

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

```
In [3]: df.head_csv(r"C:\Users\Owner\Downloads\titanic.csv")

Out[3]:
  survived  pclass  sex  age  sibsp  parch  fare  embarked  class  who  adult_male  deck  embark_town  alive  alone  Unnamed: 15
0         0         3  male  22.0  1         0  7.2500      S  Third  man      True  NaN  Southampton  no  False  NaN
1         1         1  female  38.0  1         0  71.2833      C  First  woman  False  C  Cherbourg  yes  False  NaN
2         1         3  female  26.0  0         0  7.9250      S  Third  woman  False  NaN  Southampton  yes  True  NaN
3         1         1  female  35.0  1         0  53.1000      S  First  woman  False  C  Southampton  yes  False  NaN
4         0         3  male  35.0  0         0  8.0500      S  Third  man      True  NaN  Southampton  no  True  NaN
```

```
In [4]: df.shape

Out[4]:
(891, 16)
```

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

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 16 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   survived              891 non-null    int64
 1   pclass               891 non-null    int64
 2   sex                  891 non-null    object
 3   age                  714 non-null    float64
 4   sibsp                891 non-null    int64
 5   parch               891 non-null    int64
 6   fare                 891 non-null    float64
 7   embarked             891 non-null    object
 8   class                891 non-null    object
 9   who                  891 non-null    object
10  adult_male           891 non-null    bool
11  deck                 283 non-null    object
12  embark_town          889 non-null    object
13  alive                891 non-null    object
14  alone                891 non-null    bool
15  Unnamed: 15          9 non-null      float64
dtypes: bool(2), float64(3), int64(4), object(7)
memory usage: 99.3+ KB#84x4, object(7)
```

```
In [6]: df.describe()

Out[6]:
   survived  pclass  age  sibsp  parch  fare  Unnamed: 15
count  891.000000  891.000000  714.000000  891.000000  891.000000  891.000000  0.0
mean  0.343582    2.806927  29.691196  0.529698  0.559094  52.242698  NaN
std    0.469082    0.900771  14.526497  1.107243  0.906007  49.694929  NaN
min    0.000000    1.000000  0.420000  0.000000  0.000000  0.000000  NaN
25%    0.000000    2.000000  20.120000  0.000000  0.000000  7.915000  NaN
50%    0.000000    3.000000  28.000000  0.000000  0.000000  14.451000  NaN
75%    1.000000    3.000000  38.000000  1.000000  0.000000  31.450000  NaN
max    1.000000    3.000000  80.000000  8.000000  6.000000  512.329200  NaN
```

```
In [7]: df.isnull().sum()

Out[7]:
survived      0
pclass        0
age          177
sex           0
parch        0
fare          0
embarked      0
class        0
who           0
adult_male    0
deck         688
embark_town   2
alive         0
alone         0
Unnamed: 15   891
dtype: int64
```

```
In [8]: df.duplicated().sum()

Out[8]:
107
```

```
In [9]: titanic = df.info().sum().sort_values(ascending=False)
titanic = titanic[titanic>107]
titanic

Out[9]:
Unnamed: 15    891
deck          688
dtype: int64
```

```
In [11]: # checking what % column contain missing values
df.isnull().sum().sort_values(ascending=False)*100/len(df)

Out[11]:
Unnamed: 15    100.000000
age            24.746989
sex            19.886520
embarked       7.746453
embark_town    0.224647
pclass         0.000000
sibsp          0.000000
parch          0.000000
fare           0.000000
class          0.000000
who            0.000000
adult_male     0.000000
alive          0.000000
alone          0.000000
dtype: float64
```

```
In [14]: df.columns.values

Out[14]:
array(['survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare',
       'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town',
       'alive', 'alone', 'Unnamed: 15'], dtype=object)
```

```
In [17]: sns.boxplot(x="survived",y="fare",data=df)

Out[17]:
<Axes: xlabel='survived', ylabel='fare'>
```



Passengers who paid over 300

```
In [18]: df[df['fare']>300]

