

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

In [2]: df = pd.read_csv(r"C:\Users\Admin\Downloads\Expanded_data_with_more_features.csv.zip")
print(df.head)

<bound method NDFrame.head of      Unnamed: 0  Gender EthnicGroup      ParentEduc  LunchType  \
0              0  female      NaN  bachelor's degree      standard
1              1  female  group C      some college      standard
2              2  female  group B      master's degree      standard
3              3  male    group A  associate's degree  free/reduced
4              4  male    group C      some college      standard
...          ...          ...          ...          ...
38636         816  female  group D      high school      standard
38637         898  male    group E      high school      standard
38638         911  female      NaN      high school  free/reduced
38639         934  female  group D  associate's degree      standard
38640         969  male    group B      some college      standard

      TestPrep  ParentMaritalStatus  PracticeSport  IsFirstChild  NrSiblings  \
0             none      married      regularly      yes          3.0
1             NaN      married      sometimes      yes          0.0
2             none      single      sometimes      yes          4.0
3             none      married      never         no          1.0
4             none      married      sometimes     yes          0.0
...          ...          ...          ...          ...
38636         none      single      sometimes     no          2.0
38637         none      single      regularly     no          1.0
38638  completed      married      sometimes     no          1.0
38639  completed      married      regularly     no          3.0
38640         none      married      never         no          1.0

      TransportMeans  WklyStudyHours  MathScore  ReadingScore  WritingScore
0      school_bus      < 5              71              71              74
1             NaN      5 - 10             69              90              88
2      school_bus      < 5              87              93              91
3             NaN      5 - 10             45              56              42
4      school_bus      5 - 10             76              78              75
...          ...          ...          ...          ...
38636  school_bus      5 - 10             59              61              65
38637      private      5 - 10             58              53              51
38638      private      5 - 10             61              70              67
38639  school_bus      5 - 10             82              90              93
38640  school_bus      5 - 10             64              60              58

[38641 rows x 15 columns]>
```

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

Out[5]:
```

	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 [6]: 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        30610 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 [8]: df.isnull().sum()

Out[8]: Unnamed: 0      0
Gender      0
EthnicGroup 1840
ParentEduc  1845
LunchType   0
TestPrep    0
ParentMaritalStatus 1190
PracticeSport    631
IsFirstChild    904
NrSiblings      1572
TransportMeans  3134
WklyStudyHours  955
MathScore       0
ReadingScore    0
WritingScore    0
dtype: int64
```

```
In [12]: print(df.columns)

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

drop unnamed column#

```
In [6]: df = df.drop("Unnamed: 0",axis = 1)
print(df.head())
```

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	71	74
female	group C	some college	standard	NaN	married	sometimes	yes	0.0	NaN	5 - 10	69	90	88
female	group B	master's degree	standard	none	single	sometimes	yes	4.0	school_bus	< 5	87	93	91
male	group A	associate's degree	free/reduced	none	married	never	no	1.0	NaN	5 - 10	45	56	42
male	group C	some college	standard	none	married	sometimes	yes	0.0	school_bus	5 - 10	76	78	75

