribution

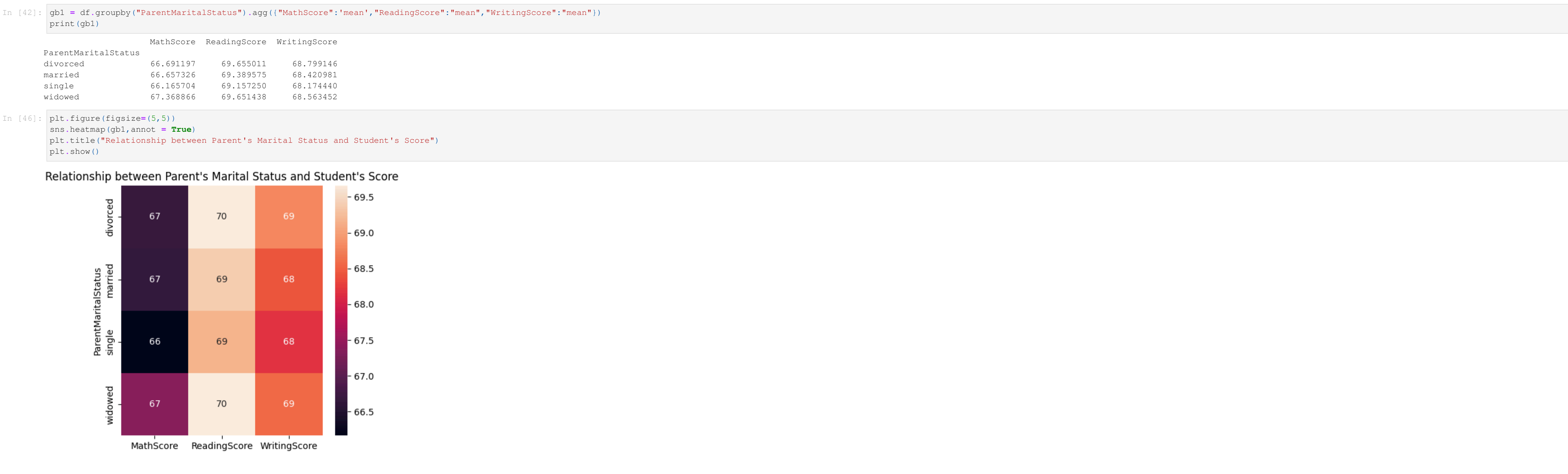


From the above chart we have analysed that:

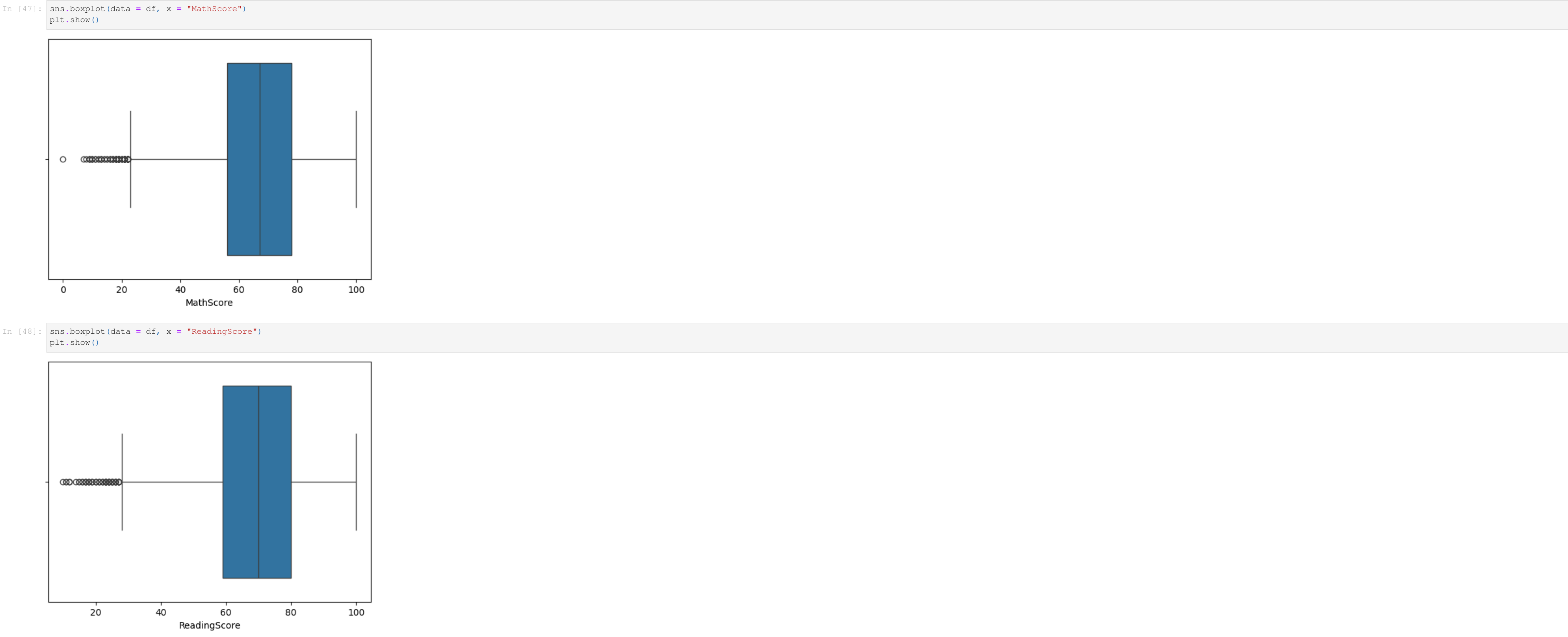
The number of females in the data is more than the number of males



From the above chart we have concluded that the education of the parents have a good impact on their scores



From the above we have concluded that there is no/negligible impact on the student's score due to their marital status



```
In [68]: ax = sns.countplot(data = df, x = "EthnicGroup")
ax.bar_label(ax.containers[0])

Out[68]:
```

Text(0, 0, '2121')
Text(0, 0, '5322')
Text(0, 0, '2219')
Text(0, 0, '7503')
Text(0, 0, '4961')