Out[18]:
   survived  pclass  sex  age  sibsp  parch  fare  embarked  class  who  adult_male  deck  embark_town  alive  alone  Unnamed: 15
256         1         1  female  35.0  0         0  512.3292      C  First  woman  False  NaN  Cherbourg  yes  True  NaN
679         1         1  male  36.0  0         1  512.3292      C  First  man      True  B  Cherbourg  yes  False  NaN
737         1         1  male  35.0  0         0  512.3292      C  First  man      True  B  Cherbourg  yes  True  NaN
```

Embarked, Pclass and Sex

```
In [21]: FaceGrid = sns.FacetGrid(df, col='embarked', height=4, aspect=1.2)
FaceGrid.map(sns.countplot, 'pclass', 'survived', 'sex', c1=0.8, palette='deep', order=None, hue_order=None)
FaceGrid.add_legend()
```

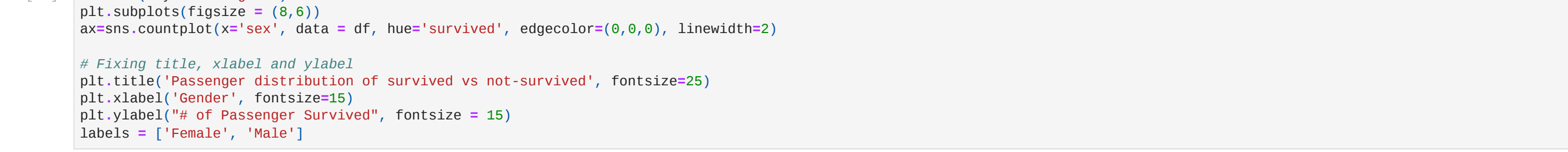
The 'ci' parameter is deprecated. Use 'errorbar=('ci', 95.0)' for the same effect.

```
func = plot_args, **plot_kwargs)
C:\Users\Owner\Desktop\A\Lib\site-packages\seaborn\axisgrid.py:848: FutureWarning:
The 'ci' parameter is deprecated. Use 'errorbar=('ci', 95.0)' for the same effect.
```

```
func = plot_args, **plot_kwargs)
C:\Users\Owner\Desktop\A\Lib\site-packages\seaborn\axisgrid.py:848: FutureWarning:
The 'ci' parameter is deprecated. Use 'errorbar=('ci', 95.0)' for the same effect.
```

```
func = plot_args, **plot_kwargs)
C:\Users\Owner\Desktop\A\Lib\site-packages\seaborn\axisgrid.py:848: FutureWarning:
The 'ci' parameter is deprecated. Use 'errorbar=('ci', 95.0)' for the same effect.
```

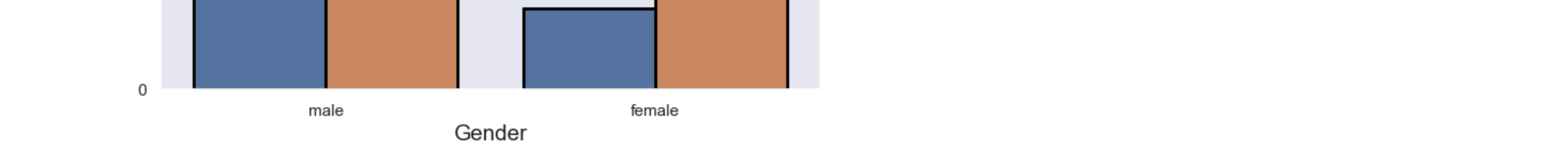
```
func = plot_args, **plot_kwargs)
C:\Users\Owner\Desktop\A\Lib\site-packages\seaborn\axisgrid.py:848: FutureWarning:
The 'ci' parameter is deprecated. Use 'errorbar=('ci', 95.0)' for the same effect.
```



Distribution of Pclass and Survived

```
In [23]: sns.set(style='darkgrid')
plt.title('Passenger distribution of survived vs not-survived', fontsize=25)
plt.xlabel('embark', fontsize=15)
plt.ylabel('% of Passenger Survived', fontsize=15)
labels = ['Female', 'Male']
```

Passenger distribution of survived vs not-survived



```
In [28]: # Fixing ticks.
plt.xticks(sorted(df.groupby('sex').mean()))

Out[28]:
[('female', 0.343582), ('male', 0.616418)]
```



```
In [29]: df.groupby(['sex'], mean())