gender distribution#

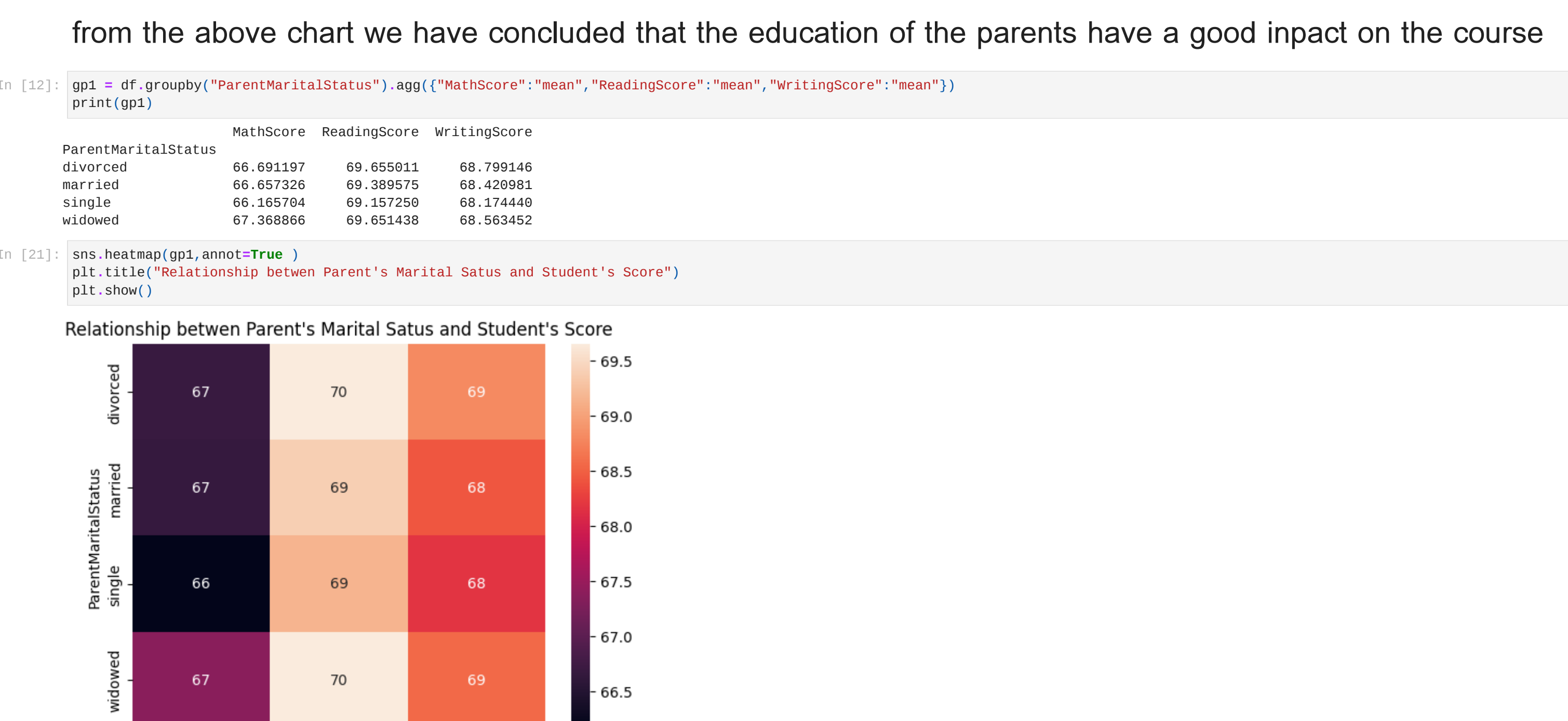


from the above chart we have analysed that:

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

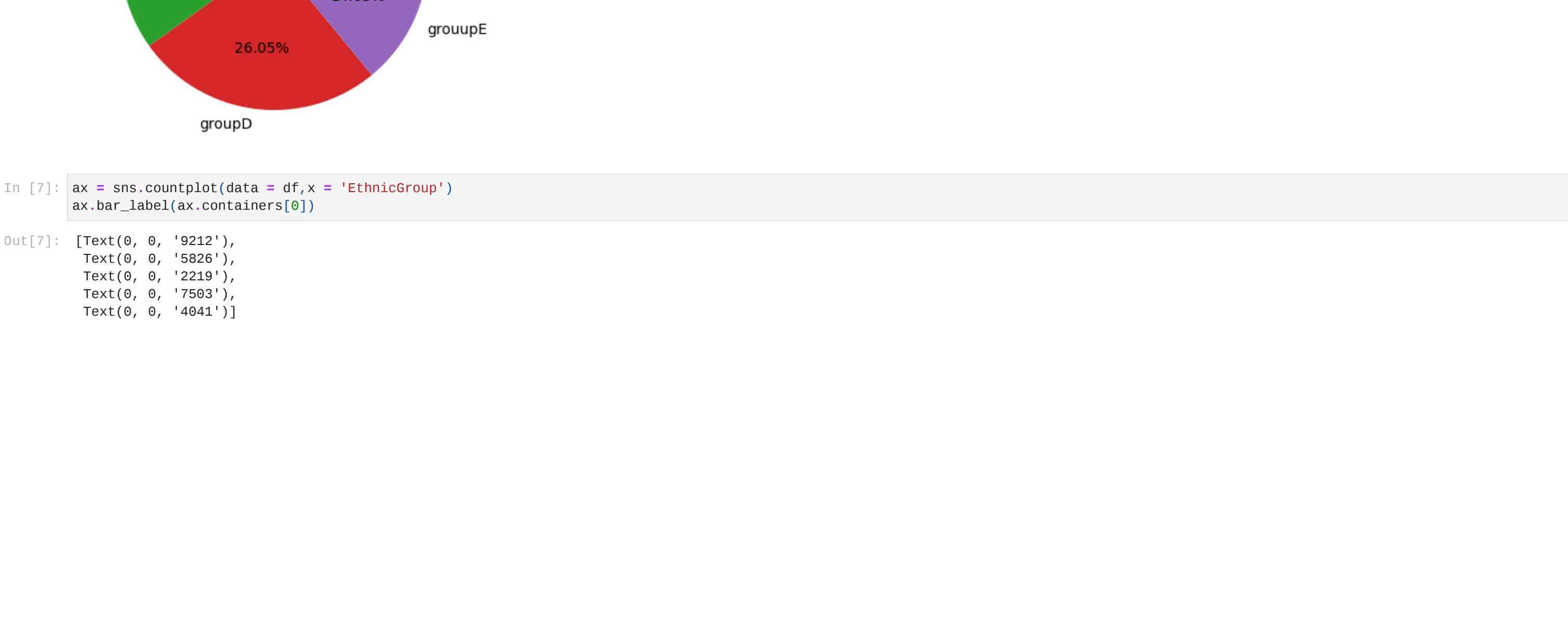
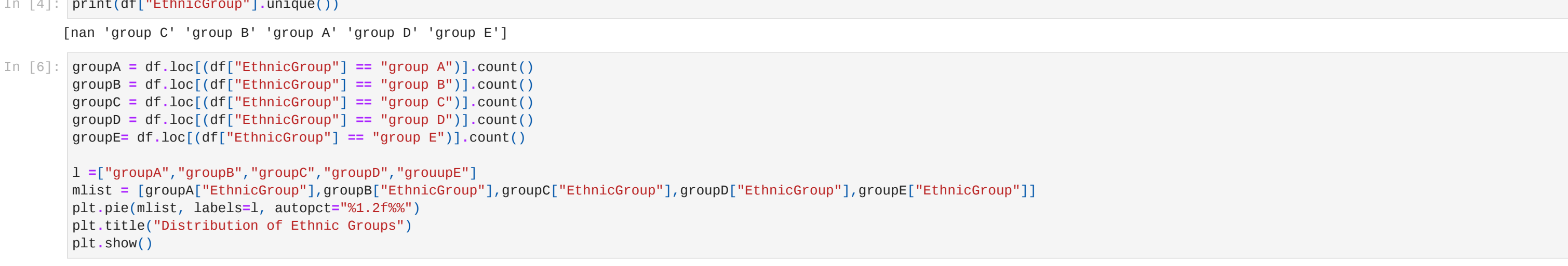
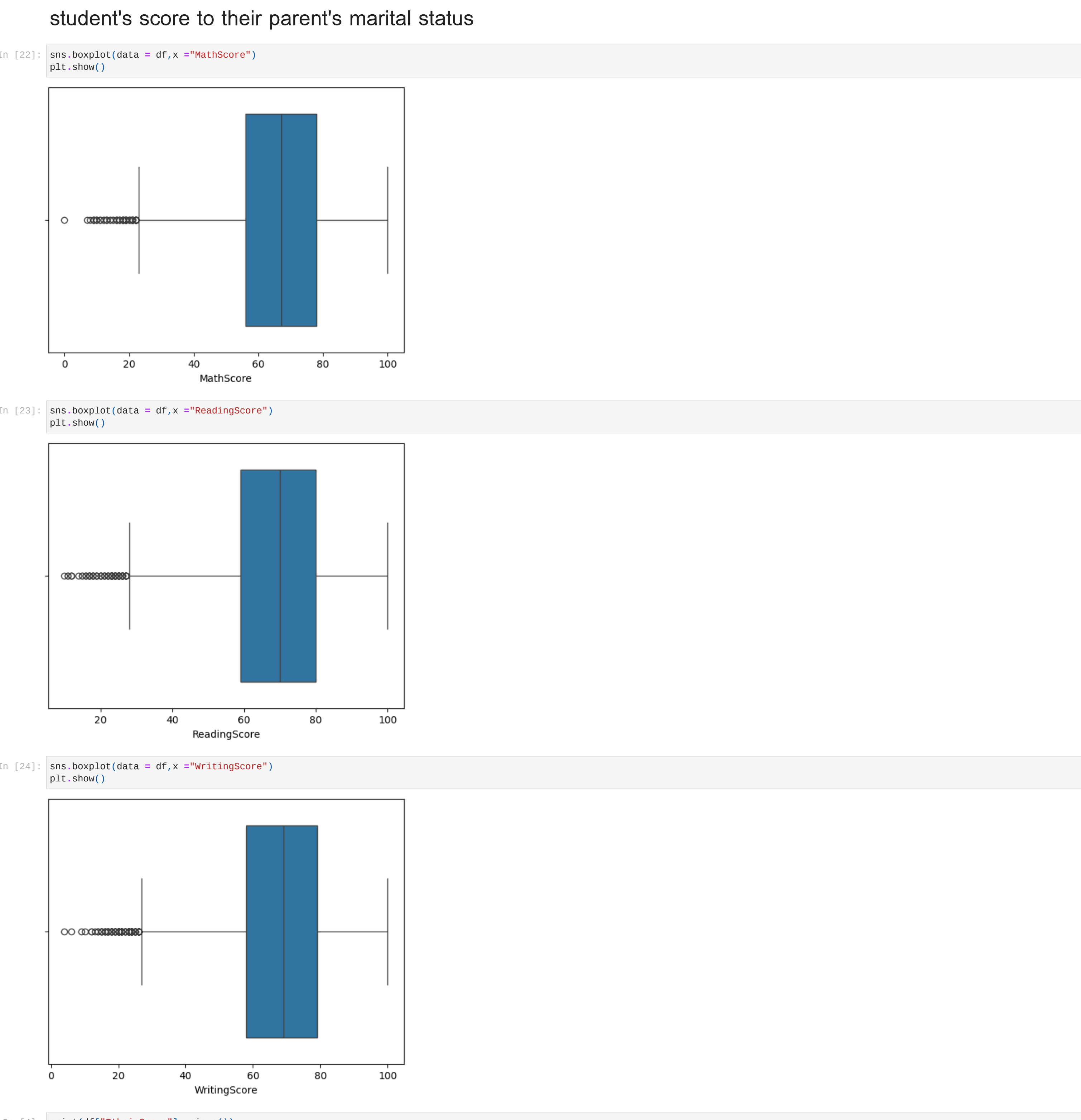


from the above chart we have concluded that the education of the parents have a good impact on the course



from the above chart we have concluded that there is no/negligible impact on the

student's score to their parent's marital status



```
In [7]: ax = sns.countplot(data = df, x = "EthnicGroup")
ax.bar_label(ax.containers[0])

Out[7]: [Text(0, 0, '9212'),
Text(0, 0, '5826'),
Text(0, 0, '2219'),
Text(0, 0, '7583'),
Text(0, 0, '4941')]
```