Out[29]:
   survived  pclass  age  sibsp  parch  fare  adult_male  alone  Unnamed: 15
sex
female  0.742538  2.105226  27.915778  0.041288  0.048962  41.079818  0.000000  0.402724  NaN
male    0.188630  2.105948  30.726649  0.422893  0.236102  25.523893  0.930676  0.712305  NaN
```

Looking deeper into differences between females and males statistics

```
In [28]: df.groupby(['sex', 'pclass']).mean()

Out[28]:
   survived  pclass  age  sibsp  parch  fare  adult_male  alone  Unnamed: 15
sex pclass
female 1 0.960805 14.311705 0.553191 0.457447 106.125798 0.000000 0.361702 NaN
      2 0.921053 28.720713 0.486842 0.650263 19.707021 0.000000 0.421053 NaN
      3 0.500000 21.750000 0.895633 0.798611 16.118810 0.000000 0.416667 NaN
male 1 0.368824 41.261386 0.314175 0.279889 67.226117 0.975420 0.614754 NaN
      2 0.457400 20.140707 0.436058 0.222222 10.747702 0.056607 0.666667 NaN
      3 0.125447 26.507589 0.488559 0.220794 12.661633 0.919328 0.706907 NaN
```

Age and Sex distributions

```
In [28]: survived = 'survived'
not_survived = 'not survived'

fig, axes = plt.subplots(nrows=1, ncols=2, figsize=(15, 5))

women = df[df['sex']=='female']
men = df[df['sex']=='male']

# Plot Female Survived vs Not-Survived distribution
ax = sns.histplot(women[women['survived']==1] age dropna(), bins=20, label = survived, ax = axes[0], color='b', kde=True)
ax = sns.histplot(women[women['survived']==0] age dropna(), bins=20, label = not_survived, ax = axes[0], color='r', kde=True)
ax.legend()
ax.set_title('female')

# Plot Male Survived vs Not-Survived distribution
ax = sns.histplot(men[men['survived']==1] age dropna(), bins=20, label = survived, ax = axes[1], color='b', kde=True)
ax = sns.histplot(men[men['survived']==0] age dropna(), bins=20, label = not_survived, ax = axes[1], color='r', kde=True)
ax.legend()
ax.set_title('male')
```

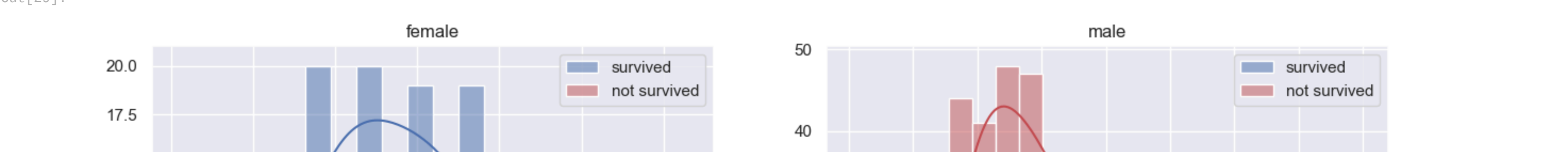


```
In [29]: df['age'].groupby(['sex', 'pclass']).mean()

Out[29]:
   survived  age  sibsp  parch  fare  adult_male  alone  Unnamed: 15
sex pclass
female 1 0.875000 14.125000 0.500000 0.875000 104.083337 0.000000 0.125000 NaN
      2 1.000000 8.333333 0.583333 1.083333 25.241667 0.000000 0.166667 NaN
      3 0.542857 8.428571 1.571429 1.057143 18.727977 0.000000 0.228571 NaN
male 1 1.180000 8.230000 0.500000 2.000000 116.073900 0.250000 0.000000 NaN
      2 0.818182 4.757273 0.727273 1.000000 25.694773 0.181818 0.191318 NaN
      3 0.222228 9.963258 2.069767 1.000000 22.752523 0.948887 0.222228 NaN
```

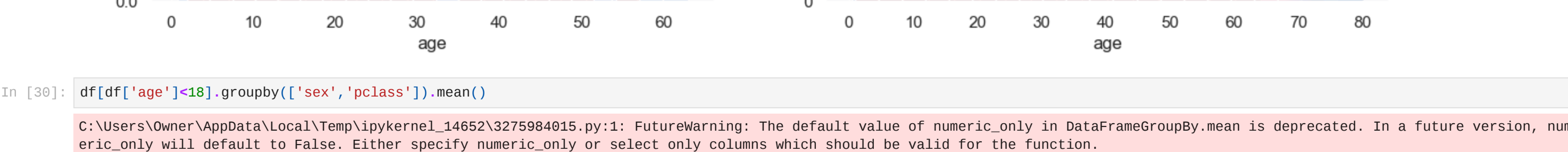
```
In [31]: plt.subplots(figsize=(8,8))
sns.countplot(x='pclass', hue='survived', data=df)
plt.title('Passenger Class Distribution - Survived vs Non-Survived', fontsize=25)
plt.xlabel('age', size=18)
plt.ylabel('Passenger Class Distribution - Survived vs Non-Survived', fontsize=25)
Text(0.5, 1.8, 'Passenger Class Distribution - Survived vs Non-Survived')
```

Passenger Class Distribution - Survived vs Non-Survived



```
In [34]: plt.subplots(figsize=(8,8))
sns.boxplot(x='pclass', y='survived', data=df)
plt.title('Passenger Class Distribution - Survived Passengers', fontsize=25)
Text(0.5, 1.8, 'Passenger Class Distribution - Survived Passengers')
```

Passenger Class Distribution - Survived Passengers

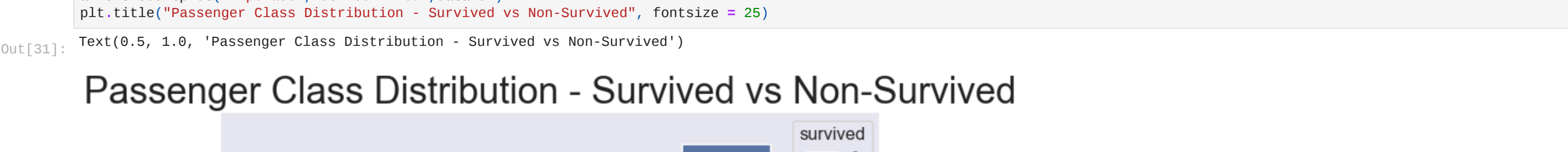


Correlation Matrix and Heatmap

```
In [37]: # Look at numeric and categorical values separately
df_num = df[['age', 'sibsp', 'parch', 'fare']]
df_cat = df[['survived', 'pclass', 'sex', 'embarked']]

In [38]: sns.heatmap(df_num.corr(), annot=True, cmap='magma')
plt.title('Correlations Among Numeric Features', fontsize=18)

Out[38]:
Text(0.5, 1.8, 'Correlations Among Numeric Features')
```

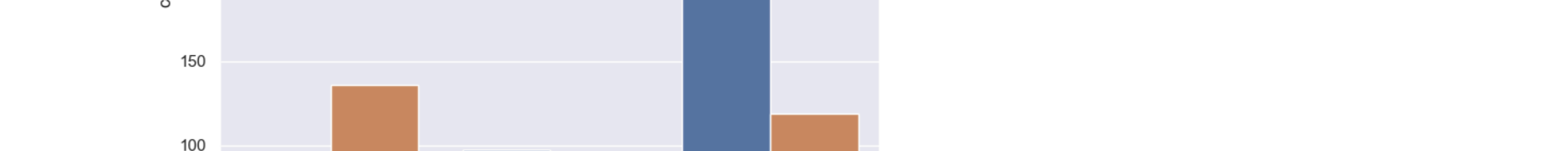


Combining SibSp and Parch

```
In [39]: data = df[
    'survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare', 'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town', 'alive', 'alone', 'Unnamed: 15', 'relatives', 'not_alone'
]

Out[39]:
0 1 537
1 384
Name: not_alone, dtype: int64

In [40]: plt.subplots(figsize=(16,4))
ax = sns.lineplot(x='relatives', y='survived', data=df)
```

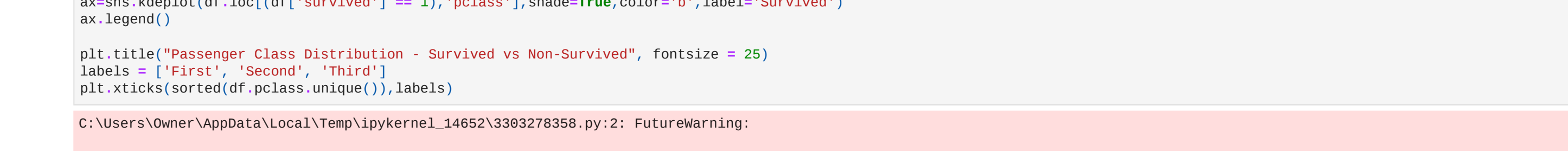


Missing Data

```
In [40]: data = df[
    'survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare', 'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town', 'alive', 'alone', 'Unnamed: 15', 'relatives', 'not_alone'
]

Out[40]:
0 1 537
1 384
Name: not_alone, dtype: int64

In [41]: plt.subplots(figsize=(16,4))
ax = sns.lineplot(x='relatives', y='survived', data=df)
```

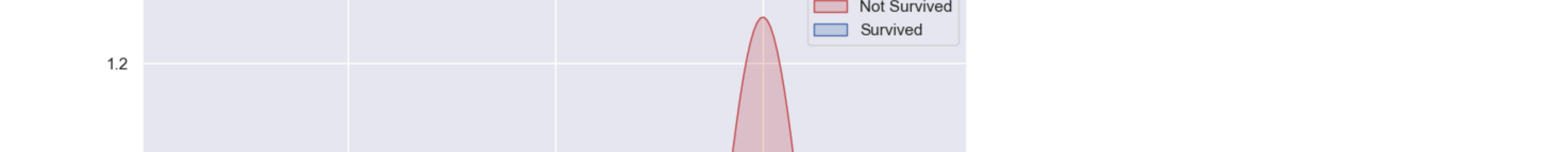


Missing Data

```
In [40]: data = df[
    'survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare', 'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town', 'alive', 'alone', 'Unnamed: 15', 'relatives', 'not_alone'
]

Out[40]:
0 1 537
1 384
Name: not_alone, dtype: int64

In [41]: plt.subplots(figsize=(16,4))
ax = sns.lineplot(x='relatives', y='survived', data=df)
```



Missing Data

```
In [40]: data = df[
    'survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare', 'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town', 'alive', 'alone', 'Unnamed: 15', 'relatives', 'not_alone'
]

Out[40]:
0 1 537
1 384
Name: not_alone, dtype: int64

In [41]: plt.subplots(figsize=(16,4))
ax = sns.lineplot(x='relatives', y='survived', data=df)
```



Missing Data

```
In [40]: data = df[
    'survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare', 'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town', 'alive', 'alone', 'Unnamed: 15', 'relatives', 'not_alone'
]

Out[40]:
0 1 537
1 384
Name: not_alone, dtype: int64

In [41]: plt.subplots(figsize=(16,4))
ax = sns.lineplot(x='relatives', y='survived', data=df)
```

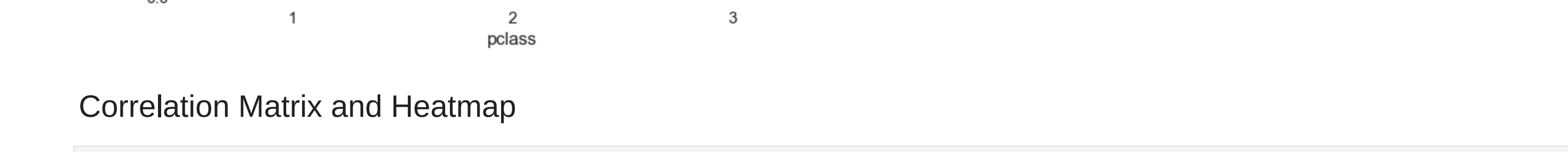


Missing Data

```
In [40]: data = df[
    'survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare', 'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town', 'alive', 'alone', 'Unnamed: 15', 'relatives', 'not_alone'
]

Out[40]:
0 1 537
1 384
Name: not_alone, dtype: int64

In [41]: plt.subplots(figsize=(16,4))
ax = sns.lineplot(x='relatives', y='survived', data=df)
```

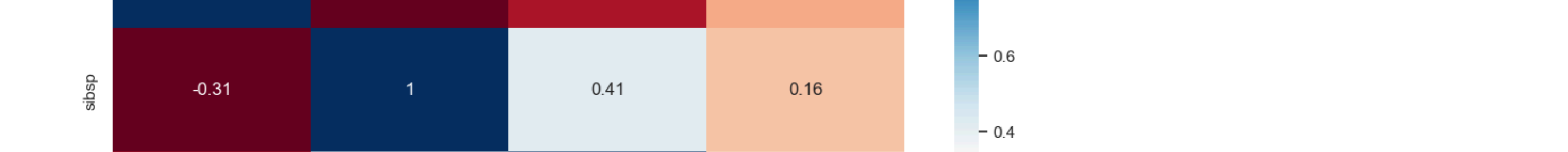


Missing Data

```
In [40]: data = df[
    'survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare', 'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town', 'alive', 'alone', 'Unnamed: 15', 'relatives', 'not_alone'
]

Out[40]:
0 1 537
1 384
Name: not_alone, dtype: int64

In [41]: plt.subplots(figsize=(16,4))
ax = sns.lineplot(x='relatives', y='survived', data=df)
```

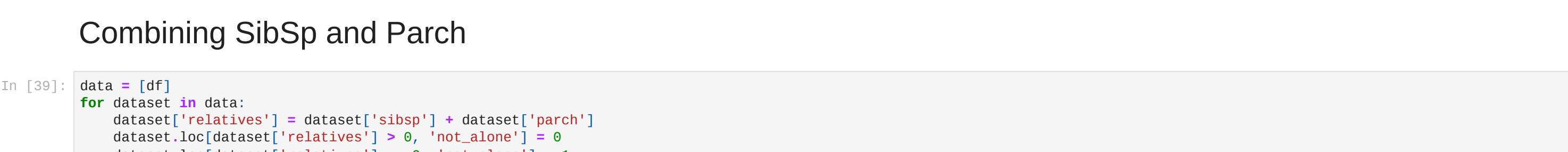


Missing Data

```
In [40]: data = df[
    'survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare', 'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town', 'alive', 'alone', 'Unnamed: 15', 'relatives', 'not_alone'
]

Out[40]:
0 1 537
1 384
Name: not_alone, dtype: int64

In [41]: plt.subplots(figsize=(16,4))
ax = sns.lineplot(x='relatives', y='survived', data=df)
```



Missing Data

```
In [40]: data = df[
    'survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare', 'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town', 'alive', 'alone', 'Unnamed: 15', 'relatives', 'not_alone'
]

Out[40]:
0 1 537
1 384
Name: not_alone, dtype: int64

In [41]: plt.subplots(figsize=(16,4))
ax = sns.lineplot(x='relatives', y='survived', data=df)
```

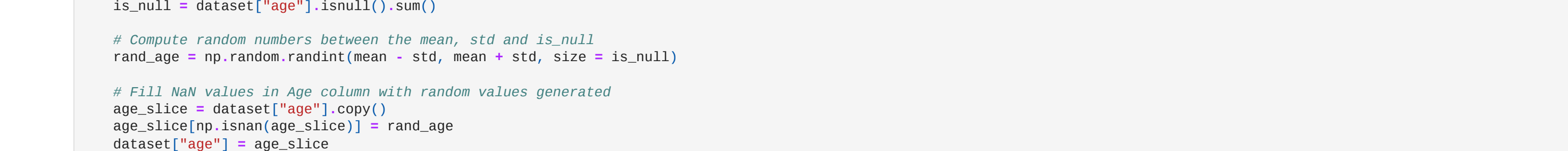


Missing Data

```
In [40]: data = df[
    'survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare', 'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town', 'alive', 'alone', 'Unnamed: 15', 'relatives', 'not_alone'
]

Out[40]:
0 1 537
1 384
Name: not_alone, dtype: int64

In [41]: plt.subplots(figsize=(16,4))
ax = sns.lineplot(x='relatives', y='survived', data=df)
```



Missing Data

```
In [40]: data = df[
    'survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare', 'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town', 'alive', 'alone', 'Unnamed: 15', 'relatives', 'not_alone'
]

Out[40]:
0 1 537
1 384
Name: not_alone, dtype: int64

In [41]: plt.subplots(figsize=(16,4))
ax = sns.lineplot(x='relatives', y='survived', data=df)
```

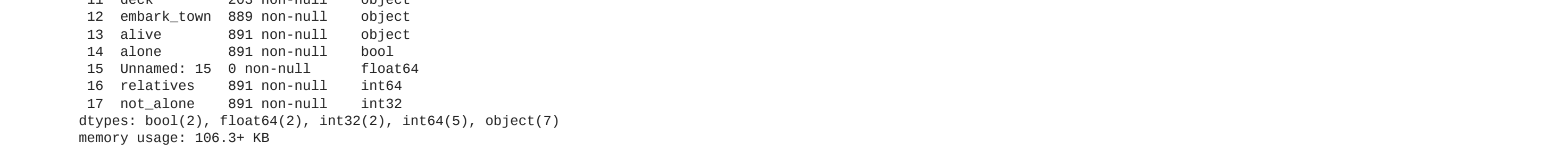


Missing Data

```
In [40]: data = df[
    'survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare', 'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town', 'alive', 'alone', 'Unnamed: 15', 'relatives', 'not_alone'
]

Out[40]:
0 1 537
1 384
Name: not_alone, dtype: int64

In [41]: plt.subplots(figsize=(16,4))
ax = sns.lineplot(x='relatives', y='survived', data=df)
```



Missing Data

```
In [40]: data = df[
    'survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare', 'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town', 'alive', 'alone', 'Unnamed: 15', 'relatives', 'not_alone'
]

Out[40]:
0 1 537
1 384
Name: not_alone, dtype: int64

In [41]: plt.subplots(figsize=(16,4))
ax = sns.lineplot(x='relatives', y='survived', data=df)
```



Missing Data

```
In [40]: data = df[
    'survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare', 'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town', 'alive', 'alone', 'Unnamed: 15', 'relatives', 'not_alone'
]

Out[40]:
0 1 537
1 384
Name: not_alone, dtype: int64

In [41]: plt.subplots(figsize=(16,4))
ax = sns.lineplot(x='relatives', y='survived', data=df)
```



Missing Data

```
In [40]: data = df[
    'survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare', 'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town', 'alive', 'alone', 'Unnamed: 15', 'relatives', 'not_alone'
]

Out[40]:
0 1 537
1 384
Name: not_alone, dtype: int64

In [41]: plt.subplots(figsize=(16,4))
ax = sns.lineplot(x='relatives', y='survived', data=df)
```

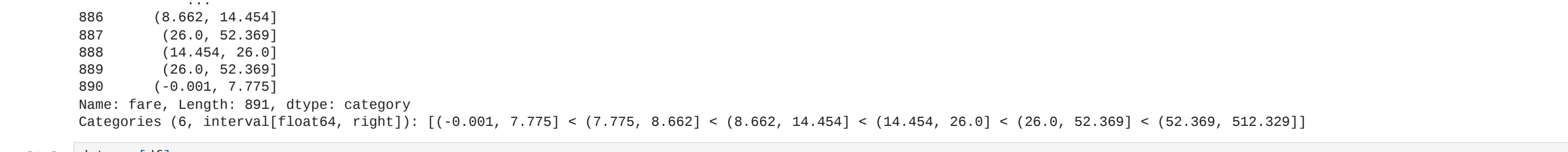


Missing Data

```
In [40]: data = df[
    'survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare', 'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town', 'alive', 'alone', 'Unnamed: 15', 'relatives', 'not_alone'
]

Out[40]:
0 1 537
1 384
Name: not_alone, dtype: int64

In [41]: plt.subplots(figsize=(16,4))
ax = sns.lineplot(x='relatives', y='survived', data=df)
```



Missing Data

```
In [40]: data = df[
    'survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare', 'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town', 'alive', 'alone', 'Unnamed: 15', 'relatives', 'not_alone'
]

Out[40]:
0 1 537
1 384
Name: not_alone, dtype: int64

In [41]: plt.subplots(figsize=(16,4))
ax = sns.lineplot(x='relatives', y='survived', data=df)
```

